# Institut National Polytechnique de Toulouse

## HABILITATION À DIRIGER DES RECHERCHES

## Aspects of Covariance Modelling in Variational Ocean Data Assimilation

soutenue le 13 décembre 2013

 $\operatorname{par}$ 

### Anthony Weaver

### CERFACS / SUC URA 1875

soutenue devant le Jury composé de :

M. Fisher	ECMWF, Reading	Rapporteur
A. M. Moore	University of California, Santa Cruz	Rapporteur
O. Talagrand	École Normale Supérieure, Paris	Rapporteur
D. L. T. Anderson	Oxford University, Oxford	Membre Invité
M. A. Balmaseda	ECMWF, Reading	Membre Invité
G. Desroziers	Météo-France, Toulouse	Membre Invité
O. Thual	INPT-IMFT, Toulouse	Membre Invité
S. Gratton	INPT-IRIT, Toulouse	Correspondant

Année universitaire 2013–2014

### Abstract

This manuscript summarizes research I have conducted in ocean data assimilation, emphasizing my work in background-error covariance modelling. First, I describe my early work in applying four-dimensional variational data assimilation to ocean models. Results from both idealized and realistic studies are presented, that highlight the fundamental importance of the background-error covariance matrix (**B**). Then, the manuscript describes specific developments I have made to improve the specification of **B** in variational ocean data assimilation, and discusses areas where further improvements can be expected. 

# Contents

1	Introduction 7			
2	Mat 2.1 2.2 2.3	thematical formulation of the 4D-Var problem         Nonlinear formulation         Incremental formulation and algorithm         Role of the $J_b$ term and $\mathbf{B}$ matrix		
3	Dev	elopm	ents in 4D-Var for the ocean	17
3.1 4D-Var in a simplified framework			r in a simplified framework	17
3.2 4D-Var in a realistic framework		r in a realistic framework	19	
		3.2.1	Validation of the linear approximation	21
		3.2.2	Flow-dependent background-error variances in 4D-Var	23
		3.2.3	Improving the background-error covariance model	24
4 Developments in background-error covariance modelling for the			ents in background-error covariance modelling for the ocean	29
4.1 Univariate correlation modelling using a diffusion equation			iate correlation modelling using a diffusion equation	30
		4.1.1	Two-dimensional correlation models	30
		4.1.2	Solid-wall boundary conditions	35
		4.1.3	Inhomogeneity and non-separability	36
		4.1.4	Three-dimensional correlation models $\ldots \ldots \ldots \ldots \ldots \ldots$	39
		4.1.5	Anisotropy	40
4.2 Multivariate covariance modelling using balance operators		ariate covariance modelling using balance operators	41	
		4.2.1	A general transformation of variables	42
		4.2.2	Temperature-salinity constraints	44
		4.2.3	Additional balance constraints	45
	4.3 Covariance estimation using an ensemble			49
		4.3.1	Ensemble estimation of background-error variances	49
		4.3.2	Ensemble estimation of background-error correlations	51
<b>5</b>	Summary and outlook 5			
$\mathbf{A}$	Selected articles			

# Chapter 1 Introduction

Data assimilation has its roots in meteorology where the need to produce accurate initial conditions for Numerical Weather Prediction (NWP) models has long been recognized as critically important (Daley 1991; Kalnay 2003). Over the past two decades, fourdimensional variational assimilation (4D-Var) has emerged as the preferred data assimilation method of most major operational NWP centres (Rabier et al. 2000; Desroziers et al. 2003; Rawlins et al. 2007; Gauthier et al. 2007). Conceptually, the basic objective of 4D-Var (in its standard strong-constraint formulation) is extremely simple and appealing: determine the model state trajectory that best fits, in a weighted least squares sense, the observations available over a given time interval. Since the model trajectory is controlled by the initial conditions, the problem can be effectively reduced to one of optimizing the model initial conditions at the beginning of the assimilation interval (Lewis and Derber 1985; Le Dimet and Talagrand 1986; Talagrand and Courtier 1987; Thacker and Long 1988). Here the term observations is used in a general sense, referring to all available information on the interval. Typically, this information consists of a prior estimate of the initial conditions (the model background state), as well as actual observations from measurements. Mathematically, 4D-Var is a constrained minimization problem involving a cost function that measures the model fit to the observations and the background state, subject to the constraint that the solution must satisfy the prognostic equations of the numerical model. The weights that control the closeness of fit to the observations and background are specified in terms of estimates of the observation- and background-error covariances. In this way, the information can be assimilated according to its perceived accuracy.

In oceanography, data assimilation systems are less advanced than those in NWP but have been significantly improved in recent years in response to important developments to the global ocean observing system<sup>1</sup> and with the advent of operational activities in short-term ocean forecasting with high-resolution ocean models<sup>2</sup> and in seasonal climate forecasting with coupled ocean-atmosphere models<sup>3</sup>. Extensive research has been conducted in ocean 4D-Var and several centres such as ECMWF and the UK Met Office have

<sup>&</sup>lt;sup>1</sup>See the OceanObs'09 Conference Proceedings available at http://www.oceanobs09.net.

<sup>&</sup>lt;sup>2</sup>See the collection of articles from the Special Issue on "The Revolution in Global Ocean Forecasting - GODAE: 10 Years of Achievement", *Oceanography*, **22**, 2009.

 $<sup>^{3}</sup>$ For example, see http://www.ecmwf.int/products/forecasts/seasonal for a description of operational seasonal forecasting activities at ECMWF.

recently implemented variational-based systems for research and operational applications (Mogensen *et al.* 2009; Mogensen *et al.* 2012; Balmaseda *et al.* 2013; Waters *et al.* 2013). There has been particular interest in oceanography in applying 4D-Var to the field of retrospective analysis, or reanalysis, which concerns the task of reconstructing the past ocean state by combining highly quality-controlled historical observations with a state-of-the-art ocean model and atmospheric forcing field, with the aim of improving our understanding of the ocean circulation and its variability<sup>4</sup>. The atmospheric forcing field used for an ocean reanalysis is itself generally produced from an atmospheric reanalysis (Uppala *et al.* 2005; Dee *et al.* 2011). Reanalysis has become an important application of data assimilation in view of its fundamental role for understanding past and present climate variability, and for calibrating climate forecast models.

Given the complexity of state-of-the-art atmospheric and ocean general circulation models (GCMs), their high cost of integration and the size of their state vectors, which is greater than several million on any given time step, the practical feasibility of 4D-Var is by no means obvious. To determine an exact global minimizing solution of the cost function is impossible for all but very simple nonlinear problems. Techniques for minimizing the cost function typically involve successive linearization steps. Even for linearized problems, however, determining an exact minimizing solution using a direct solution method leads to matrix equations that are far too big to manipulate without major simplifications.

Certain algorithmic features of 4D-Var are key to its feasibility. Rather than trying to solve the problem exactly, only an approximate solution is sought using iterative methods. As such, only operators (vector-to-vector transformations) are required for the solution algorithm and there is thus no need to manipulate large matrices explicitly. Limitedmemory quasi-Newton methods (Nocedal 1980; Gilbert and Lemaréchal 1989), which belong to the class of so-called Krylov subspace methods (Golub and van Loan 1996; Nocedal and Wright 1999), are iterative algorithms that are specially designed for solving nonlinear minimization problems with a large number of adjustable parameters. To apply these methods to 4D-Var requires, on each iteration, the computation of the gradient of the cost function with respect to the initial conditions. As initially pointed out by Le Dimet and Talagrand (1986), the gradient vector can be computed efficiently by integrating (in reverse time) the adjoint of the tangent-linear of the underlying numerical model. Much of the early work in 4D-Var was devoted to the practical design of accurate adjoint models and their application within quasi-Newton minimization algorithms (Courtier and Talagrand 1987; Courtier and Talagrand 1990; Thépaut and Courtier 1991).

The introduction of the incremental formulation was an important milestone in the practical development of 4D-Var (Courtier *et al.* 1994). In the incremental approach, the full nonlinear 4D-Var problem defined by the minimization of a nonquadratic cost function is transformed into a sequence of minimizations of simpler, linearized 4D-Var problems involving quadratic cost functions. Linearization is appropriate if nonlinear processes are sufficiently 'weak', so care must be taken to ensure that the width of the assimilation window does not exceed the time-scale of validity of the linear approximation for increments generated during the minimization. Though less apparent, the linear hypothesis is also important in nonincremental 4D-Var since adjoint models used for computing the gradient are constructed through transposition of the tangent-linear of the nonlinear model.

<sup>&</sup>lt;sup>4</sup>Global ocean renalysis is promoted and coordinated at the international level by the CLIVAR Global Synthesis and Observations Panel (GSOP). See http://www.clivar.org/organization/gsop.

The quadratic minimization in incremental 4D-Var is usually referred to as the *inner* loop and involves computing corrections, or increments, to a 'guess' estimate of the initial conditions. Krylov-subspace methods based on conjugate gradient or Lanczos iterative algorithms are usually employed in the inner loop as they are particularly effective for minimizing quadratic cost functions (Golub and van Loan 1995; Meurant 2006). Nonlinear computations involving the 'guess' estimate are performed on the so-called *outer loop*. The complete incremental 4D-Var algorithm involves a feedback process to allow the basic state trajectory of the linearized model to be updated using the most recent guess estimate of the state trajectory from the nonlinear model. On the first outer iteration, the guess estimate is defined as the background state. The solution, or analysis, is taken to be the updated guess estimate on the final outer iteration. Incremental 4D-Var is in fact a variant of the classical optimization method known as Gauss-Newton (GN) in the optimization literature (Gratton et al. 2007). In particular, Gratton et al. (2007) refer to incremental 4D-Var as a *perturbed* and *truncated* GN algorithm since the linearized operators are generally (and sometimes crude) approximations to the exact tangent-linear operators, and since the inner-loop minimization is only solved approximately. In addition, only a very small number of outer iterations is affordable for very large problems.

Incremental 4D-Var should be viewed as a practical minimization algorithm for approximately solving the complete (nonincremental) 4D-Var problem. In particular, it provides a more flexible framework than nonincremental 4D-Var for simplifying technical development and reducing computational costs. This is a major advantage for applications with complex, high-resolution GCMs. A hierarchy of linearized models can be used in the inner loop thereby providing a clear development path from simplified to more complex incremental 4D-Var systems. In practice, approximations usually involve simplifying the linearization of highly nonlinear physical parameterizations and/or reducing the resolution of the increment used in the inner loop. An extreme case is when the linearized model is approximated by the identity matrix. The resulting algorithm is effectively a three-dimensional variational (3D-Var) algorithm, since the temporal dimension is removed from the inner loop. It is referred to as 3D-Var 'First-Guess at Appropriate Time' (FGAT) in meteorology and oceanography, since the full nonlinear model is still retained in the outer loop for computing the model misfit with the observations. 3D-Var FGAT is significantly cheaper than incremental 4D-Var. Despite its simplicity, 3D-Var FGAT can be effective and is a natural starting point in the development of incremental 4D-Var.

My research in variational ocean data assimilation started during my doctoral work at the University of Oxford where I studied 4D-Var in a simplified framework. As a postdoctoral fellow at the Laboratoire d'Océanographie Dynamique et de Climatologie  $(LODYC)^5$  in Paris, I began working on the development of an incremental 4D-Var system for the OPA GCM (Madec *et al.* 1999), this model now being part of the more general framework known as NEMO (Nucleus for European Modelling of the Ocean; Madec 2008). I pursued this work as a senior researcher at CERFACS and it eventually resulted in what became

<sup>&</sup>lt;sup>5</sup>Now known as the Laboratoire d'Oceanographie et du Climat: Expérimentation et Approches Numériques (LOCEAN).

known as the OPAVAR system. The main focus of my research in data assimilation has been on improving the scientific and technical development of OPAVAR and its recent successor NEMOVAR.

The remainder of the manuscript is organized as follows. A mathematical formulation of the 4D-Var problem and incremental algorithm is given in Chapter 2. This chapter will provide a helpful reference for the specific research topics described in Chapters 3 and 4. Chapter 3 summarizes some of my early research in 4D-Var which involved both idealized and real-data applications. In particular, I focus on aspects of that work that illustrate the importance of the background-error covariance matrix and led me to develop a specific research activity in covariance modelling. In Chapter 4, I summarize my main contributions in this area, much of which has resulted from collaborative work with PhD students. The first part of Chapter 4 describes theoretical and numerical issues in the design of correlation operators for variational assimilation. In this work, a diffusion algorithm is proposed as a practical and flexible approach to represent the action of a correlation operator in complex boundary domains such as those encountered in ocean models. While much of this work has been published, some aspects are original. In the second part of Chapter 4, I discuss *multivariate* aspects of background-error covariance modelling. Here the physical nature of covariance modelling is emphasized, through the need to produce model-variable corrections that are in appropriate dynamical balance. The problem of es*timating* background-error covariances using a combined ensemble-variational approach is discussed in the third part of Chapter 4. Here the emphasis is on exploiting error information from ensembles to calibrate *parameters* of a prescribed background-error covariance model, rather than to specify the covariances directly from a sample estimate. Finally, in Chapter 5, I describe ongoing work and some of the future challenges in covariance modelling. My curriculum vitae is given Appendix A. Selected articles of which I am co-author and on which I have based this manuscript are listed in Appendix B. In the text I have referenced these articles, and PhD theses I have supervised, using a boldface font to clearly distinguish them from other references.

## Chapter 2

# Mathematical formulation of the 4D-Var problem

### 2.1 Nonlinear formulation

Let  $\mathbf{x}(t_i) \in \mathbb{R}^n$  denote the model state vector at the discrete time  $t_i$ , which is obtained by propagating the state vector  $\mathbf{x}(t_{i-1})$  from  $t_{i-1}$  to  $t_i$  using the nonlinear model operator  $M(t_i, t_{i-1}) : \mathbb{R}^n \to \mathbb{R}^n$ ;

$$\mathbf{x}(t_i) = M(t_i, t_{i-1})[\mathbf{x}(t_{i-1})].$$
(2.1)

For an ocean model, the elements of  $\mathbf{x}(t_i)$  are typically potential temperature (T), salinity (S), sea-surface height  $(\eta)$  and the horizontal components of velocity  $(\mathbf{u}^h = (u, v))$  at discrete points on the three-dimensional (3D) model grid. Assume that the ocean is observed over a period  $t_0 \leq t_i \leq t_N$  and let  $\mathbf{y}^o \in \mathbb{R}^p$  denote the observation vector where

$$\mathbf{y}^{\mathrm{o}} = \begin{pmatrix} \mathbf{y}_{0}^{\mathrm{o}} \\ \vdots \\ \mathbf{y}_{i}^{\mathrm{o}} \\ \vdots \\ \mathbf{y}_{N}^{\mathrm{o}} \end{pmatrix}, \qquad (2.2)$$

 $\mathbf{y}_i^{o} \in \mathbb{R}^{p_i}$  being the observation vector at time  $t_i$ , and  $\sum_{i=0}^{N} p_i = p$ . Let  $\mathbf{x}^{b}(t_0)$  be the model background initial state. The background estimate of the state at future times in the period is obtained through successive applications of Eq. (2.1):

$$\mathbf{x}^{\mathsf{b}}(t_i) = M(t_i, t_0)[\mathbf{x}^{\mathsf{b}}(t_0)]$$
(2.3)

where

$$M(t_i, t_0) \equiv M(t_i, t_{i-1}) \cdots M(t_1, t_0)$$
(2.4)

is the  $\mathbb{R}^n \to \mathbb{R}^n$  model propagator from  $t_0$  to  $t_i$ . Equation (2.3) shows explicitly the role of the initial state vector in determining the model state trajectory. We refer to the initial state vector as the *control* vector of the assimilation problem.

For notational simplicity, let  $\mathbf{x} \equiv \mathbf{x}(t_0)$  and  $\mathbf{x}^{\mathrm{b}} \equiv \mathbf{x}^{\mathrm{b}}(t_0)$ . The nonlinear 4D-Var problem is then defined as

$$\min_{\mathbf{x}\in\mathbb{R}^n} J[\mathbf{x}] = \underbrace{\frac{1}{2} (\mathbf{x} - \mathbf{x}^{\mathrm{b}})^{\mathrm{T}} \mathbf{B}^{-1} (\mathbf{x} - \mathbf{x}^{\mathrm{b}})}_{J_{\mathrm{b}}} + \underbrace{\frac{1}{2} (G(\mathbf{x}) - \mathbf{y}^{\mathrm{o}})^{\mathrm{T}} \mathbf{R}^{-1} (G(\mathbf{x}) - \mathbf{y}^{\mathrm{o}})}_{J_{\mathrm{o}}}$$
(2.5)

where  $G : \mathbb{R}^n \to \mathbb{R}^p$  is a nonlinear operator mapping the initial vector into the space of the observation vector, and  $\mathbf{B} \in \mathbb{R}^{n \times n}$  and  $\mathbf{R} \in \mathbb{R}^{p \times p}$  are symmetric, positive-definite matrices containing estimates of the background- and observation-error covariances, respectively. By splitting the background and observation terms  $(J_b \text{ and } J_o)$ , we are tacitly assuming that the errors in the background state are uncorrelated with those in the observation vector.

The nonlinear operator  $G = G(\mathbf{x})$  is given by

$$G(\mathbf{x}) = \begin{pmatrix} G_0(\mathbf{x}) \\ \vdots \\ G_i(\mathbf{x}) \\ \vdots \\ G_N(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} H_0(\mathbf{x}) \\ \vdots \\ H_i(M(t_i, t_0)(\mathbf{x})) \\ \vdots \\ H_N(M(t_N, t_0)(\mathbf{x})) \end{pmatrix}$$
(2.6)

where  $H_i : \mathbb{R}^n \to \mathbb{R}^{p_i}$  is the observation operator that transforms  $\mathbf{x}(t_i)$  into the measured quantity at the appropriate measurement location, and  $M(t_i, t_0)$  is given by Eq. (2.4). We refer to  $G_i(\mathbf{x}) = H_i(M(t_i, t_0)(\mathbf{x})) : \mathbb{R}^n \to \mathbb{R}^{p_i}$  as a generalized (nonlinear) observation operator for measurements at time  $t_i$ .

The optimal trajectory through the assimilation window  $t_0 \leq t_i \leq t_N$  is given by

$$\mathbf{x}^{\mathbf{a}}(t_i) = M(t_i, t_0)(\mathbf{x}^{\mathbf{a}})$$

where  $\mathbf{x}^{a} \equiv \mathbf{x}^{a}(t_{0})$  is the global minimizing solution of Eq. (2.5) and is referred to as the *analysis*.

### 2.2 Incremental formulation and algorithm

Let  $\mathbf{x}^{(k-1)}$  be a reference state and let  $\delta \mathbf{x}^{(k)}$  be an increment to the reference state such that

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)}.$$

Here we consider the special case when the resolution of the state vector is the same as that of the increment vector. In terms of the increment, the minimization problem defined by Eq. (2.5) can be rewritten as

$$\min_{\delta \mathbf{x}^{(k)} \in \mathbb{R}^{n}} J[\delta \mathbf{x}^{(k)}] = \frac{1}{2} \left( \mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)} - \mathbf{x}^{b} \right)^{\mathrm{T}} \mathbf{B}^{-1} \left( \mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)} - \mathbf{x}^{b} \right) 
+ \frac{1}{2} \left( G(\mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)}) - \mathbf{y}^{o} \right)^{\mathrm{T}} \mathbf{R}^{-1} \left( G(\mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)}) - \mathbf{y}^{o} \right). (2.7)$$

The minimization problems (2.5) and (2.7) are equivalent.

The incremental algorithm is an iterative algorithm defined by the minimization of a sequence,  $k = 1, ..., K_0$ , of *quadratic* cost functions. The quadratic cost function on iteration k is obtained by linearizing G about the estimate obtained from the previous iteration k - 1:

$$G(\mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)}) \approx G(\mathbf{x}^{(k-1)}) + \widetilde{\mathbf{G}}^{(k-1)} \delta \mathbf{x}^{(k)}$$
(2.8)

where

$$\widetilde{\mathbf{G}}^{(k-1)} = \begin{pmatrix} \widetilde{\mathbf{G}}_{0}^{(k-1)} \\ \vdots \\ \widetilde{\mathbf{G}}_{i}^{(k-1)} \\ \vdots \\ \widetilde{\mathbf{G}}_{N}^{(k-1)} \end{pmatrix}$$

is an approximation to the tangent-linear operator

$$\mathbf{G}^{(k-1)} \equiv \begin{pmatrix} \partial G_0 / \partial \mathbf{x} |_{\mathbf{x}=\mathbf{x}^{(k-1)}} \\ \vdots \\ \partial G_i / \partial \mathbf{x} |_{\mathbf{x}=\mathbf{x}^{(k-1)}} \\ \vdots \\ \partial G_N / \partial \mathbf{x} |_{\mathbf{x}=\mathbf{x}^{(k-1)}} \end{pmatrix} = \begin{pmatrix} \mathbf{H}_0^{(k-1)} \\ \vdots \\ \mathbf{H}_i^{(k-1)} \mathbf{M}^{(k-1)}(t_i, t_0) \\ \vdots \\ \mathbf{H}_N^{(k-1)} \mathbf{M}^{(k-1)}(t_N, t_0) \end{pmatrix}$$

where  $\mathbf{M}^{(k-1)}(t_i, t_0)$  and  $\mathbf{H}_i^{(k-1)}$  denote the tangent-linear of the model propagator and observation operator, respectively. Using the approximation Eq. (2.8), Eq. (2.7) can be transformed into the quadratic minimization problem

$$\min_{\delta \mathbf{x}^{(k)} \in \mathbb{R}^{n}} J^{(k)}[\delta \mathbf{x}^{(k)}] = \frac{1}{2} (\delta \mathbf{x}^{(k)} - \delta \mathbf{x}^{\mathbf{b},(k-1)})^{\mathrm{T}} \mathbf{B}^{-1} (\delta \mathbf{x}^{(k)} - \delta \mathbf{x}^{\mathbf{b},(k-1)}) 
+ \frac{1}{2} (\widetilde{\mathbf{G}}^{(k-1)} \delta \mathbf{x}^{(k)} - \delta \mathbf{y}^{\mathbf{o},(k-1)})^{\mathrm{T}} \mathbf{R}^{-1} (\widetilde{\mathbf{G}}^{(k-1)} \delta \mathbf{x}^{(k)} - \delta \mathbf{y}^{\mathbf{o},(k-1)}) (2.9)$$

where

$$\delta \mathbf{x}^{\mathbf{b},(k-1)} = \mathbf{x}^{\mathbf{b}} - \mathbf{x}^{(k-1)} \tag{2.10}$$

can be interpreted as the background estimate of  $\delta \mathbf{x}^{(k)}$ , and

$$\delta \mathbf{y}^{\mathrm{o},(k-1)} = \begin{pmatrix} \delta \mathbf{y}_{0}^{\mathrm{o},(k-1)} \\ \vdots \\ \delta \mathbf{y}_{i}^{\mathrm{o},(k-1)} \\ \vdots \\ \delta \mathbf{y}_{N}^{\mathrm{o},(k-1)} \end{pmatrix} = \begin{pmatrix} \mathbf{y}_{0}^{\mathrm{o}} - G_{0}(\mathbf{x}^{(k-1)}) \\ \vdots \\ \mathbf{y}_{i}^{\mathrm{o}} - G_{i}(\mathbf{x}^{(k-1)}) \\ \vdots \\ \mathbf{y}_{N}^{\mathrm{o}} - G_{N}(\mathbf{x}^{(k-1)}) \end{pmatrix}$$

is a four-dimensional vector of observation-minus-reference state differences which, using Eqs. (2.2) and (2.6), can be written in compact form as

$$\delta \mathbf{y}^{\mathbf{o},(k-1)} = \mathbf{y}^{\mathbf{o}} - G(\mathbf{x}^{(k-1)}).$$

The iterations indexed by  $k = 1, ..., K_0$  are called *outer* iterations while the quadratic minimization iterations performed on each outer iteration are called *inner* iterations.

On the first outer iteration, the reference state is set to the background state,

$$\mathbf{x}^{(0)} = \mathbf{x}^{\mathbf{b}},\tag{2.11}$$

and hence from Eq. (2.10)

$$\delta \mathbf{x}^{\mathbf{b},(0)} = \mathbf{0}.$$

On the first inner iteration of each outer iteration, the increment is initialized to zero,

$$\delta \mathbf{x}_{(0)}^{(k)} = \mathbf{0}$$

where the subscript (j) denotes the inner iteration counter, with  $j = 0, \ldots, m_k$ . The subscript k on  $m_k$  indicates that the total number of inner iterations may be different from one outer iteration to the next. After the final inner iteration on each outer iteration, the reference state is updated with the increment,

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \delta \mathbf{x}^{(k)}_{(m_k)}.$$
(2.12)

For the next outer iteration k, the background estimate of the increment can be written, using Eqs (2.10)–(2.12), in terms of the sum of the increments from all previous outer iterations,

$$\delta \mathbf{x}^{\mathbf{b},(k-1)} = -\sum_{l=1}^{k-1} \delta \mathbf{x}_{(m_l)}^{(l)}.$$

The *analysis*  $\mathbf{x}^{a}$  for the incremental problem is taken to be the reference state after the final outer iteration, which can be written as

$$\mathbf{x}^{\mathrm{a}} = \mathbf{x}^{\mathrm{b}} + \delta \mathbf{x}^{\mathrm{a}}$$

where

$$\delta \mathbf{x}^{\mathbf{a}} = \sum_{k=1}^{K_{\mathbf{o}}} \delta \mathbf{x}_{(m_k)}^{(k)} = -\delta \mathbf{x}^{\mathbf{b},(K_{\mathbf{o}})}$$

is the analysis increment.

### 2.3 Role of the $J_{\rm b}$ term and B matrix

In a cycled data assimilation experiment, the background state on each cycle represents the integrated response of all information (observations, forcing fields,...) used on previous cycles. Clearly this is valuable information that should be assimilated along with the observations. Moreover, in the usual case when the number of observations is less than the number of control variables, the background state is fundamental information for regularizing the assimilation problem. The assimilation of the background state is done through the  $J_{\rm b}$  term in Eq. (2.5). In this term, an estimate of the inverse of the background-error covariance matrix,  $\mathbf{B}^{-1}$ , is required to weight the background state departures. In practice, the **B** matrix is used as a preconditioner which allows the minimization problem to be reformulated in such a way that the inverse **B** matrix (or operator) never needs to be specified explicitly. **B**-preconditioning is commonly achieved via a change of variables which, although avoiding the need to specify  $\mathbf{B}^{-1}$ , requires the availability of 'square-root' factors **U** and  $\mathbf{U}^{\rm T}$  such that  $\mathbf{B} = \mathbf{U}\mathbf{U}^{\rm T}$  (Parrish and Derber 1992, Derber and Bouttier 1999). **B**-preconditioning can also be achieved without the need for a square root, but specially adapted conjugate gradient (CG) algorithms are required (**Gürol et al. 2013**).

The optimal preconditioner is the inverse of the Hessian matrix in the sense that it results in convergence of the minimization in a single iteration with a gradient descent method. The **B** preconditioner corresponds to the inverse Hessian of the  $J_{\rm b}$  term only. In the **B**-preconditioned space, the eigenvalues of the Hessian matrix are bounded below by 1 and have a cluster of eigenvalues at 1 with size at least max (0, n - p) (**Tshimanga** et al. 2008). When  $p \ll n$ , which is typically the case in ocean data assimilation, this preconditioner can significantly improve the convergence properties of the minimization. For example, with a single observation (p = 1), the minimization will converge in one iteration of a CG method.

The **B** matrix plays an important role in determining the spatial structure and amplitude of the analysis increments. To see this, consider the exact solution of the inner loop problem (Eq. (2.9)). The outer iteration superscript (k) will be dropped for clarity of notation. Setting the gradient of J to zero yields the solution

$$\delta \mathbf{x}^{a} = \delta \mathbf{x}^{b} + \left( \mathbf{B}^{-1} + \widetilde{\mathbf{G}}^{T} \mathbf{R}^{-1} \widetilde{\mathbf{G}} \right)^{-1} \widetilde{\mathbf{G}}^{T} \mathbf{R}^{-1} \left( \delta \mathbf{y}^{o} - \widetilde{\mathbf{G}} \delta \mathbf{x}^{b} \right), \qquad (2.13)$$

which using the Sherman-Morrison-Woodbury formula (Golub and van Loan 1996) can be rewritten in the alternative form

$$\delta \mathbf{x}^{\mathbf{a}} = \delta \mathbf{x}^{\mathbf{b}} + \mathbf{B} \underbrace{\widetilde{\mathbf{G}}^{\mathrm{T}}}_{\boldsymbol{\alpha}} \underbrace{\left( \underbrace{\widetilde{\mathbf{G}} \, \mathbf{B} \, \widetilde{\mathbf{G}}^{\mathrm{T}} + \mathbf{R} \right)^{-1} \left( \delta \mathbf{y}^{\mathrm{o}} - \widetilde{\mathbf{G}} \delta \mathbf{x}^{\mathrm{b}} \right)}_{\boldsymbol{\beta}}}_{\boldsymbol{\alpha}}. \tag{2.14}$$

Each column vector  $\mathbf{b}_i$  of  $\mathbf{B}$  corresponds to the error covariance of the background field at grid-point *i* with the background field at all grid points i = 1, ..., n. Letting  $\alpha_i$  denote the *i*th component of the *n*-dimensional vector  $\boldsymbol{\alpha}$  highlighted by the underbrace in Eq (2.14) then the analysis increment can be written as

$$\delta \mathbf{x}^{\mathbf{a}} = \delta \mathbf{x}^{\mathbf{b}} + \sum_{i=1}^{n} \alpha_i \mathbf{b}_i.$$

Since  $\delta \mathbf{x}^{\mathbf{b}}$  is a linear combination of increments from previous outer iterations, it too can be written as a linear combination of  $\mathbf{b}_i$ . In other words, the analysis increment is in the column space of  $\mathbf{B}$ , thus clearly illustrating the importance of the latter.

If we back up one operation in Eq. (2.14) then we can also write the analysis increment as a linear combination of the columns  $\mathbf{r}_j$ ,  $j = 1, \ldots, p$ , of the generally rectangular matrix  $\mathbf{B}\widetilde{\mathbf{G}}^{\mathrm{T}}$ ,

$$\delta \mathbf{x}^{\mathrm{a}} = \delta \mathbf{x}^{\mathrm{b}} + \sum_{j=1}^{p} \beta_{j} \mathbf{r}_{j}$$

where  $\beta_j$  denotes the *j*th element of the *p*-dimensional vector  $\boldsymbol{\beta}$  highlighted in Eq (2.14). Each column vector  $\mathbf{r}_j$  corresponds to the error covariance of the observation at location *j* with the background field at all grid points i = 1, ..., n. The  $\mathbf{r}_j$  vectors are often referred to as the *representers* (Bennett 2002). This expression clearly illustrates the role of the adjoint operator  $\tilde{\mathbf{G}}^{\mathrm{T}}$ , as well as **B**, in determining the analysis increment.

A diagonal **B** matrix, which would contain information only about the backgrounderror variances, can help control the amplitude of the analysis increment but is ineffective at controlling small-scale noise that can result from the assimilation of sparse observations. This is particularly true in 3D-Var which, with a diagonal **B** matrix, would produce bullet-like analysis increments near observation locations. In 3D-Var, the backgrounderror correlations, which correspond to the off-diagonal elements of **B**, will be primarily responsible for spreading the influence of an observation away from its measurement location and between model variables. In 4D-Var, the adjoint dynamics will also contribute to the propagation of observational information, the extent to which will depend on the width of the assimilation window as well as the underlying dynamical processes.

## Chapter 3

# Developments in 4D-Var for the ocean

### **3.1** 4D-Var in a simplified framework

Satellite altimeters provide near-global, time-continuous measurements of the ocean surface topography and are a key component of the ocean observing system (Fu and Cazenave 1991). Indeed, much of the early interest in ocean data assimilation was inspired by advances in satellite altimetry. A particularly intriguing question for oceanographers was the following. Since large-scale variations in the ocean surface topography are a manifestation of large-scale dynamical processes occurring below the surface, could altimeter measurements combined with a data assimilation method be sufficient to determine the subsurface ocean circulation? In many ways, the problem is analogous to the early attempts made by meteorologists to determine the 3D circulation of the atmosphere from surface pressure measurements (Bengtsson 1979). The problem of extracting 3D information from surface measurements is arguably more important in oceanography than in meteorology, however, since direct measurements that probe the fluid's vertical structure are much more sparse and difficult to make for the ocean than for the atmosphere.

It was with this question in mind that Weaver and Anderson (1997) explored the potential of 4D-Var to project altimeter data onto the subsurface fields in an ocean model. The region of interest was the tropical Pacific Ocean. It is in this region that ocean data assimilation is important for initializing climate forecasts on seasonal-to-interannual time scales. The experimental framework was kept deliberately simple in order to answer a few fundamental questions. How much information about the ocean state can we extract given a time sequence of perfect altimeter data, an accurate numerical model and forcing field, and an assimilation scheme that makes optimal use of the dynamics? In particular, how well are the subsurface fields constrained in models with different vertical resolution, and how sensitive are the results to the lengths of the assimilation and data sampling period?

The ocean model was a linear reduced-gravity model with an arbitrary number of active layers. The prognostic model variables were the layer heights and the horizontal velocity components in each layer. The equatorial waves supported by the model equations, most importantly baroclinic Rossby and Kelvin waves, play a major role in the seasonal and interannual variability of the tropical ocean circulation. As a result, when forced by a realistic windstress field, simple models of this type are capable of reproducing much of the observed tropical ocean variability on these timescales (Cane 1979; Busalacchi and O'Brien 1981). Close variants of this model were used as the ocean component in simple coupled atmosphere-ocean models in the pioneering efforts to forecast the El Niño– Southern Oscillation (ENSO) (McCreary 1983; McCreary and Anderson 1984).

The experiments were of the simulated-data (identical-twin) type. The assimilated sea surface height (SSH) 'observations' consisted of complete maps on regular intervals. The cost function measured the sum of quadratic ('energy') misfits between the model and 'observed' SSH fields at different times. In the model, SSH is a diagnostic variable defined as a linear combination of the perturbation layer depths. The adjoint equations used to compute the gradient of the cost function with respect to the initial height and velocity fields were derived with respect to an 'energy' inner product. With this choice of inner product, the inverse of the matrix of 'energy' weights can be interpreted as a preconditioning matrix for the gradient derived with respect to the canonical inner product. In these experiments there was no background term in the cost function since it was our intent to concentrate on the efficiency of the dynamics for improving the subsurface fields in the absence of prior statistical information.

As an example, Fig. 3.1 shows the height-field errors in the third layer of a threeactive layer model, from experiments without (upper panel) and with (lower panel) 4D-Var assimilation of SSH. The errors are shown at the end of a one-year period. In the 4D-Var experiment, the SSH 'observations' were assimilated every ten days using a oneyear assimilation window. This figure illustrates very clearly that information about the subsurface ocean state is more easily extracted from SSH near the equator than in offequatorial regions. The crucial factor governing the success of the assimilation is the phase separation that develops between the different baroclinic modes over the observation interval. In off-equatorial regions, where the variability is dominated by slowly propagating baroclinic Rossby waves, the highest modes can only produce significant phase separations after several years or even decades depending on how many modes are present in the system. On the other hand, SSH information is much more easily transferred to depth near the equator where the faster propagation speeds of large-scale waves enable greater phase separations to develop between the higher modes on a much shorter timescale. Decreasing the data sampling period to one day had little effect other than to reduce small-scale errors primarily close to the western boundary where the model variability is dominated by strongly dissipated, short Rossby waves.

While the results from this study showed, albeit in highly favorable conditions, that 4D-Var is an effective method for dynamically projecting altimeter data to depth within the equatorial belt, they also indicated that the task becomes significantly more difficult as soon as the vertical resolution is extended to include more than a couple of layers (baroclinic modes), especially away from the equator. Without additional information, assimilation experiments of considerably longer duration than a year must be considered to reconstruct the subsurface fields adequately from altimeter data alone. The problem would clearly be further complicated in real-data applications due to the presence of errors in the model, observations and forcing fields.

There are two obvious ways where the effectiveness of the altimeter assimilation could be improved. The first way is to assimilate altimeter data simultaneously with other data, such as temperature and salinity profiles which provide direct measurements of the ocean's



Figure 3.1: The layer 3 height field error in the control (upper panel) and assimilation (lower panel) experiments at the end of a one-year 4D-Var experiment with full-field SSH observations every 10 days. The contour interval is 5 metres. (From Weaver and Anderson (1997).)

density structure. From this perspective, the altimeter assimilation problem should be viewed within the more general framework of operational data assimilation or reanalysis; both aim at determining the best possible ocean analysis by combining all available information. The second way is to include a background term in the 4D-Var cost function in order to exploit prior knowledge about the multivariate statistical properties of the background state. This information is contained in the estimate of the background-error covariance matrix and complements the time-dependent dynamical information brought by the model through the observation term. The impact of a multivariate formulation of the background-error covariance matrix on the assimilation of altimeter data is discussed in chapter 4.

### 3.2 4D-Var in a realistic framework

Whereas simple models such as the one used by **Weaver and Anderson (1997)** are helpful for exploring new assimilation methods and understanding basic processes, they have severe limitations when applied to real data. Real-data assimilation applications should employ state-of-the-art GCMs which contain the most complete and accurate representations of the physical processes governing the ocean circulation. To employ an advanced assimilation method such as 4D-Var with a GCM in a real-data application requires significantly more development and validation than with a simple model in a simulated-data framework. In addition to the model and assimilation method, procedures are required for data handling and quality control, for *cycling* the model from one assimilation window to the next, and for evaluating the assimilation performance through a comprehensive set of diagnostics. Together, the different components comprise the data assimilation *system*.

Motivated by the growing interest in ENSO forecasting and the importance of initializing coupled ocean-atmosphere GCMs for making ENSO forecasts, I began the development of an incremental 4D-Var system for the OPA ocean GCM (Madec et al. 1999) as a postdoctoral fellow at LODYC in Paris. At the time, OPA was a relatively new GCM and had mainly been applied to process studies in regional basins such as the tropical Pacific Ocean and Mediterranean Sea. Central to a 4D-Var system are tangent-linear (TL) and adjoint models. These needed to be developed for the latest version of OPA (version 7.0 at the time). While I spent considerable effort deriving the TL and adjoint models for OPA, this was an excellent way to learn about the details of the ocean GCM and ultimately was beneficial when it came to develop other components of the assimilation system. By the time the first version of the incremental 4D-Var system, originally referred to as OPAVAR, was ready for real-data experimentation (nearly 4 years later), OPA had gone through two major upgrades (versions 8.0 and 8.1). A third major upgrade (version 8.2) was introduced near year 2000 to include, among others, important developments for the global configuration (ORCA), such as the introduction of a free surface (Roullet and Madec 2000). Substantial changes to the OPAVAR system were required after each model upgrade, such as rederiving the TL and adjoint models. This delayed scientific experimentation, but kept OPAVAR in phase with model releases, thereby enabling it to benefit from the latest model improvements and to be attractive to other developers and users.

Initially the 4D-Var system was developed for a tropical Pacific configuration of OPA and applied to the assimilation of *in situ* temperature observations from the Global Temperature and Salinity Pilot Programme. A description of both the system and the results from a 6-year cycled 4D-Var experiment are given in the two-part article by **Weaver** et al. (2003) and Vialard et al. (2003) (a more detailed presentation is given in the technical report of **Weaver** et al. (2002)). These articles also present results from a 3D-Var FGAT system which was derived as a by-product of the incremental 4D-Var system by simply replacing the tangent-linear model with a persistence model. The 3D-Var system provided a useful reference for evaluating the benefits of the more expensive 4D-Var system.

A global version of the 4D-Var system was developed as part of the ENACT<sup>1</sup> project, and adapted to assimilate salinity profiles and along-track altimeter data as well as temperature profiles. The system also supported 3D-Var, and was later extended to produce *ensembles* of 3D-Var analyses in the ENSEMBLES<sup>2</sup> project (**Daget** *et al.* **2009**). Due to

 $<sup>^1\</sup>mathrm{ENhanced}$  ocean data assimilation and ClimaTe prediction: a 3-year project financed under EC-Framework 5.

<sup>&</sup>lt;sup>2</sup>ENSEMBLE-based predictions of climate changes and their impactS: a 5-year project financed under

the high computational cost of running the global 4D-Var system, only the cheaper 3D-Var system was used extensively for analysis production over long periods. The 4D-Var system was used for two specific studies that required only single-cycle experimentation. In particular, Weaver *et al.* (2005) (chapter 4) used the global 4D-Var system to study the impact of multivariate balance operators in 4D-Var, and Tshimanga *et al.* (2008) used it to investigate the effectiveness of different limited-memory preconditioners in a multi-outer iteration framework.

### 3.2.1 Validation of the linear approximation

To assess the validity of the linear approximation in the incremental formulation, Weaver et al. (2003) compared the time-evolution of an initial perturbation in the nonlinear model with its evolution in the linear models: persistence in 3D-Var, and the TL model in 4D-Var. In the assimilation experiments described by Weaver et al. (2003) and Vialard et al. (2003), the width of the assimilation window was taken to be 10 days for 3D-Var and 30 days for 4D-Var, so the linear approximation was verified on these time-scales. Ideally, the initial perturbation should have structure and amplitude typical of actual background errors. Actual background errors are not known so in that study, as a proxy, the initial perturbation was taken to be the difference between two model background states valid at the same time but produced from different initial conditions: one initial state was taken from an analysis 30 days into the past whereas the other state was taken from an analysis 60 days into the past. Figure 3.2a shows a meridional-vertical section of the perturbation in the eastern Pacific. The field is characterized by large perturbations of up to 4°C, appearing as a result of temperature observations that were assimilated between day 60 and day 30 into the past. Figures 3.2b and c show the nonlinear perturbation after 10 days and 30 days respectively, and Fig. 3.2d shows the TL perturbation after 30 days. In the 3D-Var system, the assumption is that the perturbation does not evolve significantly over the 10-day window. By comparing Figs 3.2a and b, we see that this is a very good approximation outside the equatorial belt but has some limitations closer to the equator where the dynamical adjustment time-scales are shorter. On the other hand, over 30-days the TL model provides a good description of large-scale perturbations in both off-equatorial and equatorial regions. The TL approximation was shown to be mainly limited by tropical instability waves that develop in the eastern Pacific particularly during the autumn months. Vertical mixing and convective processes, both of which were highly simplified in the TL model, also tend to limit the TL approximation in the small vertical scales.

The linear approximation can be examined from a different perspective by comparing the values of the incremental (quadratic) and nonincremental (nonquadratic) cost functions at various stages during minimization. Figure 3.3 shows the behaviour of the incremental and nonincremental cost functions in two experiments using the tropical Pacific 4D-Var system of **Weaver** *et al.* (2003). In these experiments, the resolution was identical in the inner and outer loops. Five outer-loop iterations were used in one experiment (solid curve), whereas only one outer-loop iteration was used in the other experiment. The assimilation window was 30 days. The jumps in the solid curves occur after an outer-loop iteration when the reference trajectory is reinitialized (here every ten inner-loop itera-

EC-Framework 6.



Figure 3.2: (a) Meridional-vertical section at  $110^{\circ}$ W of a temperature perturbation used to check the validity of the linear approximation. The perturbation after (b) 10 days and (c) 30 days evolution in the nonlinear model, after (d) 30 days evolution in the TL model. The contour interval is 0.5°C. (From Weaver *et al.* (2003).)

tions). The final values of the nonincremental cost function are plotted with an asterix and plus symbol for the experiment with and without multiple outer-loop iterations, respectively. This figure provides a clear illustration of the positive impact of the outer-loop iterations. The final value of the nonincremental cost function with outer-loop iterations is about half that without multiple outer-loop iterations. Furthermore, it is very close to the final value obtained with the incremental cost function (solid curve) and thus provides a good measure of the consistency of the incremental approach. In the absence of multiple outer-loop iterations, however, this consistency is lost as illustrated by the large discrepancy between the final values of the nonincremental and incremental cost functions (cf. plus symbol and dashed curve).

A similar verification was given in **Tshimanga** *et al.* (2008). Their experiments were conducted on a (shorter) 10-day window and with a *global* version of the 4D-Var system with three outer-loop iterations and ten inner-loop iterations per outer-loop iteration. As in the 4D-Var system of **Weaver** *et al.* (2003), identical resolution was used in the outer and inner loops. As can be seen from Fig. 3.4, the differences between the values of the incremental and nonincremental cost functions at the outer loop iterations (k = 1) where the relative error is 4.5% (for k = 2 the relative error is less than 0.1%). This suggests that the linear approximation is also quite accurate in this 4D-Var experiment. It is worth pointing out that in these experiments the linear approximation also included simplifications in the TL model, the most important being in the representation of vertical and isopycnal slopes



Figure 3.3: The value of the cost function (J) as a function of the minimization iteration for 4D-Var experiments with five outer-loop iterations (solid curve) and one outer-loop iteration (dashed curve). The plus (asterix) symbol indicates the final value of the nonincremental cost function for the experiment with one (five) outer iteration(s). For clarity, these symbols have been displaced slightly to the left of the right border. J has been normalized by its respective value at the start of minimization and plotted on a logarithmic vertical axis. (From Weaver *et al.* (2003).)

were neglected.

### 3.2.2 Flow-dependent background-error variances in 4D-Var

It is well known that, in the limit of a perfect, linear model, variational assimilation is equivalent to the Kalman filter in that, given identical inputs, they produce the same analysis at the end of the assimilation window (Courtier *et al.* 1994). Weaver *et al.* (2003) exploited this theoretical equivalence to assess how in 4D-Var the dynamical model acts to modify the background-error *variances*. In particular, for a specific assimilation window  $t_0 \leq t_i \leq t_N$ , they estimated the diagonal of the matrix

$$\mathbf{P}^{b}(t_{N}) = \mathbf{M}(t_{N}, t_{0}) \mathbf{B} \mathbf{M}(t_{N}, t_{0})^{\mathrm{T}}$$
(3.1)

which results from propagating the initial background-error covariance matrix  $\mathbf{P}^{b}(t_{0}) = \mathbf{B}$ using the tangent-linear propagator and its adjoint. In an extended Kalman filter (with a perfect model), an equation of the form (3.1) would be used explicitly to transport the covariances forward in time. In incremental variational assimilation, this propagation is implicit in the global minimization process.

Figure 3.5, left panel, shows vertical profiles of the background-error standard deviations  $\sigma^b$  for temperature at the equator in the central Pacific (140°W). The dashed-dotted curve is the prior value of  $\sigma^b$ , which has been estimated, albeit rather crudely, from the model climatology in an experiment without data assimilation. It has a rather broad structure throughout the upper ocean and displays a maximum around the depth of the climatological thermocline (near 170 m). In 3D-Var these  $\sigma^b$  are effectively used to weight



Figure 3.4: The values of the quadratic cost function (solid curve) and nonquadratic cost function (open circles) as a function of the inner-loop (conjugate gradient) iteration number in each of the three outer-loop iterations of the unpreconditioned 4D-Var experiment. The curves are placed one after the other in sequence and the inner-loop iterations are cumulated. (From **Tshimanga** *et al.* (2008).)

the background state at all times within the assimilation window, whereas in 4D-Var they are used to weight the background state only at the beginning of the window. The solid curve in Fig. 3.5 shows an example of the effective  $\sigma^b$  used in 4D-Var at the end of a 30-day window. The tangent-linear dynamics tend to reduce the  $\sigma^b$  in the mixed layer and to sharpen the profile around the level of maximum background-error variance. The maximum occurs at the level of the thermocline as confirmed by comparing  $\sigma^b$  to a corresponding 30-day mean profile of the background  $|\partial T/\partial z|$  (Fig. 3.5, right panel). This tendency is physically sensible since the level of maximum variability of the thermal field, and thus of maximum likely error, is located at the level of the thermocline.

### 3.2.3 Improving the background-error covariance model

Comparisons of the 3D-Var and 4D-Var analyses in the tropical Pacific study of **Weaver** et al. (2003) and Vialard et al. (2003) indicated that 4D-Var was superior to 3D-Var in several areas. The fit to the assimilated temperature observations was consistently better in 4D-Var than in 3D-Var. Furthermore, the impact on the state variables not directly observed (the velocity and salinity fields) was generally better in 4D-Var than in 3D-Var. One of the distinguishing features of the 3D-Var analyses was a large bias in the velocity field which was associated with a spurious circulation cell that developed along the equatorial strip (strong downwelling in the eastern Pacific, weaker but broader upwelling in the central/western Pacific, eastward surface currents). Salinity also exhibited a bias, with a drift towards lower values. A bias was present in the temperature field but was greatly reduced on each assimilation cycle by the direct assimilation of subsurface



Figure 3.5: Left panel: vertical profiles of the background-error standard deviations (in °C) used in 3D-Var and 4D-Var at the equator at 140°W. The dashed-dotted curve corresponds to the standard deviations specified at the beginning of the assimilation window. In 3D-Var, these are also the effective standard deviations used at all future times within the window. The solid curve corresponds to the effective standard deviations used implicitly in 4D-Var at the end of the 30-day window in a particular cycle of 4D-Var. Right panel: the corresponding profile of the background  $|\partial T/\partial z|$  at this location. The values of  $|\partial T/\partial z|$  have been multiplied by a factor of ten in order to be plotted with the same horizontal scale as in panel a. (From Weaver *et al.* (2003)).

temperature data.

Many of the deficiencies with the 3D-Var analyses were not fundamental to 3D-Var but could be attributed to an inadequate treatment of the background-error covariances. The **B** matrix used in the 3D-Var experiments was *univariate* in temperature. The error correlations were assumed to be approximately Gaussian and modelled using a diffusion algorithm (Weaver and Courtier 2001; chapter 4). The zonal (meridional) length scales of the Gaussian function were slightly increased (reduced) near the equator to account for simple anisotropic effects due to equatorial dynamics. The vertical scales were set to be a scalar multiple of the local grid depth to provide adequate smoothing between model levels. The variances were allowed to vary with each grid point and specified according to model climatology as explained in the previous section. A similar (univariate) **B** matrix was also used in the 4D-Var experiments, but extended to include additional spatial correlation functions for salinity and the horizontal components of velocity.

The univariate treatment of the background-error covariances was obviously a weak point and was one of the suspected reasons for the large bias that developed in the 3D-Var analyses. Since only the temperature field was corrected, problems occurred as the other fields adjusted to the temperature field during the forward integration of the model. While the **B** matrix was also univariate in 4D-Var, the TL model constraint resulted in salinity and velocity field increments that were partially in balance with the temperature increments, leading to better surface currents and reduced salinity drift. There were still some areas, however, where 4D-Var was not as good as the model simulation *without* data assimilation (underestimated equatorial undercurrent, and unrealistic variability in salinity) which suggested that improvements to **B** could be beneficial to 4D-Var analyses as well.

Several improvements were subsequently made to the **B** model. A fully multivariate formulation was developed which included constraints between, for example, temperature and salinity, and density and currents. Some of these constraints, such as those between temperature and salinity, were state-dependent, and thus were able to account for some (weak) flow dependency in **B**. Inspired by the results of Fig. 3.5, the background-error variances were also made flow-dependent by parameterizing them in terms of the vertical gradient of the background state. The positive impact of these developments to  $\mathbf{B}$  in 3D-Var is illustrated in Fig. 3.6. Shown are the mean equatorial surface currents from the 4D-Var and 3D-Var experiments of Weaver et al. (2003) and Vialard et al. (2003), which used the simpler (univariate, flow-independent) **B** formulation (Figs 3.6b and c), and from a 3D-Var experiment which used an improved (multivariate, flow-dependent) **B** formulation (Fig 3.6d). For comparison, Fig. 3.6a shows the Reverdin *et al.* (1994) climatology which was derived directly from surface drifter and current-meter observations. A striking feature in Fig. 3.6c is the large eastward bias in the equatorial surface currents, as already pointed out above. The main impact of the improved covariance model in 3D-Var is to eliminate this eastward bias and bring the currents much closer to the observed and 4D-Var climatologies.

In this and the previous chapter, I have outlined the role of the **B** matrix and illustrated how improved **B** models can have a significant impact on the quality of ocean analyses. In the next chapter, I provide a more detailed account of the different aspects of **B** modelling related to my published work and to the methods being developed specifically for NEMOVAR.



Figure 3.6: Surface zonal current climatologies in the tropical Pacific Ocean. a) The Reverdin *et al.* (1994) climatology. Climatolgies from the (b) 4D-Var and (c) 3D-Var experiments of **Weaver** *et al.* (2003) and **Vialard** *et al.* (2003). (d) Climatology from a 3D-Var experiment with an improved **B** formulation. The contour interval is  $0.1 ms^{-1}$ , and the blue (yellow) shaded regions indicate westward (eastward) currents.

## Chapter 4

## Developments in background-error covariance modelling for the ocean

In real-life applications the true covariances of background error are never known to great accuracy since there is insufficient information to estimate them. Even if they could be estimated accurately, representing them in a full-rank covariance matrix would not be possible due to its enormous size. Practical techniques for representing backgrounderror covariances in variational data assimilation require simplifying assumptions based on physical insight and computational considerations. Statistical estimation is still important but can only be used sensibly to calibrate a relatively small number of covariance parameters. Useful background-error covariance information can be extracted from an adequately perturbed ensemble data assimilation system (Daget et al. 2009; section 4.3). In variational assimilation, the challenge is then how to synthesize this information effectively in a covariance model that can be applied efficiently on each iteration of the minimization algorithm. Ideally, the covariance model should be designed to capture robust features in the available covariance estimates, such as multivariate relationships, geographical and wavenumber dependencies of the variances and correlations, and anisotropic variations which often occur near regions of pronounced density or topographic gradients, or where the data distribution is discontinuous.

Multivariate relationships between variables give rise to a cross-variable component in the covariances. In variational assimilation, it is common to specify the multivariate component of the covariances through a balance operator that captures known physical or statistical relationships between model variables (Weaver et al. 2005, section 4.2). The inverse of the balance operator can be interpreted as a transformation to a new set of variables whose cross-variable covariances are much weaker than those of the untransformed variables. The cross-covariances of the transformed variables are usually considered sufficiently small that they can be neglected. In doing so, the full multivariate covariance matrix can be decomposed into a sum of simpler univariate covariance matrices acting on each of the transformed variables. Each univariate covariance matrix can in turn be factored into a (symmetric and positive-definite) correlation matrix multiplied to its left and right by a diagonal matrix of standard deviations. The remaining challenge is how to represent adequately and evaluate efficiently the product of a generally full-rank, non-diagonal correlation matrix with an arbitrary vector. This can be done using a diffusion equation (Weaver and Courtier 2001; Weaver and Ricci 2004; Mirouze and Weaver 2010; Weaver and Mirouze 2013) as discussed next.

### 4.1 Univariate correlation modelling using a diffusion equation

The continuous analogue of a correlation matrix–vector product is an integral operator, the kernel of which is a symmetric and positive-definite correlation function. When evaluated at discrete points, the correlation function defines a (full-rank) correlation matrix although this definition is of limited practical interest for large matrix problems. The evaluation of integral equations involving the background-error correlations is among the most costly steps in the variational data assimilation algorithm.

Several methods have been developed in meteorology and oceanography for representing univariate correlation operators (Bannister 2008). Methods based on spectral or wavelet transforms have been developed extensively for atmospheric data assimilation (Fisher 2003; Deckmyn and Berre 2005; Pannekoucke *et al.* 2007). Correlation operators based on physical-space models have also been developed. Physical-space formulations are more convenient than spectral or wavelet formulations in complex-boundary domains such as those encountered in ocean data assimilation. Gaspari and Cohn (1999) and Gneiting (2002) derive families of parameterized homogeneous and isotropic correlation functions with the important property of compact support so that they can be integrated efficiently in grid-point space. Gaspari *et al.* (2006) describe extensions to these functions to account for spatially varying length-scales. Correlation integrals can also be represented efficiently using grid-point filters. The recursive filter (Lorenc 1992; Purser *et al.* 2003a,b) and the diffusion equation (Derber and Rosati 1989; Egbert *et al.* 1994; Weaver and Courtier 2001; Mirouze and Weaver 2010; Weaver and Mirouze 2013) fall into this class of correlation model.

Derber and Rosati (1989) proposed the use of an iterative Laplacian grid-point smoother in order to approximate a Gaussian correlation operator. Egbert et al. (1994) and Bennett et al. (1997) described a close variant of the algorithm in which the Laplacian smoothing could be interpreted as a pseudo-time-step integration of a diffusion equation. Weaver and Courtier (2001) described the algorithm in more detail and proposed various extensions to account for more general correlation functions than the quasi-Gaussian of the original Derber and Rosati (1989) algorithm. Weaver and Ricci (2004) and Mirouze and Weaver (2010) studied implicit versions of the diffusion algorithm. In preparing this manuscript, I embarked on providing a comprehensive review of the theory underpinning the explicit- and implicit-diffusion methods for modelling correlation functions. The review material ultimately became part of an article (Weaver and **Mirouze 2013)** in the Quarterly Journal of the Royal Meteorological Society. The details are not repeated here (the reader is referred to Appendix B where a copy of the article can be found). The remainder of this section will focus mainly on the practical construction of 2D and 3D correlation models using a diffusion operator formulated with an implicit scheme, with reference to methods developed specifically for NEMOVAR as part of the PhD thesis of Mirouze (2010).

### 4.1.1 Two-dimensional correlation models

Coordinate systems of global ocean models are referenced to the spherical-shell geometry of the ocean. From a mathematical perspective, this leads naturally to consider 2D 'hor-

izontal' correlation functions on the spherical space  $\mathbb{S}^2$ . The product of a 2D correlation function on  $\mathbb{S}^2$  and a 1D correlation function on the bounded subset of the Euclidean space  $\mathbb{R}^1$  is commonly used to construct 3D correlation functions on the spherical-shell subspace of  $\mathbb{R}^3$ . This approach of separating the horizontal and vertical correlation functions is usually justified by the fact that the global ocean circulation is characterized by scales that are much larger in the horizontal direction (along geopotential surfaces) than in the vertical direction (perpendicular to geopotential surfaces).

To fix the ideas, we focus on the use of a diffusion operator to represent a 2D background-error correlation model on a discretized grid on  $\mathbb{S}^2$ . Let  $i = 1, \ldots, N$  denote the number of grid-points and let  $\boldsymbol{\psi}$  be the  $N \times 1$  vector containing the grid-point values of a background variable  $\psi = \psi(\lambda, \phi)$  where  $\lambda$  is longitude ( $0 \leq \lambda \leq 2\pi$ ) and  $\phi$  is latitude ( $-\pi/2 \leq \phi \leq \pi/2$ ). The background-error covariance matrix **B** for  $\boldsymbol{\psi}$  can be formulated as

$$\mathbf{B} = \mathbf{D}^{1/2} \, \mathbf{\Gamma}^{1/2} \, \mathbf{L} \, \mathbf{W}^{-1} \mathbf{\Gamma}^{1/2} \, \mathbf{D}^{1/2} \tag{4.1}$$

where  $\mathbf{D} = \mathbf{D}^{1/2} \mathbf{D}^{1/2}$  is a diagonal matrix of background-error variances ( $\mathbf{D}^{1/2}$  is the diagonal matrix of background-error standard deviations),  $\mathbf{L}$  is the matrix representation of a positive-definite smoothing operator, and  $\mathbf{\Gamma} = \mathbf{\Gamma}^{1/2} \mathbf{\Gamma}^{1/2}$  is a diagonal matrix of normalization factors that transforms  $\mathbf{C} = \mathbf{\Gamma}^{1/2} \mathbf{L} \mathbf{W}^{-1} \mathbf{\Gamma}^{1/2}$  into a matrix with 1s along the diagonal in accordance with a correlation matrix. The smoothing operator depends on the model coordinate system and grid resolution. It is constructed to be self-adjoint with respect to the inner product  $\langle \cdot, \cdot \rangle_{\mathbf{W}} = (\cdot)^{\mathrm{T}} \mathbf{W}(\cdot)$ , where  $\mathbf{W}$  is a diagonal matrix of discrete metric coefficients. Denoting the adjoint of  $\mathbf{L}$  by  $\mathbf{L}^*$  then the self-adjointness property implies that

$$\mathbf{L} = \mathbf{L}^* = \mathbf{W}^{-1} \mathbf{L}^{\mathrm{T}} \mathbf{W}, \qquad (4.2)$$

and hence that **B** is symmetric in the usual sense  $(\mathbf{B} = \mathbf{B}^{\mathrm{T}})$ .

For some applications, such as preconditioning (**Tshimanga** *et al.* 2008) or randomization (**Weaver and Courtier 2001**), it is desirable to have access to a "square-root" operator **U** such that  $\mathbf{B} = \mathbf{U}\mathbf{U}^{\mathrm{T}}$ . Providing the smoothing operator can be factored as  $\mathbf{L} = \mathbf{L}^{1/2}\mathbf{L}^{1/2}$  then from Eqs (4.1) and (4.2) it follows that

$$\mathbf{U} = \mathbf{D}^{1/2} \, \mathbf{\Gamma}^{1/2} \, \mathbf{L}^{1/2} \, \mathbf{W}^{-1/2}.$$

Weaver and Courtier (2001) discuss how a diffusion operator L can be used to represent the action of a positive-definite covariance matrix. Ignoring continental boundaries, the solution of the 2D diffusion equation on  $S^2$ ,

$$\frac{\partial \psi}{\partial s} - \kappa \nabla^2 \psi = 0 \quad ; \quad \psi(\phi, \lambda, 0) = \widehat{\psi}(\phi, \lambda) \tag{4.3}$$

where

$$\nabla^2 \psi \equiv \frac{1}{a^2 \cos \phi} \frac{\partial}{\partial \phi} \left( \cos \phi \frac{\partial \psi}{\partial \phi} \right) + \frac{1}{a^2 \cos^2 \phi} \frac{\partial^2 \psi}{\partial \lambda^2},$$

has the general form

$$\psi(\lambda,\phi,s) = \int_{\mathbb{S}^2} c_{\mathbf{s}}(\theta;\kappa s) \,\widehat{\psi}(\lambda',\phi') \,a^2 \cos\phi' \,\mathrm{d}\lambda' \,\mathrm{d}\phi'. \tag{4.4}$$

Here s is to be interpreted as a dimensionless pseudo-time coordinate. The diffusion coefficient  $\kappa$  then has physical units of length squared.

The kernel  $c(\theta; \kappa s)$  of the integral operator (4.4) is an *isotropic* covariance function that approximates a Gaussian function:

$$c_{\rm s}(\theta;\kappa s) \approx c_{\rm g}(r,D_{\rm g}) \propto \exp\left(-r^2/2D_{\rm g}^2\right)$$

$$(4.5)$$

where

$$= r(\theta) = a\sqrt{2(1-\cos\theta)}$$

r

is the chordal distance between points  $(\lambda, \phi)$  and  $(\lambda', \phi')$  separated by an angle  $\theta$   $(0 \le \theta \le$  $\pi$ ) on the sphere of radius a. The length-scale,  $D_{\rm g}$ , of the Gaussian function is defined here in terms of the local curvature of the covariance function at its peak. This is a standard definition in meteorology and oceanography, and is sometimes referred to as the Daley length-scale in reference to Daley's book (Daley 1991) where it is discussed within the context of data assimilation. In Eq. (4.5),  $D_{\rm g}$  is related to the product of the diffusion coefficient and the total action time s of the diffusion:

$$D_{\rm g} \approx \sqrt{2\kappa s}.$$

Therefore, to approximate the action of a Gaussian correlation operator on  $S^2$ , Eq. (4.3) can be solved numerically and normalized by  $c_s(0)$  to obtain the correct (unit) amplitude. In this case, the normalization matrix in Eq. (4.1) is simply a constant multiple of the identity matrix,  $\Gamma \approx \gamma_{\rm s} \mathbf{I}$ , where

$$\gamma_{\rm s} = (c(0))^{-1} \approx 2\pi D_{\rm g}^2.$$

Weaver and Ricci (2004) and Weaver and Mirouze (2013) discuss how a "time"-implicit the diffusion model on  $\mathbb{S}^2$  can be used to define a more general and robust covariance model than the "time"-explicit formulation proposed by Weaver and **Courtier (2001)**. Define the linear operator  $\mathcal{A}: \psi \mapsto \psi$  where

$$\mathcal{A}\psi \equiv \left(1 - L^2 \nabla^2\right)\psi,\tag{4.6}$$

L being a constant scale parameter. Equation (4.6) is a roughening operator and can be derived by discretizing the "temporal" derivative in Eq. (4.3) with an Euler-backward implicit scheme and setting the diffusion coefficient  $\kappa = L^2$ . The inverse operator  $\mathcal{A}^{-1}$ , on the other hand, is a *smoothing* operator. Applying the  $\mathcal{A}$  operator M times can be interpreted as inverting an *M*-step implicitly-formulated diffusion operator  $\mathcal{L} \equiv (\mathcal{A}^M)^{-1} =$  $(\mathcal{A}^{-1})^M$  where  $\mathcal{A}^{-1}: \widetilde{\psi} \mapsto \psi$ . The elliptic equation  $\mathcal{A}^M \psi = \widehat{\psi}$  yields a general solution of the form

$$\psi(\lambda,\phi) = \int_{\mathbb{S}^2} c_{\rm h}(\theta;L,M) \,\widehat{\psi}(\lambda',\phi') \,a^2 \cos \phi' \,\mathrm{d}\lambda' \,\mathrm{d}\phi', \tag{4.7}$$

where the kernel  $c_{\rm h}(\theta; L, M)$  is an isotropic covariance function with characteristics very similar to those of a Whittle-Matérn or Matérn function  $c_{\rm w}(r; L, M)$  on  $\mathbb{R}^2$  (Guttorp and Gneiting 2006):

$$c_{\rm h}(\theta; L, M) \approx c_{\rm w}(r; L, M) \propto \left(\frac{r}{L}\right)^{M-1} K_{M-1}\left(\frac{r}{L}\right)$$

where  $K_{M-1}(r/L)$  is the modified Bessel function of the second kind of order M-1. The parameters L and M control the length-scale and spectral-decay rate of the covariance function. From the standard (Daley) definition of length-scale, it can be shown that

$$D_{\rm h} \approx D_{\rm w} = \sqrt{2M - 4} L \tag{4.8}$$

where M > 2 to exclude non-differentiable Matérn functions. Furthermore, the normalization matrix is  $\Gamma \approx \gamma_{\rm h} \mathbf{I}$  where

$$\gamma_{\rm h} \approx 4\pi (M-1)L^2 = \frac{4\pi (M-1)}{2M-4}D_{\rm w}^2.$$
 (4.9)

The integral equation (4.7) can be evaluated numerically by solving the linear system of equations

$$\begin{array}{rcl}
\mathbf{A}\boldsymbol{\psi}_{1} &= \boldsymbol{\psi} \\
\mathbf{A}\boldsymbol{\psi}_{2} &= \boldsymbol{\psi}_{1} \\
&\vdots \\
\mathbf{A}\boldsymbol{\psi}_{M} &= \boldsymbol{\psi}_{M-1}
\end{array}$$

$$(4.10)$$

where **A** is a self-adjoint, positive-definite matrix constructed from a discretized version of the linear operator  $\mathcal{A}$ . The self-adjointness of **A** is defined with respect to the **W**inner product where, assuming a constant zonal and meridional resolution  $\Delta\lambda$  and  $\Delta\phi$ ,  $\mathbf{W} = \text{diag}(a^2 \cos \phi_i \Delta \lambda \Delta \phi)$ , and hence  $\mathbf{A} = \mathbf{A}^* = \mathbf{W}^{-1} \mathbf{A}^{\mathrm{T}} \mathbf{W}$ . The smoothing operator in (4.1) is then

$$\mathbf{L} = (\mathbf{A}^{-1})^M$$

where  $\mathbf{L}\hat{\boldsymbol{\psi}} = \boldsymbol{\psi}_M$ . Symmetric matrices are preferable to self-adjoint matrices when used with standard linear system solvers. The system matrix in (4.10) can be easily made symmetric by left-multiplying both sides of the equations by  $\mathbf{W}$ . This yields the smoothing operator

$$\mathbf{L} = (\widehat{\mathbf{A}}^{-1} \mathbf{W})^M$$

where  $\widehat{\mathbf{A}} = \mathbf{W}\mathbf{A} = \widehat{\mathbf{A}}^{\mathrm{T}}$ . Notice that by restricting M to be even, it is straightforward to define a square-root operator  $\mathbf{L}^{1/2}$  as a diffusion operator over M/2 steps:  $\mathbf{L}^{1/2} = (\widehat{\mathbf{A}}^{-1}\mathbf{W})^{M/2}$ .

Applying a direct matrix solver to a 2D implicit diffusion problem such as (4.10) requires access to algebraic software libraries specially designed for large, sparse matrices. For example, **Weaver and Ricci (2004)** and **Ricci (2004)** used the HSL Mathematical Software Library to solve a generalized 2D implicit diffusion problem for a data assimilation application with a regional ocean configuration, run on a single processor machine. However, adapting this method to the massively-parallel environment of the global-ocean configuration of OPA/NEMO was not feasible, so alternative methods were explored.

For 1D implicit diffusion problems, the linear system (4.10) can be solved efficiently using a standard matrix solver such as Cholesky decomposition (Golub and van Loan 1996). The covariance kernels implied by 1D implicit diffusion in  $\mathbb{R}$  belong to the class of autoregressive (AR) functions (**Mirouze and Weaver 2010**). Inspired by the work of Purser *et al.* (2003a,b) with the 1D recursive filter, **Mirouze and Weaver (2010)** and Mirouze (2010) suggest approximating the 2D implicit diffusion operator by a product of 1D implicit diffusion operators,

$$\mathbf{L} \approx \mathbf{L}_{\phi} \mathbf{L}_{\lambda}$$
  
=  $(\mathbf{A}_{\phi}^{-1})^{M} (\mathbf{A}_{\lambda}^{-1})^{M}$   
=  $(\mathbf{A}_{\phi}^{-1} \mathbf{A}_{\lambda}^{-1})^{M}$  (4.11)

$$= (\mathbf{A}_{\lambda} \mathbf{A}_{\phi})^{-M} \tag{4.12}$$

where the matrices  $\mathbf{A}_{\lambda}$  and  $\mathbf{A}_{\phi}$  are discrete representations of the 1D operators

$$\mathcal{A}_{\lambda}\psi \equiv \left(1 - \frac{L^2}{a^2\cos^2\phi}\frac{\partial^2}{\partial\lambda^2}\right)\psi,\tag{4.13}$$

$$\mathcal{A}_{\phi}\psi \equiv \left(1 - \frac{L^2}{a^2\cos\phi}\frac{\partial}{\partial\phi}\left(\cos\phi\frac{\partial}{\partial\phi}\right)\right)\psi.$$
(4.14)

In terms of symmetric matrices  $\widehat{\mathbf{A}}_{\phi} = \mathbf{W}_{\phi}\mathbf{A}_{\phi}$  and  $\widehat{\mathbf{A}}_{\lambda} = \mathbf{W}_{\lambda}\mathbf{A}_{\lambda}$ , where  $\mathbf{W}_{\lambda} = \operatorname{diag}(a\cos\phi_{i}\Delta\lambda)$ ,  $\mathbf{W}_{\phi} = \operatorname{diag}(a\Delta\phi)$ , and  $\mathbf{W} = \mathbf{W}_{\phi}\mathbf{W}_{\lambda}$ , the smoothing operator becomes

$$\mathbf{L} = (\widehat{\mathbf{A}}_{\phi}^{-1} \mathbf{W}_{\phi} \, \widehat{\mathbf{A}}_{\lambda}^{-1} \mathbf{W}_{\lambda})^{M} = (\widehat{\mathbf{A}}_{\phi}^{-1} \, \widehat{\mathbf{A}}_{\lambda}^{-1} \mathbf{W})^{M}.$$
(4.15)

The matrix product  $\mathbf{A}_{\lambda}\mathbf{A}_{\phi}$  in Eq. (4.12) is a discrete representation of the 2D roughening operator (cf. Eq. (4.6))

$$\mathcal{A}_{\lambda}\mathcal{A}_{\phi}\psi = \mathcal{A}\psi + \mathcal{A}_{\lambda\phi}'\psi$$

where

$$\mathcal{A}'_{\lambda\phi}\psi = \frac{L^4}{a^4\cos^3\phi}\frac{\partial}{\partial\phi}\left(\cos\phi\frac{\partial}{\partial\phi}\right)\frac{\partial^2\psi}{\partial\lambda^2}.$$

The  $\mathcal{A}'_{\phi\lambda}$  term can be interpreted as the error resulting from approximating  $\mathcal{A}$  by a product of 1D operators. Notice that the error is different if we change the order of the operations:

$$\mathcal{A}_{\phi}\mathcal{A}_{\lambda}\psi = \mathcal{A}\psi + \mathcal{A}_{\phi\lambda}^{\prime}\psi$$

where

$$\mathcal{A}_{\phi\lambda}^{\prime}\psi = \frac{L^4}{a^4\cos\phi}\frac{\partial}{\partial\phi}\left(\cos\phi\frac{\partial}{\partial\phi}\left(\frac{1}{\cos^2\phi}\frac{\partial^2\psi}{\partial\lambda^2}\right)\right) \neq \mathcal{A}_{\lambda\phi}^{\prime}\psi$$

since, unlike in the Euclidean space  $\mathbb{R}^2$ , the component derivatives do not commute in a non-Euclidean space such as  $\mathbb{S}^2$ . As a consequence, the self-adjointness of **L** on  $\mathbb{S}^2$  is not preserved by this construction.

Another problem with this construction is that it gives rise to a spurious anisotropic response characterized by correlations that are diamond-shaped instead of circular. (This occurs in  $\mathbb{R}^2$  as well as  $\mathbb{S}^2$ ). However, since the problematic error terms are of order  $M^{-2}$ 

for a fixed length-scale  $D_{\rm h}$  (Eq. (4.8)), they can be effectively suppressed by making M sufficiently large; i.e., by seeking an approximate Gaussian response since  $c_{\rm h} \approx c_{\rm g}$  and  $D_{\rm h} \approx D_{\rm g}$  for large M. In practice,  $M \sim 10$  gives an adequate isotropic response. A 2D correlation model based on a product of 1D implicit diffusion operators has been developed for NEMOVAR by **Mirouze (2010)**. Some of the issues involved in the practical implementation of this method are discussed in the following sections.

### 4.1.2 Solid-wall boundary conditions

Now consider Eq. (4.3) in the presence of solid-wall boundaries such as coastlines in an ocean model. Different boundary conditions (BCs) have been considered by Weaver and Courtier (2001) and Mirouze and Weaver (2010): setting the normal derivative of the field to be zero at the boundary (Neumann BCs); setting the field itself to be zero at the boundary (Dirichlet BCs); or a mixture of the two (Robin BCs). Regardless of the type of BC employed, the resulting response of the diffusion operator near the boundary leads to a large change in the amplitude of the covariance function. As a consequence, the constant normalization factor  $\gamma_h$  becomes a poor estimate of the correct normalization factors near the boundaries. In numerical applications, the correct normalization factor at a particular grid-point can be computed *exactly* by applying  $\mathbf{L}$  to the discrete vector representation of a Dirac  $\delta$ -function centred at that point; i.e., a vector  $\delta_i$  equal to zero everywhere except at the grid-point i where it is defined by the inverse of the local area element. (The spatial integral of  $\delta_i$  would then produce a value of 1 as required by the formal definition of a  $\delta$ function.) The value of the smoothed field at grid-point i is the inverse of the normalization factor at that point. This algorithm, which requires as many applications of the diffusion operator as there are grid-points, is expensive and of limited practical interest for large problems. Alternatively, the normalization factors can be *approximated* using a cheaper algorithm based on an ensemble of random vectors (Weaver and Courtier 2001). The accuracy of this method is determined by the number of random vectors employed.

While the normalization process can correct the amplitude of the covariance function, it results in a distortion of the shape and effective length-scale of the covariance functions near the boundaries. Mirouze and Weaver (2010) provide examples in 1D that illustrate this effect. They go on further to propose a formulation of the diffusion operator that effectively renders the operator *transparent* to solid-wall boundaries. This can be done by defining L to be the average of two diffusion operators, one employing Neumann (N) BCs and the other employing Dirichlet (D) BCs:

$$\mathbf{L} = \frac{1}{2} (\mathbf{L}_{\mathrm{N}} + \mathbf{L}_{\mathrm{D}})$$
$$= \left(\frac{1}{\sqrt{2}} \mathbf{L}_{\mathrm{N}}^{1/2} \frac{1}{\sqrt{2}} \mathbf{L}_{\mathrm{D}}^{1/2}\right) \left(\begin{array}{c} \frac{1}{\sqrt{2}} \mathbf{L}_{\mathrm{N}}^{1/2} \\ \frac{1}{\sqrt{2}} \mathbf{L}_{\mathrm{D}}^{1/2} \end{array}\right)$$
(4.16)

where the second expression illustrates that the square-root associated with  $\mathbf{L}$  is now a rectangular matrix operator. A mathematical proof of this result is given in **Mirouze** and Weaver (2010) for the Gaussian case described by the 1D diffusion equation. Equation (4.16) follows from an intuitive generalization of this result to higher dimensions and to account for the (finite M) Matérn-like functions represented by the implicit diffusion operator. Its validity has been demonstrated numerically in several examples in 1D as well as in global configurations of NEMOVAR. While this formulation doubles the cost of the correlation model, it may be important for applications requiring accurate analyses near the boundary. Note that with **L** defined by Eq. (4.16), the constant normalization factor  $\gamma_{\rm h}$  is a good approximation of the correct normalization factor near boundaries as well as far from boundaries.

#### 4.1.3 Inhomogeneity and non-separability

The assumption that background-error correlation length-scales are constant is very restrictive in ocean data assimilation. Ocean variability occurs on a wide-range of time and spatial scales, and has a strong regional dependence. For example, ocean variability is dominated by mesoscale eddies in western boundary current regions such as the Gulf Stream, whereas it is dominated by larger scale, linear wave dynamics in the equatorial regions. Furthermore, the ocean observing system consists of a largely inhomogeneous distribution of measurements. For example, there is an abundance of measurements in certain regions of strategic or economic interest, such as the Gulf Steam or tropical Pacific, but far fewer measurements in the more hostile and isolated environment of the Southern oceans. The inhomogeneous nature of the ocean observing system and the regional diversity of ocean dynamics will result in a geographical dependence in the background-error correlations which should be accounted for.

Now consider the case when we allow for spatial variations in the scale parameter  $L = L(\lambda, \phi)$ . Local estimates of the Daley length-scales  $D = D(\lambda, \phi)$  can be obtained from statistics of ensemble-forecast differences (Belo Pereira and Berre 2006; Pannekoucke and Massart 2008; **Daget 2008**; see section 4.3). Given estimates of D and a prescribed value of M, the local diffusion scale L can be determined from Eq. (4.8). To account for spatially-varying scale parameters in the implicit diffusion model, the parameter  $L^2$  must be introduced within the derivatives in order to ensure that  $\mathcal{A}$  remains self-adjoint with respect to the inner product  $\langle \psi_1, \psi_2 \rangle_W = \iint \psi_1 \psi_2 a^2 \cos \phi \, d\lambda \, d\phi$ , and hence that the correlation functions implied by  $\mathcal{L} = (\mathcal{A}^{-1})^M$  are symmetric. Equation (4.6) then becomes

$$\mathcal{A}\psi \equiv \left(1 - \nabla \cdot L^2 \nabla\right)\psi \tag{4.17}$$

where  $\nabla$  denotes the horizontal gradient operator and  $\nabla$  the horizontal divergence operator.

A consequence of varying the scale parameters is that the normalization factors are no longer constant. When the scale variations are sufficiently slow compared to the scale itself, a reasonable approximation to the normalization factor at each grid-point can be obtained from Eq. (4.9) using the local value of  $L^2$  (Pannekoucke and Massart 2008; **Mirouze and Weaver 2010**). A better approximation may be obtained using the same expression but with filtered estimates of  $L^2$  (Purser *et al.* 2003b; **Mirouze and Weaver 2010**; Yaremchuk and Carrier 2012). Alternatively, the more expensive randomization or  $\delta$ function methods described earlier can be applied to obtain more accurate estimates.

Let us now return to the formulation (4.15) where the 2D implicit diffusion operator is separated into a product of 1D implicit diffusion operators, and let the matrix operators
$\mathbf{A}_{\lambda}$  and  $\mathbf{A}_{\phi}$  denote discrete representations of

$$\mathcal{A}_{\lambda}\psi = \left(1 - \frac{1}{a^{2}\cos^{2}\phi}\frac{\partial}{\partial\lambda}\left(L^{2}(\lambda,\phi)\frac{\partial}{\partial\lambda}\right)\right)\psi$$
  
and 
$$\mathcal{A}_{\phi}\psi = \left(1 - \frac{1}{a^{2}\cos\phi}\frac{\partial}{\partial\phi}\left(L^{2}(\lambda,\phi)\cos\phi\frac{\partial}{\partial\phi}\right)\right)\psi,$$

which are the 1D implicit diffusion operators derived from Eq. (4.17) (cf. Eqs (4.13) and (4.14)).

While  $\mathbf{A}_{\lambda}$  and  $\mathbf{A}_{\phi}$  are self-adjoint with respect to the  $\mathbf{W}_{\lambda}$ - and  $\mathbf{W}_{\phi}$ -inner products, respectively, their product  $\mathbf{A}_{\lambda}\mathbf{A}_{\phi}$  (or  $\mathbf{A}_{\phi}\mathbf{A}_{\lambda}$ ) is not, in general, self-adjoint with respect to the **W**-inner product:

$$\begin{aligned} \mathbf{A}_{\lambda} \mathbf{A}_{\phi} &= \mathbf{W}_{\lambda}^{-1} \mathbf{A}_{\lambda}^{\mathrm{T}} \mathbf{W}_{\lambda} \mathbf{W}_{\phi}^{-1} \mathbf{A}_{\phi}^{\mathrm{T}} \mathbf{W}_{\phi} \\ &= \mathbf{W}^{-1} \mathbf{A}_{\lambda}^{\mathrm{T}} \mathbf{A}_{\phi}^{\mathrm{T}} \mathbf{W} \\ &= \mathbf{W}^{-1} (\mathbf{A}_{\phi} \mathbf{A}_{\lambda})^{\mathrm{T}} \mathbf{W} \\ &= (\mathbf{A}_{\phi} \mathbf{A}_{\lambda})^{*} \\ &\neq (\mathbf{A}_{\lambda} \mathbf{A}_{\phi})^{*} \quad \text{(in general).} \end{aligned}$$

This means that self-adjointness of the associated smoothness operator  $\mathbf{L}$  has to be forced explicitly. There are several ways to do this as outlined below, but some have more attractive properties than others.

To simplify notation, let

$$\mathbf{L}_{\phi\lambda} = \mathbf{L}_{\phi} \mathbf{L}_{\lambda}$$
$$\mathbf{L}_{\phi\lambda}^{*} = \mathbf{L}_{\lambda}^{*} \mathbf{L}_{\phi}^{*} = \mathbf{L}_{\lambda} \mathbf{L}_{\phi}$$
$$\mathbf{L}_{\phi\lambda}^{1/2} = \mathbf{L}_{\phi}^{1/2} \mathbf{L}_{\lambda}^{1/2}.$$
$$\left. \right\}$$
(4.18)

The formulation

where

$$\mathbf{L}_{(1)} = \frac{1}{2} \left( \mathbf{L}_{\phi\lambda} + \mathbf{L}_{\phi\lambda}^* \right), \qquad (4.19)$$

obtained by averaging  $\mathbf{L}_{\phi\lambda}$  with its adjoint, is an obvious way to impose self-adjointness. However, it doubles the number of diffusion operations and does not have a simple squareroot factorization. It also has unsatisfactory smoothness properties near complex boundaries as discussed later and illustrated by **Mirouze (2010)** 

More convenient formulations can be derived using square-root factorizations of  $\mathbf{L}_{\phi\lambda}$ and  $\mathbf{L}_{\phi\lambda}^*$ . Using the relations (4.18) in Eq. (4.19) and rearranging terms leads to an alternative (self-adjoint) formulation

$$\mathbf{L}_{(2)} = \frac{1}{2} \left( \mathbf{L}_{\phi\lambda}^{1/2} \mathbf{L}_{\phi\lambda}^{1/2*} + \mathbf{L}_{\phi\lambda}^{1/2*} \mathbf{L}_{\phi\lambda}^{1/2} \right)$$
  
=  $\mathbf{L}_{(2)}^{1/2} \mathbf{L}_{(2)}^{1/2*}$  (4.20)

where

$$\mathbf{L}_{(2)}^{1/2} = \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbf{L}_{\phi\lambda}^{1/2} & \mathbf{L}_{\phi\lambda}^{1/2*} \end{pmatrix}.$$

Equation (4.20) also doubles the number of diffusion operations but, contrary to Eq. (4.19), possesses a relatively simple "square-root"  $\mathbf{L}_{(2)}^{1/2}$ . The rectangular nature of  $\mathbf{L}_{(2)}^{1/2}$ , however, means that the size of the control vector needed for minimization with a square-root preconditioner will effectively be doubled. (This situation also arises from the treatment of the boundary conditions following Eq. (4.16)). A simple way to avoid this is to average  $\mathbf{L}_{\phi\lambda}^{1/2}$  and  $\mathbf{L}_{\phi\lambda}^{1/2}$ \* directly within the definition of the square root. This leads to the following formulation:

$$\mathbf{L}_{(3)} = \frac{1}{4} \left( \mathbf{L}_{\phi\lambda}^{1/2} \mathbf{L}_{\phi\lambda}^{1/2} + \mathbf{L}_{\phi\lambda}^{1/2} \mathbf{L}_{\phi\lambda}^{1/2*} + \mathbf{L}_{\phi\lambda}^{1/2*} \mathbf{L}_{\phi\lambda}^{1/2} + \mathbf{L}_{\phi\lambda}^{1/2*} \mathbf{L}_{\phi\lambda}^{1/2*} \right), = \mathbf{L}_{(3)}^{1/2} \mathbf{L}_{(3)}^{1/2*}$$
(4.21)

where

$$\mathbf{L}_{(3)}^{1/2} = \frac{1}{2} \left( \mathbf{L}_{\phi\lambda}^{1/2} + \mathbf{L}_{\phi\lambda}^{1/2}^{*} \right).$$

Eq. (4.21) contains the two terms in Eq. (4.20) plus two extra terms associated with the cross product in Eq. (4.21). When implemented using the factorization (4.21),  $\mathbf{L}_{(3)}$  has an equivalent number of diffusion operations as  $\mathbf{L}_{(2)}$ . Dobricic and Pinardi (2008) proposed a similar construction for representing 3D covariances from the product of a horizontal covariance operator based on the recursive filter and a vertical covariance matrix derived from Empirical Orthogonal Functions.

The additional diffusion terms required by formulations  $\mathbf{L}_{(1)}$ ,  $\mathbf{L}_{(2)}$  and  $\mathbf{L}_{(3)}$  will increase the computational cost of the correlation operator. This may be a problem for some applications. Two cheaper formulations can be derived by considering each self-adjoint component of Eq. (4.20) separately:

$$\mathbf{L}_{(4)} = \mathbf{L}_{\phi\lambda}^{1/2} \left( \mathbf{L}_{\phi\lambda}^{1/2} \right)^* \tag{4.22}$$

and 
$$\mathbf{L}_{(5)} = \left(\mathbf{L}_{\phi\lambda}^{1/2}\right)^* \mathbf{L}_{\phi\lambda}^{1/2}.$$
 (4.23)

Both Eqs (4.22) and (4.23) are valid correlation operators, and have an easily accessible square root, but choosing between one or the other is somewhat arbitrary.

Formulations  $\mathbf{L}_{(p)}$ ,  $p = 1, \ldots, 5$ , are identical when  $\mathbf{L}_{\lambda}$  and  $\mathbf{L}_{\phi}$  commute, which occurs for the special case when the correlation kernels of  $\mathbf{L}_{\lambda}$  and  $\mathbf{L}_{\phi}$  are independent of  $\phi$  and  $\lambda$ , respectively, and M is large. We can therefore expect them to produce similar results in regions where the correlation function  $C_{\lambda\phi}(\lambda,\phi,\lambda',\phi')$  implied by  $\mathbf{L}_{\lambda}\mathbf{L}_{\phi}$  can be well approximated by a separable function  $C_{\lambda}(\lambda,\lambda')C_{\phi}(\phi,\phi')$ . They can, however, produce quite different results in regions where non-separability is important. For example, this can occur near complex coastlines where a particular order of the diffusion operations may result in some regions being "smoothed" more than others. In data assimilation, this can lead to unphysical gradients in the analysis increments, particularly when the correlation scale is much larger than the scale of the local geometry. As illustrated by **Mirouze (2010)**,

38



Figure 4.1: Point correlations obtained using a different number  $(m_s)$  of alternating directions in a 2D correlation operator based on Eqs (4.22) and (4.24) but applied on a uniform Cartesian grid. The correlations are computed with respect to a point immediately to the left of an island of the size of one grid-point. The Daley length-scale is equal to 10 grid-points. The number of implicit diffusion iterations (M) equals 10.

the effect is especially acute with the formulation  $\mathbf{L}_{(1)}$ , and least problematic with the formulations  $\mathbf{L}_{(2)}$  and  $\mathbf{L}_{(3)}$  which have the greatest number of diffusion applications with alternating directions. Even these latter formulations may not be sufficient to remove spurious features. In this case, it would be desirable to increase the number of alternating diffusion directions within the individual square-root operators. This can be done by formulating the square-root operators as

$$\mathbf{L}_{\phi\lambda}^{1/2} = \left(\mathbf{L}_{\phi}^{1/2m_{\rm s}}\mathbf{L}_{\lambda}^{1/2m_{\rm s}}\right)^{m_{\rm s}}$$
(4.24)

where  $m_s$  is a positive integer that is chosen to be a factor of M. The maximum possible number of alternating diffusion directions would be produced by choosing  $m_s = M/2$ . Figures 4.1a and b illustrate the effect of the parameter  $m_s$  on minimizing numerical artifacts in the correlations near an isolated boundary point, here taken to be a small island whose size is 10 times smaller than the scale of the correlation structures. In practice, an appropriate choice of  $m_s$  will also depend on the details of the implicit solver and its implementation on massively parallel machines. For example, the solution algorithm currently used in NEMOVAR involves a Cholesky algorithm, combined with a procedure to remap the rectangular domains of the parallel decomposition used in NEMO into banded domains that are suitable with the forward-elimination and backward-substitution steps of the Cholesky algorithm (**Mirouze 2010**). For this algorithm, it is desirable to keep  $m_s$  sufficiently small to limit the number of costly remapping steps needed to go from one decomposition to the next.

#### 4.1.4 Three-dimensional correlation models

A 3D correlation operator that includes the vertical dimension z can be constructed as a self-adjoint product of three 1D implicit diffusion operators  $\mathbf{L}_{\lambda}$ ,  $\mathbf{L}_{\phi}$  and  $\mathbf{L}_{z}$ , following the basic approach described above. In the 3D case, there are six possible combinations of  $\mathbf{L}_{\lambda}$ ,  $\mathbf{L}_{\phi}$  and  $\mathbf{L}_{z}$  and therefore six possible self-adjoint square-root formulations (cf. Eqs (4.22) and (4.23) in the 2D case). As in the 2D case, to avoid spurious gradients in the correlations near complex bathymetry, it may be necessary to perform additional smoothing between the horizontal and vertical planes, using an "interleaving" formulation analogous to (4.24). In NEMO, the parallel decomposition is done in the horizontal domain, so that a switch from either  $\mathbf{L}_{\lambda}$  or  $\mathbf{L}_{\phi}$  to  $\mathbf{L}_{z}$  (or vice-versa) does not require a reorganization of processors as required by a switch from  $\mathbf{L}_{\lambda}$  to  $\mathbf{L}_{\phi}$  (or vice-versa).

The numerical artifacts near coastlines and bathymetry can be eliminated naturally through the use of a non-separable smoothing operator constructed from a true 3D implicit diffusion operator  $\mathcal{L}_{3D} \equiv (\mathcal{A}_{3D}^{-1})^M$  where

$$\mathcal{A}_{3\mathrm{D}}\psi \equiv \left(1 - \nabla \cdot L_{\mathrm{h}}^{2}(\lambda,\phi,z) \nabla - \frac{\partial}{\partial z} \left(L_{\mathrm{z}}^{2}(\lambda,\phi,z) \frac{\partial}{\partial z}\right)\right)\psi.$$
(4.25)

Here, the horizontal and vertical scale parameters,  $L_{\rm h}$  and  $L_{\rm v}$ , are treated as functions of the three spatial coordinates. The operator  $\mathcal{A}_{\rm 3D}$  is positive definite and self-adjoint with respect to the inner product  $\langle \psi_1, \psi_2 \rangle = \int \int \int \psi_1 \psi_2 a^2 \cos \phi \, d\lambda \, d\phi \, dz$ . Applying a fully 3D implicit diffusion-based correlation operator is computationally challenging for large systems, and to my knowledge has not been attempted yet in practical data assimilation systems. It is an attractive possibility for future development however, especially in view of its importance for representing anisotropic correlations as discussed next.

### 4.1.5 Anisotropy

Isotropic correlation models, or quasi-isotropic correlation models that allow stretching of the correlations only in the direction of the computational coordinates, are commonly used in data assimilation algorithms because of their simplicity and computational convenience. There is no reason, however, to expect actual background-error correlations to be isotropic or quasi-isotropic in geophysical fluids such as the ocean. On the contrary, one would expect them to be strongly anisotropic, particularly near coastlines, bathymetry, ocean fronts or in unevenly observed regions. Anisotropic correlation models allow for preferential stretching or shrinking of the correlation functions along arbitrary directions. Anisotropy can be taken into account in the diffusion equation by replacing the diffusion coefficient with a diffusion tensor. With reference to the 3D implicit diffusion (cf. Eq. (4.25)), the  $\mathcal{A}_{3D}$  operator becomes

$$\mathcal{A}_{3D}\psi \equiv (1 - \nabla \cdot \boldsymbol{L}(\lambda, \phi, z)\nabla)\psi \qquad (4.26)$$

where  $\nabla$  is the 3D divergence operator,  $\nabla$  is the 3D gradient operator, and L is a symmetric (diffusion) *tensor* which is, in general, a function of the three spatial coordinates.

In 3D the symmetric diffusion tensor contains six independent elements. Four of the tensor elements account for anisotropy between the horizontal and vertical directions. The importance of these terms compared to the diagonal terms is related to the choice of vertical coordinate in the correlation model. In an ocean model, for example, a natural vertical coordinate is a hybrid coordinate involving a standard geopotential (z) coordinate in unstratified regions such as the mixed layer, an isopycnal  $(\rho)$  coordinate in strongly stratified regions, and a terrain-following (s) coordinate near the ocean bottom, the latter

being particularly important in shallow coastal regions (Haidvogel and Beckmann 1999). In this hybrid coordinate system, the flow is more naturally decoupled into 'horizontal' and 'vertical' processes. If the same coordinate system is adopted for a background-error correlation model then it is reasonable to assume, at least from a physical viewpoint, that the non-diagonal tensor elements are small and can be neglected. However, anisotropy in background-error correlations can also arise from the assimilation of data, especially near the transition between poorly observed and well observed regions. In general, the relative importance of the diagonal and non-diagonal terms of the tensor can only be determined after a thorough diagnostic study involving the direct estimation of the elements of the tensor.

Many ocean models used for global- and basin-scale circulation studies employ a z coordinate. Weaver and Courtier (2001) showed how a standard isopycnal diffusion tensor used to parameterize mixing of unresolved processes in a z-coordinate ocean model could also be used to transform the coordinates of a background-error correlation model formulated as an explicit 3D diffusion operator. The resulting operator produces correlations that are strongest along the background isopyncal surfaces but fall off rapidly across these surfaces. This is illustrated in Fig. 4.2 which compares correlation structures from a 3D diffusion operator in z coordinates and isopycnal coordinates. A correlation model based on Fig. 4.2c would clearly be less destructive to the background density profile than a correlation model based on Fig. 4.2b. Furthermore, being state-dependent, isopycnal coordinates would allow the background-error correlation model to evolve from one cycle to the next, rather than being fixed as in z coordinates.

An analogous coordinate transformation was proposed within the framework of Optimal Interpolation by Balmaseda *et al.* (2008). While the isopycnal correlation model has appealing features, the implementation based on the explicit scheme proposed by **Weaver and Courtier (2001)** is too expensive for routine applications since a prohibitively high number of iterations is required to maintain numerical stability in regions of strong isopycnal gradients. Moreover, the specification of the length-scales must be performed in isopycnal space, which makes estimating them more difficult in a *z*-coordinate model. **Weaver and Mirouze (2013)** propose alternative methods for defining anisotropic correlations, which involve estimating the tensor directly in the model coordinate system. This is discussed within the context of ensemble data assimilation in section 4.3.2.

### 4.2 Multivariate covariance modelling using balance operators

Ocean state vectors involve a mixture of variables. The primary variables of an OGCM are potential temperature (T), salinity (S), SSH  $(\eta)$  and the two components (u and v) of the horizontal velocity vector. These variables are highly coupled through the physical relationships described by the governing equations of the OGCM. Background errors can thus be expected to be highly correlated between variables. Accounting for cross-variable covariances in the background-error formulation is important in allowing the assimilation system to extract information about unobserved variables from observed variables, and, as already illustrated in Chapter 3, in helping to minimize unphysical adjustment processes that can occur when the model is initialized from an analysis.



Figure 4.2: (a) A meridional section of a typical background potential temperature field in the eastern tropical Pacific. (b) The auto-correlation field at a depth of 55 metres and latitude 8°N generated by a 3D diffusion operator defined with respect to the geopotential coordinate system. (c) The corresponding auto-correlation field obtained using the 3D diffusion operator defined with respect to an isopycnal coordinate system based on the background isopycnal surfaces associated with (a). (From Weaver and Courtier 2001).

The next subsection describes how the general 4D-Var problem described in Chapter 2 can be modified to take into account multivariate relationships in the background-error formulation. This will then be followed by a brief description of the specific multivariate relationships proposed by **Ricci** *et al.* (2005) and **Weaver** *et al.* (2005) for variational ocean data assimilation, and which form the basis of the multivariate covariance formulation used in NEMOVAR. A close variant of this formulation is also used in the 4D-Var system for the ROMS model (Moore *et al.* 2011a,b,c).

### 4.2.1 A general transformation of variables

Central to the multivariate covariance formulation is the assumption that the initial state variables can be transformed into a new set of variables whose cross-covariances are sufficiently small that they can be ignored. The basic idea was originally proposed by Derber and Bouttier (1999) for atmospheric data assimilation and developed in an oceanographic context by **Weaver** et al. (2005).

We refer to the notation of Chapter 2 where the general 4D-Var problem was introduced. The transformation of the initial state vector  $\mathbf{x} = \mathbf{x}(t_0)$  into a vector  $\mathbf{w}$  of approximately uncorrelated variables will be denoted as

$$\mathbf{w} = K^{-1}(\mathbf{x})$$

where  $K^{-1} : \mathbb{R}^n \to \mathbb{R}^n$  is assumed to be invertible and possibly nonlinear. The forward operator K is usually called a balance operator. In terms of **w**, the nonlinear 4D-Var

problem (2.5) can be recast as finding  $\mathbf{w}^a = \arg \min J[\mathbf{w}]$  where

$$J[\mathbf{w}] = \frac{1}{2} \left( \mathbf{w} - \mathbf{w}^{\rm b} \right)^{\rm T} \mathbf{B}_{(\mathbf{w})}^{-1} \left( \mathbf{w} - \mathbf{w}^{\rm b} \right) + \frac{1}{2} \left( G[K(\mathbf{w})] - \mathbf{y}^{\rm o} \right)^{\rm T} \mathbf{R}^{-1} \left( G[K(\mathbf{w})] - \mathbf{y}^{\rm o} \right)$$
(4.27)

where

$$\mathbf{w}^{\mathbf{b}} = K^{-1}(\mathbf{x}^{\mathbf{b}}) \tag{4.28}$$

is the background estimate of  $\mathbf{w}$ . The analysis in model space is given by

$$\mathbf{x}^{\mathbf{a}} = K(\mathbf{w}^{\mathbf{a}}). \tag{4.29}$$

By virtue of the 'decorrelation' transformation (4.28), the error covariance matrix  $\mathbf{B}_{(\mathbf{w})}$ of the transformed background-state vector  $\mathbf{w}^{b}$  can be assumed to be block-diagonal (univariate). The individual blocks correspond to the error covariance matrices of each of the transformed variables. These univariate covariance matrices can be modelled separately using a diffusion operator as described in the previous section. Whereas the model propagator  $M(t_i, t_0)$  in the general expression for G (Eq. (2.6)) constrains the sequence of model states between  $t_0$  and  $t_i$ , the balance operator K in Eq. (4.27) provides a complementary constraint on the state at initial time. Although K can be absorbed into the generalized observation operator G, it is convenient here to keep them separate in Eq. (4.27) in order to clarify an approximation made in NEMOVAR, as outlined below.

The incremental approximation of (4.27) leads to the sequence  $k = 1, ..., K_{o}$  of quadratic minimization problems (cf. (2.9)) for finding  $\delta \mathbf{w}^{a} = \arg \min J[\delta \mathbf{w}]$  where

$$J^{(k)}[\delta \mathbf{w}^{(k)}] = \frac{1}{2} (\delta \mathbf{w}^{(k)} - \delta \mathbf{w}^{\mathbf{b},(k-1)})^{\mathrm{T}} \mathbf{B}_{(\mathbf{w})}^{-1} (\delta \mathbf{w}^{(k)} - \delta \mathbf{w}^{\mathbf{b},(k-1)})$$

$$+ \frac{1}{2} (\widetilde{\mathbf{G}}^{(k-1)} \widetilde{\mathbf{K}}^{(k-1)} \delta \mathbf{w}^{(k)} - \delta \mathbf{y}^{\mathbf{o},(k-1)})^{\mathrm{T}} \mathbf{R}^{-1} (\widetilde{\mathbf{G}}^{(k-1)} \widetilde{\mathbf{K}}^{(k-1)} \delta \mathbf{w}^{(k)} - \delta \mathbf{y}^{\mathbf{o},(k-1)})$$
(4.30)

where

$$\delta \mathbf{w}^{\mathbf{b},(\mathbf{k}-1)} = \mathbf{w}^{\mathbf{b}} - \mathbf{w}^{(k-1)},$$
  

$$\delta \mathbf{y}^{\mathbf{o},(\mathbf{k}-1)} = \mathbf{y}^{\mathbf{o}} - G[K(\mathbf{w}^{(k-1)})]$$
(4.31)

and  $\widetilde{\mathbf{K}}^{(k-1)}$  is an approximation of the tangent-linear of the balance operator,  $\partial K/\partial \mathbf{w}|_{\mathbf{w}=\mathbf{w}^{(k-1)}}$ . The analysis increment in control space can be written as a sum of the minimizing increments from all outer iterations (see Chapter 2),

$$\delta \mathbf{w}^{\mathrm{a}} = \sum_{k=1}^{K_{\mathrm{o}}} \delta \mathbf{w}_{(m_k)}^{(k)} = -\delta \mathbf{w}^{\mathrm{b},(K_{\mathrm{o}})}.$$

The analysis in model space can then be obtained from Eq. (4.29) with

$$\mathbf{w}^{\mathrm{a}} = \mathbf{w}^{\mathrm{b}} + \delta \mathbf{w}^{\mathrm{a}}.$$

By transforming the problem back to model space, it is easy to see that the effective formulation of the background-error covariance matrix for  $\mathbf{x}^{b}$  has the form

$$\mathbf{B}^{(k-1)} = \widetilde{\mathbf{K}}^{(k-1)} \mathbf{B}_{(\mathbf{w})} (\widetilde{\mathbf{K}}^{(k-1)})^{\mathrm{T}}.$$
(4.32)

The dependence of  $\widetilde{\mathbf{K}}^{(k-1)}$  on the reference state (in the case of a nonlinear balance operator) implies that  $\mathbf{B}^{(k-1)}$  will be updated from one outer iteration to the next.

The nonlinear balance operator is required in Eq. (4.31) to compute the model counterpart of the observation vector given the reference state  $\mathbf{w}^{(k-1)}$ . Through successive linearizations about  $\mathbf{w}^{(l)}$ , l = 0, ..., k - 2, this operator can be approximated by

$$K(\mathbf{w}^{(k-1)}) \approx K(\mathbf{w}^{(0)}) + \sum_{l=1}^{k-1} \mathbf{K}^{(l-1)} \delta \mathbf{w}_{(m_l)}^{(l)}$$
  
=  $\mathbf{x}^{b} + \sum_{l=1}^{k-1} \delta \widetilde{\mathbf{x}}_{(m_l)}^{(l)}$  (4.33)

where

$$\delta \widetilde{\mathbf{x}}_{(m_k)}^{(k)} = \mathbf{K}^{(k-1)} \delta \mathbf{w}_{(m_k)}^{(k)}$$

is an approximation of the model-space increment

$$\delta \mathbf{x}_{(m_k)}^{(k)} = K \Big( \mathbf{w}^{(k-1)} + \delta \mathbf{w}_{(m_k)}^{(k)} \Big) - K \Big( \mathbf{w}^{(k-1)} \Big) \,.$$

With this approximation, only the sequence of *linearized* balance operators  $\mathbf{K}^{(k-1)}$  are required to iterate the incremental algorithm. This approximation has been adopted in NEMOVAR for practical convenience and will be discussed further below.

### 4.2.2 Temperature-salinity constraints

Many of the early applications of data assimilation in OGCMs focussed on the tropical oceans, particularly the Pacific, where there was great interest in producing accurate estimates of the upper ocean state for improving the initialization of forecasts of climate anomalies such as ENSO with coupled GCMs. The backbone of the ENSO observing system, which was established in the mid-1980s to mid-1990s as part of the Tropical Ocean Global Atmosphere (TOGA) programme, consisted of temperature profiles from both XBTs deployed by volunteer observing ships and from ATLAS moorings in the Tropical Atmosphere Ocean (TAO) array. Much of the early efforts in tropical ocean data assimilation were dedicated to assimilating these temperature data using relatively simple *univariate* schemes.

While these schemes led to vastly more accurate estimates of upper ocean thermal state than without data assimilation, they tended to have a detrimental effect on the ocean-state variables (salinity, currents) not directly constrained by the data. The main reason for this was attributed to artificial changes of water masses in the model caused by spurious mixing when the temperature field is changed while leaving salinity, the other primary thermodynamic variable controlling density, unchanged. Multivariate schemes, in particular those that produced corrections to salinity in response to corrections to temperature, were proposed to alleviate this problem. Troccoli and Haines (1999) and Troccoli *et al.* (2002) proposed a two-step nonlinear scheme (hereafter referred to as the TH scheme) that involved first correcting temperature through assimilation and then searching for the analyzed temperature value in the background density profile in order to

retrieve the corresponding background salinity value for defining the salinity "analysis" at that point. The method acts to preserve the background water masses. While the scheme worked remarkably well when assimilating only temperature data in a univariate OI-type system, it was difficult to apply in a more general context such as 4D-Var or with the simultaneous assimilation of multiple data-types.

Ricci et al. (2005) proposed a linearized variant of the TH scheme to fit within the multivariate covariance formulation of Derber and Bouttier (1999). In their scheme, temperature was taken as the primary variable whereas salinity was split into a "balanced" term dependent on temperature, and an "unbalanced" term to account for that part of salinity that is independent of temperature. They derived their linearized scheme from the second-term of a Taylor series expansion of the (assumed) background S(T) relation at each grid-point. A local salinity balance coefficient was then defined from a finitedifference approximation of the background value of  $\partial S/\partial T$  at each grid-point. Only the contribution of the vertical component of the derivative was taken into account in approximating  $\partial S/\partial T$ . This is analogous to the TH scheme which accounted for salinity changes arising from vertical advective processes only. In well-mixed regions there is little or no reason to expect temperature and salinity to be correlated. In these regions, which are characterized by small T or S gradients, the balance coefficient was set to zero, and the salinity analysis was dependent entirely on the unbalanced salinity component. The incremental formulation (4.30) suggests that the linearized scheme should be used in the inner loop only and that the original nonlinear TH scheme should be retained in the outer loop. While this may have had a beneficial impact on the analyses, it was decided instead to employ the approximation (4.33) for simplicity.

Figure 4.3 illustrates the impact of the multivariate T-S scheme on the mean salinity profile between 0 and 2000m in the TAO region  $(10^{\circ}\text{S}-10^{\circ}\text{N}; 160^{\circ}\text{E}-70^{\circ}\text{W})$  of the tropical Pacific. The results are from a multi-year 3D-Var experiment assimilating only temperature data. Compared to climatology (red curve), when no T-S constraint is applied, the water is too fresh above 400m and too salty below (green curve), and produces a notable reduction in the salinity maximum near 100m. The mean salinity profile in this experiment is clearly degraded compared to the experiment without data assimilation (black curve). With the T-S constraint, the mean salinity profile is improved, through a reduction of the artificial freshening and saltening above and below 400m (blue curve). The currents at the surface and below the core of the Equatorial Undercurrent were also improved, even though this field was not directly constrained by the multivariate analysis. **Ricci** et al. (2005) went on further to perform a detailed analysis of the heat and salt budgets in the model. This original aspect of their work provided useful physical insight of the processes at work in the model in response to the temperature data assimilation, with and without the T-S constraint.

### 4.2.3 Additional balance constraints

Weaver et al. (2005) extended the work of Ricci et al. (2005) to include balance relationships for the other ocean state variables,  $\eta$ , u and v. The balance was still conditioned by temperature in the sense that this was the primary variable used as the starting point to establish the balanced part of the other variables. Given the increment to the state vector  $\delta \mathbf{x}^{(k)} = (\delta T, \delta S, \delta \eta, \delta u, \delta v)^{\mathrm{T}}$ , where the components are understood to be column



Figure 4.3: The 1996 mean salinity profile between 0 and 2000m averaged over the TAO region: Levitus climatology (solid red curve), control (no data assimilation; dashed black curve), assimilation without T-S constraint (dashed-dotted green curve), and assimilation with T-S constraint (dashed-dotted blue curve). (From **Ricci** et al. (2005)).

vectors containing the value of the variable at each grid-point, then the control vector was taken to be  $\delta \mathbf{w}^{(k)} = (\delta T, \delta S_{\mathrm{U}}, \delta \eta_{\mathrm{U}}, \delta u_{\mathrm{U}}, \delta v_{\mathrm{U}})^{\mathrm{T}}$  where the subscript U means the unbalanced component of that variable. In general operator form, the sequence of linearized balance transformations comprising  $\mathbf{K}^{(k-1)}$  had the following structure,

$$\begin{split} \delta T^{(k)} &= \delta T^{(k)} \\ \delta S^{(k)} &= \mathbf{K}_{ST}^{(k-1)} \, \delta T^{(k)} \, + \, \delta S_{\mathrm{U}}^{(k)} \, = \, \delta S_{\mathrm{B}}^{(k)} \, + \, \delta S_{\mathrm{U}}^{(k)} \\ \delta \eta^{(k)} &= \mathbf{K}_{\eta\rho} \, \delta \rho^{(k)} \, + \, \delta \eta_{\mathrm{U}}^{(k)} \, = \, \delta \eta_{\mathrm{B}}^{(k)} \, + \, \delta \eta_{\mathrm{U}}^{(k)} \\ \delta u^{(k)} &= \mathbf{K}_{up} \, \delta p^{(k)} \, + \, \delta u_{\mathrm{U}}^{(k)} \, = \, \delta u_{\mathrm{B}}^{(k)} \, + \, \delta u_{\mathrm{U}}^{(k)} \\ \delta v^{(k)} &= \mathbf{K}_{vp} \, \delta p^{(k)} \, + \, \delta v_{\mathrm{U}}^{(k)} \, = \, \delta v_{\mathrm{B}}^{(k)} \, + \, \delta v_{\mathrm{U}}^{(k)} \end{split}$$

where

$$\delta \rho^{(k)} = \mathbf{K}_{\rho T}^{(k-1)} \delta T^{(k)} + \mathbf{K}_{\rho S}^{(k-1)} \delta S^{(k)}$$
  
$$\delta p^{(k)} = \mathbf{K}_{p\rho} \delta \rho^{(k)} + \mathbf{K}_{pn} \delta \eta^{(k)}$$
(4.34)

are diagnostic quantities corresponding to increments of density and pressure, respectively, and  $\mathbf{K}_{xy}$  represents the balance transformation from variable(s) y to x. The subscript B denotes the balanced part of that variable. The lower triangular form of  $\mathbf{K}^{(k-1)}$  means that the inverse operator  $(\mathbf{K}^{(k-1)})^{-1}$  can be computed trivially using forward elimination. This is convenient for estimating background-error covariance statistics of the control variables ( $\mathbf{w}$ ) from samples of background error with respect to the model variables ( $\mathbf{x}$ ) which is the information available in practice (e.g., from ensemble perturbations, see section 4.3).

The balanced component of salinity is defined from the T-S constraint described earlier. The resulting total salinity increment (balanced + unbalanced component) is then used with the temperature increment in a linearized version of the equation of state to compute a density increment (the first of the diagnostic equations in (4.34)). Note that the T-S and density balance operators include a superscript (k-1) to indicate that they are dependent on the reference state, which is not the case for the other balance relations. The density increment is then used to compute a baroclinic contribution to the SSH through the dynamic height relation, assuming a "level of no motion" at some reference depth. (An alternative balance relation that does not require this assumption is described in **Weaver** et al. (2005) but involves the solution of an elliptic equation and thus is more costly and more complicated to implement). The unbalanced component of SSH is attributed to a barotropic signal. The total SSH increment and the density increment are then used with the hydrostatic relation to compute a pressure increment (the second of the diagnostic equations in (4.34)). Finally, the pressure increment is used in the geostrophic relation to compute a balanced component for velocity. Near the equator the Coriolis parameter goes to zero and thus the standard f-plane geostrophic relation breaks down. In this region the meridional component of velocity is reduced to zero while a  $\beta$ -plane approximation is used to compute a geostrophically balanced zonal component of velocity. The unbalanced part of velocity is then attributed to the ageostrophic component.

The coupling of the balance operator and its transpose with the univariate error covariance matrix  $\mathbf{B}_{(\mathbf{w})}$  via Eq. (4.32) results in a complex multivariate error covariance matrix with respect to the model variables, as illustrated in the Appendix of Weaver *et al.* (2005). Although the explicit matrix form is never needed in practice, it is helpful for understanding how the different elements combine to determine the expression for the



Figure 4.4: Horizontal section of the SSH analysis increments generated by the 4D-Var assimilation of a single temperature observation located on the equator in the central Pacific at 100m depth and at day 10 into an assimilation window. The increments are displayed on day 10 for a 4D-Var experiment (a) without and (b) with the balance operator activated. The contour interval is 0.02m. (From Weaver *et al.* (2005)).

analysis increment in the presence of a single observation-type. Several illustrations were given in Weaver et al. (2005). While the balance operator is clearly fundamental in establishing a physically sensible (multivariate) analysis in 3D-Var, it also plays an important role in 4D-Var. This is illustrated in Fig. 4.4 which shows the SSH increments produced from two 4D-Var single temperature observation experiments performed without and with the balance operator activated. The observation is located in the thermocline on the equator and at the end of a 10-day assimilation window. The SSH increments shown are those produced at the observation time (day 10) using the tangent-linear model to propagate forward the analysis increment at initial time. The SSH increment produced without the balance operator has a localized structure similar to that obtained by a multivariate 3D-Var, whereas the increment produced with the balance operator results in the temperature observation projecting much more effectively onto large-scale equatorial waves-modes.

The balance operator can be considered effective if the variance of background error of the balanced variables explains a substantial part of the variance of background error of the full variables. If this is not the case then the balance operator would provide little useful information for the analysis. Using simulation errors produced with a NMC-type method as a proxy for actual background errors, **Weaver** et al. (2005) provide evidence that their proposed balance operator explains a significant percentage of background-error variance. Further evidence was given by **Daget** (2008) who studied the validity of the balance operator using simulation errors computed from the 9-member ensemble 45-year reanalysis for the ENSEMBLES project. His work also highlighted, however, certain limitations in the T-S balance and the equatorial velocity balance where a non-negligible cross-covariance was obtained between the balanced and unbalanced components of the error, contradicting the underlying hypothesis that these components are uncorrelated.

### 4.3 Covariance estimation using an ensemble

The general purpose of a covariance model is to incorporate in it as much prior knowledge as possible about the statistical structures of the errors in order to limit the number of covariance parameters that need to be estimated from actual statistics. The balance operator and the diffusion-based correlation operator described in this chapter are the core components of the covariance model developed for NEMOVAR. In the remainder of this chapter, the problem of estimating parameters of the covariance model is discussed within the context of ensemble data assimilation.

The next subsection deals with the variance estimation problem as studied by **Daget (2008)** and **Daget et al. (2009)** using an ensemble technique applied to a globalocean 3D-Var system. The final subsection discusses some recent work by **Weaver and Mirouze (2013)** to apply an ensemble technique to the estimation of parameters that control the length-scales and anisotropic response of a correlation model based on a diffusion operator.

### 4.3.1 Ensemble estimation of background-error variances

An important feature of an ensemble data assimilation system is its capacity to provide flow-dependent information on analysis and background error. This information can be exploited in a cycled assimilation system to improve the estimate of the backgrounderror covariance matrix on each cycle. A theoretical justification of the method within the context of variational assimilation is provided by Berre *et al.* (2006) and **Daget** *et* al. (2009). In particular, they show how perturbing the input parameters of a cycled analysis/forecast system leads to linearized evolution equations for the analysis and forecast state perturbations which are identical to those for the true errors. Furthermore, assuming that the perturbations to the input parameters are random samples drawn from the probability distribution of the true errors, then the evolved analysis and forecast perturbations from the cycled ensemble will also be random samples from the distribution of the true errors. The covariance matrices estimated from a sample of perturbed-minusunperturbed analysis and forecast differences then provide accurate estimates of the true analysis- and forecast-error covariance matrices. In practice, these covariance matrices will only be approximate due to the finite sample of the ensemble and due to inaccuracies in the specification of the error covariance matrix of the input parameters.

**Daget** et al. (2009) investigated the potential of using the spread of backgroundstates from a 9-member ensemble 3D-Var system to provide flow-dependent estimates of the background-error variances of temperature and (unbalanced) salinity in a low resolution global model (the ORCA2 configuration). (Only temperature and salinity data were assimilated so the only background-error variables to specify in the 3D-Var system were temperature and unbalanced salinity). The cycling and estimation procedure adopted in that study is summarized schematically in Fig. 4.5 (see caption for details). In order to reduce sampling error, a sliding window was used to include the ensemble of background states from the previous 9 cycles (90-days) in the computation of the variances for the current cycle. This effectively increased the ensemble size to 81, but at the expense of filtering out background-error variations on intraseasonal time-scales and strongly damping those on seasonal time-scales.



Figure 4.5: Schematic illustration of the ensemble 3D-Var system. The ensemble of analysis states  $\mathbf{x}_{l,c-1}^{a}(t_{N})$ ,  $l = 0, \ldots L - 1$ , at the end of cycle c - 1 are used to initialize the background trajectories of each ensemble member on the next cycle c. The background trajectory of each member l is produced by integrating the model with a perturbed set of forcing fields (wind-stress, heat flux, PmE),  $\mathbf{f}_{l,c,i} = \mathbf{f}_{c,i} + \tilde{\boldsymbol{\epsilon}}_{l,c,i}^{f}$ , from the initial condition  $\mathbf{x}_{l,c}^{b}(t_{0}) = \mathbf{x}_{l,c-1}^{a}(t_{N})$ . Each background trajectory is compared with a set of perturbed observations  $\mathbf{y}_{l,c,i}^{o} = \mathbf{y}_{c,i}^{o} + \tilde{\boldsymbol{\epsilon}}_{l,c,i}^{o}$  to produce an innovation vector for each member l. A 3D-Var (FGAT) analysis is then performed for each ensemble member using the appropriate innovation vector and a background-error variance matrix  $\mathbf{D}_{(\widehat{\mathbf{w}}),c}$  that has been estimated from the ensemble of background initial states  $\mathbf{x}_{l,c}^{b}(t_{0})$ . Ensemble member l = 0is unperturbed:  $\tilde{\boldsymbol{\epsilon}}_{0,c,i}^{f} = \mathbf{0}$  and  $\tilde{\boldsymbol{\epsilon}}_{0,c,i}^{o} = \mathbf{0}$ . The resulting analysis increment is then used to produce an analysis state trajectory using Incremental Analysis Updates. (From **Daget** et al. 2009).

The 3D-Var analyses produced with the ensemble variances showed improvements over those produced with an empirical parameterization of the variances. In particular, there was a reduction of error-growth between assimilation cyles and a better fit to nonassimilated observations such as equatorial current measurements in the Pacific, suggesting that the dynamical balance in the analyses was improved with the ensemble variances. However, statistical consistency diagnostics indicated that the ensemble variances were largely underestimated, especially in the upper ocean, which pointed to deficiencies in the ensemble-generation procedure and the need for variance inflation (which was not applied). Even more troubling was that both sets of 3D-Var analyses were in some regions and for some variables worse than those from a control analysis in which no data were assimilated. This pointed to more general problems in the data assimilation system. Model systematic error associated with poor resolution and inaccuracies in the forcing fluxes and vertical mixing parameterization scheme was viewed as an important factor contributing to the deficiencies in the data assimilation.

### 4.3.2 Ensemble estimation of background-error correlations

The study of **Daget** *et al.* (2009) examined the use of an ensemble to estimate only the diagonal elements of **B**. If  $N_{\rm g}$  denotes the total number of grid-points and  $N_{\rm v}$  the number of background-error variables then this requires the estimation of  $N = N_{\rm g} \times N_{\rm v}$ elements, which is much smaller than the  $(N^2 + N)/2$  independent elements required to estimate the full (symmetric) **B** matrix. As such, the variance estimation problem will be less exposed to problems of sampling error than the correlation estimation problem. (In practice, however, some degree of spatial or temporal filtering is needed even for the variances since the number of ensemble members  $(N_{\rm e})$  is typically much smaller than N).

The ensemble estimation of background-error correlations can also be reduced to a problem of order N by making certain simplifying assumptions about the form of the correlation function that we wish to estimate. This is discussed by **Weaver and Mirouze (2013)** within the context of a diffusion-based correlation model. In particular, assume that the correlation function is locally homogeneous and at least twice differentiable. This latter requirement is satisfied by both the Gaussian function and the Matérn functions, provided that for the latter M > 1 in  $\mathbb{R}$  and M > 2 in  $\mathbb{R}^2$  and  $\mathbb{R}^3$ . These are the functions that can be approximately modelled by applying an M-step implicitlyformulated diffusion operator.

For anisotropic and homogeneous versions of these functions, the distance between two points  $\mathbf{x}$  and  $\mathbf{x}'$  in  $\mathbb{R}^d$  is measured according to

$$r = \sqrt{\left(\mathbf{x} - \mathbf{x}'\right)^{\mathrm{T}} \boldsymbol{L}^{-1} \left(\mathbf{x} - \mathbf{x}'\right)}$$

where L is the aspect tensor of the correlation function. The aspect tensor defines the scale (diffusion) tensor in the diffusion equation. For the implicit diffusion kernels,

$$L = \frac{1}{2M - d - 2} H^{-1}$$
(4.35)

where H is the tensor of second-derivatives of the correlation function evaluated at zero separation. In the geostatistical literature H is known as the correlation Hessian tensor. Weaver and Mirouze (2013) called its inverse  $H^{-1}$ , which appears in Eq. (4.35), the Daley tensor in view of its analogy with the Daley length-scale in the isotropic case. In particular, they discuss how ensemble perturbations can be used to obtain a sample estimate of the elements of H by exploiting the diagnostic formulae of Belo Pereira and Berre (2006). In compact form, the expression reads (Michel 2013),

$$\boldsymbol{H}_{e} = \sum_{l=1}^{N_{e}} \nabla \widehat{\epsilon}_{l}^{b} \left( \nabla \widehat{\epsilon}_{l}^{b} \right)^{T}$$
(4.36)

where

$$\widehat{\epsilon}_l^{\rm b} = \frac{\epsilon_l^{\rm b}}{\sigma_{\rm e}\sqrt{N_{\rm e}-1}} \tag{4.37}$$

is the background perturbation  $\epsilon_l^{\rm b}$  normalized by  $\sqrt{N_{\rm e}-1}$  times the sample estimate of its standard deviation  $\sigma_{\rm e}$ . The estimate of  $H_{\rm e}$  can be done at each grid-point and the corresponding scale tensor  $L_{\rm e}$  estimated from Eq. (4.35). The resulting correlation estimates are both anisotropic and inhomogeneous. In 2D (3D) this requires the estimation of 3N (6N) independent elements, which is of the same order as the variance estimation problem. Equations (4.36) and (4.37) illustrate that the tensor- and variance-estimation problems are tightly connected, which suggests that these parameters should be defined consistently in the data assimilation system. This contrasts the approach of **Daget** *et al.* (2009) where only the variances were estimated from the ensemble while the lengthscales were parameterized empirically.

Rather than using the ensemble to calibrate the parameters (variances and aspect tensor) of a pre-defined **B** model, an alternative approach is to use it to provide a direct sample estimate of **B** and then to localize the sample estimate using a Schur (element-by-element) product in order to remove remote covariances associated with sampling noise. This is a common approach used in practical applications of the Ensemble Kalman filter (e.g., Houtekamer and Mitchell 2005) and Ensemble Variational assimilation (e.g., Bishop *et al.* 2011). Restricting the localization to the sample *correlation* matrix,  $\mathbf{C}_{e} = \mathbf{X}\mathbf{X}^{T}$  where  $\mathbf{X}$  is a  $N_{g} \times N_{e}$  matrix whose columns are the normalized perturbations  $\hat{\boldsymbol{\epsilon}}_{l}^{b} = \mathbf{D}_{e}^{-1/2}\boldsymbol{\epsilon}_{l}^{b}$ (i.e., the discrete representation of Eq. (4.37)) then the localized sample correlation matrix has the form

$$\mathbf{B} = \mathbf{D}_{\mathrm{e}}^{1/2} [\mathbf{C}_{\mathrm{e}} \circ \mathbf{C}_{\mathrm{loc}}] \mathbf{D}_{\mathrm{e}}^{1/2}$$
(4.38)

where  $\mathbf{C}_{e} \circ \mathbf{C}_{loc}$  denotes the Schur product of  $\mathbf{C}_{e}$  with a prescribed  $N_{g} \times N_{g}$  localized correlation matrix  $\mathbf{C}_{loc}$ . The length-scales in  $\mathbf{C}_{loc}$  determine the distance beyond which the correlations in  $\mathbf{C}_{e}$  should be significantly damped or explicitly set to zero.

An equivalent but more convenient form of Eq. (4.38) for variational assimilation is (e.g., see Buehner 2012)

$$\mathbf{B} = \mathbf{D}_{e}^{1/2} \left[ \sum_{p=1}^{N_{e}} \mathbf{D}_{\hat{\boldsymbol{\epsilon}}_{p}} \mathbf{C}_{\text{loc}} \mathbf{D}_{\hat{\boldsymbol{\epsilon}}_{p}} \right] \mathbf{D}_{e}^{1/2}$$
(4.39)

where  $\mathbf{D}_{\hat{\boldsymbol{\epsilon}}_p} = \operatorname{diag}(\hat{\boldsymbol{\epsilon}}_p)$ . In Eq. (4.39) the Schur product is now replaced by a standard matrix multiplication. When implemented in a **B**-preconditioned conjugate gradient (CG) algorithm (**Gürol et al. 2013**),  $\mathbf{C}_{\text{loc}}$  can be applied as an operator. The diffusion operator is an obvious candidate for defining  $\mathbf{C}_{\text{loc}}$ . Furthermore, the scale tensor estimated from  $\boldsymbol{H}_e$  can be used as the basis for selecting a spatially-dependent localization tensor  $\boldsymbol{L}_{\text{loc}}$ , thus generalizing the concept of a localization scale to account for directionality. A natural choice of the tensor would be  $\boldsymbol{L}_{\text{loc}} = \alpha \boldsymbol{L}_e$  for some constant  $\alpha > 1$ .

Figure 4.6 illustrates the effectiveness of the two methods for estimating correlations from a small number of ensemble perturbations. The analytical model used to generate the anisotropic and inhomogeneous "true" correlations, displayed in panel (a) at selected points, is described in Weaver and Mirouze (2013). The corresponding patterns obtained from a sample estimate of the correlation matrix ( $C_e$ ) with a 10-member ensemble (panel (b)) are heavily contaminated by sampling noise. The signal is clearly recognizable with 100 ensemble members (panel (c)) but sampling noise still leads to non-negligible spurious correlations at large separation distances. Panels (d) and (e) show that both the correlation model and explicit correlation localization are reasonably effective at reproducing the true correlations even with only 10 ensemble members. (See the figure caption for further details). The diffusion-based correlation model requires the estimation of order  $N_{\rm g}$  diffusion tensor elements for each of the analysis variables. Accounting for general non-diagonal diffusion tensors requires the solution of non-trivial large-scale elliptic equations for which computationally efficient methods are essential. Correlation localization requires  $N_{\rm e}$  applications of the diffusion operator on each CG iteration but can be performed using simple (e.g., isotropic) diffusion formulations while still allowing complex correlation information to be extracted from the ensemble. This may result in important computational savings compared to the direct diffusion modelling of anisotropic and inhomogeneous correlations. Localization also provides more flexibility with regard to the specification of multivariate covariances. A promising avenue for future development is to consider hybrid covariance formulations that linearly combine the two representations of **B**; i.e., a covariance model with ensemble-estimated parameters and an explicitly localized sample estimate of the covariance matrix.



Figure 4.6: (a) True correlations at selected points, computed from an anisotropic and inhomogeneous correlation matrix. (b) Raw correlations estimated from a 10-member ensemble drawn from the true distribution. (c) As panel (b) but with a 100-member ensemble. (d) Reconstructed correlations using a diffusion operator with the diffusion tensor estimated from the same 10-member ensemble. (e) Reconstructed correlations obtained by explicitly localizing the sample correlations with a diffusion operator whose tensor is defined as twice the estimated diffusion tensor. The latter approach is equivalent to an adaptive Schur-product localization. (Panels (a) and (c) are from Weaver and Mirouze (2012)); panel (d) is analogous to panel (g) in Weaver and Mirouze (2012) but without local spatial averaging.)

# Chapter 5

## Summary and outlook

This manuscript has provided a summary of research I have conducted in ocean data assimilation, focussing on my work in background-error covariance modelling. Since the beginning of my career, I have placed great importance on demonstrating the practical benefits of my research, including the research of PhD students and post-doctoral scientists that I have supervised. To this end, much of my work has been devoted to the development of a data assimilation system for the community ocean model NEMO. Early in my career I developed a variational data assimilation system known as OPAVAR, which was tailored to NEMO's predecessor OPA. In collaboration with other groups, it was rewritten to fit within the framework of NEMO and extended to include new features and operational capabilities. The revised system is called NEMOVAR and is the basis of the current operational ocean data assimilation systems at ECMWF and the Met Office.

An overriding message in this manuscript is that the development of an effective data assimilation system requires a good understanding of the physics and characteristics of the underlying problem. Nowhere is this more true than in the specification of the background-error covariance matrix  $(\mathbf{B})$ . The importance of multivariate relationships in the covariance matrix was one aspect that was highlighted. Computational efficiency is another key aspect of the assimilation problem, especially in oceanographic applications where the dimension of the state vector is enormous and the cost of applying the different assimilation components is high. Spatial correlation operators used in background-error covariance models or for localizing ensemble-estimated background-error covariances are particularly demanding of computational resources, which explains why a great deal of research by the data assimilation community has been aimed at trying to reduce the cost of this operation. As discussed in this manuscript, diffusion operators can be used to define general and computationally efficient correlation operators, and are particularly well suited for variational assimilation with grid-point models in complex boundary domains such as the ocean. Much of my research has been dedicated to studying diffusion-based correlation operators both from a general theoretical perspective and with the specific goal of implementing them in NEMOVAR. The minimization algorithm is another component of the assimilation system where computational efficiency and design constraints are crucial. Although only briefly discussed in this manuscript, this is an area where I have collaborated with the Parallel Algorithms group at CERFACS to develop methods with important practical benefits for ocean data assimilation.

Improving the specification of  $\mathbf{B}$  remains a key challenge in ocean data assimilation

and will continue to be a focus of my research in the coming years. Ensemble methods provide the appropriate framework for obtaining flow-dependent estimates of background error and naturally link the problems of data assimilation and probabilistic forecasting. Ensemble methods are also well suited for massively parallel computations required by modern-day computers. To exploit ensemble methods effectively in ocean data assimilation will require a significant research effort. In this manuscript I have described some initial work I have conducted in this area in collaboration with colleagues and students. This work needs to be further developed to improve both the ensemble-perturbation strategy needed for sampling the major sources of uncertainty in the ocean model and data assimilation system, and the methods needed for synthesizing ensemble-covariance information in practical **B** formulations for variational assimilation, as discussed at the end of Chapter 4. These issues will be addressed through my involvement in current and future research projects (LEFE-MANU, RTRA-FILAOS/AVENUE, ERA-CLIM2) and through collaborative work with ECMWF, the UK Met Office and INRIA aimed at developing NEMOVAR.

### Acknowledgements

I would like to thank my "correspondant", Serge Gratton, the referees and the other examiners for accepting to be part of my HDR committee. I would also like to thank the many students and colleagues from CERFACS and other institutes with whom I have had the opportunity of collaborating and without whom this work would not have been possible.

## References

- Balmaseda, M. A., Vidard A., Anderson D. L. T., 2008: The ECMWF Ocean Analysis System: ORA-S3. *Mon. Weather Rev.*, **136**, 3018–3034.
- Balmaseda, M. A., Mogensen, K. and A. T. Weaver, 2013: Evaluation of the ECMWF Ocean Reanalysis ORAS4. Q. J. R. Meteorol. Soc., 139, 1132–1161.
- Bannister, R. N., 2008: A review of forecast error covariance statistics in atmospheric variational data assimilation. II: Modelling the forecast error covariance statistics. Q. J. R. Meteorol. Soc., 134, 1971–1996.
- Belo Pereira, M. and L. Berre, 2006: The use of an ensemble approach to study the background error covariances in a global NWP model. *Mon. Weather Rev.*, **134**, 2466–2489.
- Bengtsson, L., 1979: On the use of a time-sequence of surface pressures in fourdimensional data assimilation. *Tellus*, **32**, 189–196.
- Bennett, A. F., 2002: Inverse Modelling of the Ocean and Atmosphere. Cambridge University Press, U. K.
- Bennett, A. F., Chua, B. S. and L. M. Leslie, 1997: Generalized inversion of a global numerical weather prediction model. II: analysis and implementation. *Meteorol. Atmos. Phys.*, 62, 129–140.
- Berre, L., Ştefănescu, S. E. and M. Belo Pereira, 2006: The representation of the analysis effect in three error simulation techniques. *Tellus*, **58A**, 196–209.
- Bishop, C. H., Hodyss, D., Steinle, P., Simes, H. Clayton, A. M., Lorenc, A. and D. M. Barker, 2011: Efficient ensemble covariance localization in variational data assimilation. *Mon. Weather Rev.*, 139, 573–580.
- Buehner, 2012: Evaluation of spatial/spectral covariance localization approach for atmospheric data assimilation. *Mon. Weather Rev.*, **140**, 617–636.
- Busalacchi, A. and J. J. O'Brien, 1981: Interannual variability of the equatorial Pacific in the 1960s. J. Geophys. Res., 86, 10901–10907.
- Cane, M. A., 1979: The response of an equatorial ocean to simple wind stress patterns: I. Model formulation and analytic results. J. Mar. Res., 57, 233–252.

- Courtier, P. and O. Talagrand, 1987: Variational assimilaton of meteorological observations with the adjoint vorticity equation. II: Numerical results. *Q. J. R. Meteorol. Soc.*, **131**, 1329–1347.
- Courtier, P. and O. Talagrand, 1990: Variational assimilation of meteorological observations with the direct and adjoint shallow-water equations. *Tellus*, **42A**, 531– 549.
- Courtier, P., Thépaut, J.-N. and A. Hollingsworth, 1994: A strategy for operational implementation fo 4D-Var, using an incremental approach. Q. J. R. Meteorol. Soc., 120, 1367–1388.
- Daget, N., 2008: "Estimation densemble des paramètres des covariances derreur débauche dans un système dassimilation variationnelle de donnés ocèaniques". PhD thesis, Université Paul Sabatier, June 2008.
- Daget, N., Weaver, A. T. and M. A. Balmaseda, 2009: Ensemble estimation of background-error variances in a three-dimensional variational data assimilation system for the global ocean. *Q. J. R. Meteorol Soc.*, **135**, 1071–1094.
- Daley, R., 1991: Atmospheric Data Analysis, Cambridge University Press, U. K.
- Deckmyn, A. and L. Berre, 2005: A wavelet approach to representing background error covariances in a limited-area model. *Mon. Weather Rev.*, **133**, 1279–1294.
- Dee, D. P. and co-authors, 2011: The ERA-interim reanalysis: configuration and performance of the data assimilation system. *Q. J. R. Meteorol. Soc.*, **131**, 2961–3012.
- Derber, J. C. and A. Rosati, 1989: A global oceanic data assimilation system. *J. Phys. Oceanogr.*, **19**, 1333–1347.
- Derber, J. C. and F. Bouttier, 1999: A reformulation of the background error covariance in the ECMWF global data assimilation system. *Tellus*, **51A**, 195–221.
- Desroziers, G., Hello, G. and J.-N. Thépaut, 2003: A 4D-Var re-analysis of the FASTEX. Q. J. R. Meteorol. Soc., 129, 1301–1315.
- Dobricic, S. and Pinardi N., 2008: An oceanographic three-dimensional variational data assimilation scheme. *Ocean Modelling*, **22**, 89–105.
- Egbert, G., Bennett, A. F. and M. Foreman, 1994. Topex/Poseidon tides estimated using a global inverse model. J. Geophys. Res., 99, 24821–24852.
- Fisher M. 2003. Background error covariance modelling. *Proceedings of the ECMWF* Seminar on Recent Developments in Data Assimilation for Atmosphere and Ocean, ECMWF, pp. 35–63.
- Fu, L.-L. and A. Cazenave, 2001: *Satellite Altimetry and Earth Sciences*. International Geophysics Servies, Volume 69.

- Gaspari, G. and S. Cohn, 1999: Construction of correlation functions in two and three dimensions. Q. J. R. Meteorol. Soc., 125, 723–757.
- Gaspari, G., Cohn, S., Guo, J. and S. Pawson, 2006: Construction and application of correlation functions with variable length-fields. Q. J. R. Meteorol. Soc., 132, 1815–1838.
- Gauthier, P., Tanguay, M., Laroche, S. and S. Pellerin, 2007: Extension of 3DVAR to 4DVAR: Implementation of 4DVAR at the Meteorological Service of Canada. *Mon. Weather Rev.*, 135, 2339–2353.
- Gilbert, J.-C. and C. Lemaréchal, 1989: Some numerical experiments with variablestorage quasi-Newton algorithms. *Math. Program.*, **45**, 407–435.
- Gneiting, T., 2002: Compactly supported correlation functions. J. Multivariate Anal., 83, 493–508.
- Golub, G. H. and C. F. van Loan, 1996: *Matrix Computations*, 3rd edn, John Hopkins University Press, Baltimore.
- Gratton, S., Lawless, A. and N. Nichols, 2007: Approximate Gauss-Newton methods for nonlinear least squares problems. *SIAM J. Optimiz.*, **18**, 106–132.
- Gürol, S., Weaver, A. T., Moore, A. M., Piacentini, A., Arango, H. and S. Gratton, 2013: B-preconditioned minimization algorithms for variational data assimilation with the dual formulation. Q. J. R. Meteorol. Soc.. DOI:10.1002/qj.2150. In press.
- Guttorp, P. and Gneiting T. 2006: Miscellanea studies in the history of probability and statistics XLIX: on the Matérn correlation family. *Biometrika*, **93**, 989–995.
- Haidvogel, D. B. and A. Beckmann, 1999: Numerical ocean circulation modeling. Imperial College Press, London.
- Houtekamer, P. L. and Mitchell H. L., 2005: Ensemble Kalman filtering. Q. J. R. Meteorol. Soc., 131, 3269–3289.
- Kalnay, E., 2003: Atmospheric Modelling, Data Assimilation and Predictability, Cambridge University Press, U. K.
- Le Dimet, F. X. and O. Talagrand, 1986: Variational algorithms for analysis and assimilation of meteorological observations: Theoretical aspects. *Tellus*, **38A**, 97–110.
- Lewis, J. M. and J. C. Derber, 1985: The use of adjoint equations to solve a variational adjustment problem with convective constraints. *Tellus*, **37A**, 309–322.
- Lorenc, A. C., 1992: Iterative analysis using covariance functions and filters. Q. J. R. Meteorol. Soc., 118, 569–591.

- Madec, G., Delecluse, P., Imbard, M. and C. Lévy, 1998: 'OPA8.1 Ocean General Circulation Model reference manual'. Technical note No. 11, LODYC/IPSL, Université Pierre et Marie Curie, Paris.
- Madec, G., 2008: 'NEMO ocean engine'. Note du Pôle de modélisation No. 27, ISSN No. 1288–1619, IPSL, Paris.
- McCreary, J. P., 1983: A model of tropical ocean-atmosphere interaction. *Mon. Weather Rev.*, **111**, 370–387.
- McCreary, J. P. and D. L. T. Anderson, 1984: A simple model of El Niño and the Southern Oscillation. *Mon. Weather Rev.*, **112**, 934–946.
- Meurant, G., 2006: The Lanczos and Conjugate Gradient Algorithms: From Theory to Finite Precision Computations. SIAM, Philadelphia.
- Michel, Y., 2013: Estimating deformations of random processes for correlation modelling: methodology and the one-dimensional case. *Q. J. R. Meteorol. Soc.*, **139**, 771–783.
- Mirouze, I., 2010: "Régularisation de problèmes inverses à laide de léquation de diffusion généralisé". PhD thesis, Université Paul Sabatier, September 2010.
- Mirouze, I. and A. T. Weaver, 2010: Representation of correlation functions in variational assimilation using an implicit diffusion operator. *Q. J. R. Meteorol. Soc.*, **136**, 1421–1443.
- Mogensen, K. S., Balmaseda, M. A., Weaver, A. T., Martin, M. and A. Vidard, 2009: 'NEMOVAR: a variational data assimilation system from the NEMO ocean model'. In ECMWF Newsletter 120 - Summer 2009, ECMWF, Reading, U. K., pp 17-21.
- Mogensen, K., M. A. Balmaseda and A. T. Weaver, 2012: The NEMOVAR ocean data assimilation system as implemented in the ECMWF ocean analysis for System 4. ECMWF Tech. Memo., No. 668, 60 pp. Also registered as a CERFACS Technical Report No. TR-CMGC-12-30.
- Moore, A. M., Arango, H. G., Broquet, G., Powell, B. S., Weaver, A. T. and J. Zavala-Garay, 2011a: The Regional Ocean Modeling System (ROMS) 4-dimensional variational data assimilation systems. Part I System overview and formulation. *Prog. in Oceanogr.*, 91, 34–49.
- Moore, A. M., Arango, H. G., Broquet, G., Powell, B. S., Weaver, A. T. and J. Zavala-Garay, 2011b: The Regional Ocean Modeling System (ROMS) 4-dimensional variational data assimilation systems. Part II Performance and applications to the California Current system. *Prog. in Oceanogr.*, **91**, 50–73.
- Moore, A. M., Arango, H. G., Broquet, G., Powell, B. S., Weaver, A. T. and J. Zavala-Garay, 2011c: The Regional Ocean Modeling System (ROMS) 4-dimensional variational data assimilation systems. Part III Observation imppact and observation sensitivity in the California Current system. Prog. in Oceanogr., 91, 74–94.

- Nocedal, J., 1980. Updating quasi-Newton matrices with limited storage. *Math. Comput.*, **35**, 773–782.
- Nocdeal, J. and S. J. Wright, 1999. Numerical Optimization. Series in Operations Research. Springer, New York.
- Pannekoucke, O., Berre, L. and G. Desroziers, 2007: Filtering properties of wavelets for local background-error correlations. *Q. J. R. Meteorol. Soc.*, **133**, 363–379.
- Pannekoucke, O. and S. Massart, 2008: Estimation of the local diffusion tensor and normalization for heterogeneous correlation modelling using a diffusion equation. *Q. J. R. Meteorol. Soc.*, **134**, 1425–1438.
- Parrish, D. F. and J. C. Derber, 1992: The National Meteorological Center's Spectral Statistical-Interpolation Analysis System. *Mon. Weather Rev.*, **120**, 1747–1763.
- Purser, R. J., Wu, W. S., Parrish, D. F. and N. M. Roberts, 2003. Numerical aspects of the application of recursive filters to variational statistical analysis. Part I: spatially homogeneous and isotropic Gaussian covariances. *Mon. Weather Rev.*, 131, 1524–1535.
- Purser, R. J., Wu, W. S., Parrish, D. F. and N. M. Roberts, 2003: Numerical aspects of the application of recursive filters to variational statistical analysis. Part II: Spatialy inhomogeneous and anisotropic general covariances. *Mon. Weather Rev.*, 131, 1536–1548.
- Rabier, F., Järvinen, H., Klinker, E., Mahfouf, J.-F. and A. Simmons, 2000: The ECMWF operational implementation of four-dimensional variational assimilation. Part I: Experimental results with simplified physics. *Q. J. R. Meteorol. Soc.*, **124**, 1809–1829.
- Rawlins, F., Ballard, S. P., Bovis, K. J., Clayton, A. M., Li, D., Inverarity, W., Lorenc, A. C. and T. J. Payne, 2007: The Met Office global four-dimensional variational data assimilation scheme. Q. J. R. Meteorol. Soc., 133, 347–362.
- Reverdin, G., E. Frankignoul, E. Kestenare, and M. J. McPhaden, 1994: Seasonal variability in the surface currents of the equatorial Pacific. *J. Geophys. Res.*, **99**, 20323-20344.
- Ricci, S., 2004: "Assimilation variationnelle océanique: modèlisation multivariée de la matrice de covariance derreur débauche". PhD thesis, Université Paul Sabatier, March 2004.
- Ricci, S., Weaver, A. T., Vialard, J. and P. Rogel, 2005: Incorporating statedependent temperature-salinity constraints in the background-error covariance of variational ocean data assimilation. *Mon. Weather Rev.*, **133**, 317–338.
- Roullet, G. and G. Madec, 2000: Salt conservation, free surface and varying volume: a new formulation for ocean GCMs. *J. Geophys. Res.*, **105**, 23927–23942.

- Talagrand, O. and P. Courtier, 1987: Variational assimilation of meteorological observation with the adjoint vorticity equation, I: Theory. Q. J. R. Meteorol. Soc., 113, 1311–1328.
- Thacker, W. C. and R. B. Long, 1988: Fitting dynamics to data. J. Geophys. Res, 97, 7479–7491.
- Thépaut, J.-N. and P. Courtier, 1991: Four-dimensional variational data assimilation using the adjoint of a multilevel primitive equation model. Q. J. R. Meteorol. Soc., 117, 1225–1254.
- Troccoli, A. and K. Haines, 1999: Use of temperature-salinity relation in a data assimilation context. J. Atmos. Oceanic. Technol., 16, 2011–2025.
- Troccoli, A., Balmaseda, M.-A., Segschneider, J., Vialard, J., D. L. T. Anderson, Stockdale, T., Haines, K. and A. D. Fox, 2002: Salinity adjustments in the presence of temperature data assimilation. *Mon. Weather Rev.*, **130**, 89–102.
- Tshimanga, J., Gratton, S., Weaver, A. T. and A. Sartenaer, 2008: Limited-memory preconditioners, with application to incremental four-dimensional variational data assimilation. *Q. J. R. Meteorol. Soc.*, **134**, 753–771.
- Uppala, S. M. and co-authors, 2005: The ERA-40 re-analysis. *Q. J. R. Meteorol. Soc.*, **131**, 2961–3012.
- Vialard, J., Weaver, A. T., Anderson, D. L. T. and P. Delecluse, 2003: Threeand four-dimensional variational assimilation with a general circulation model of the tropical Pacific Ocean. Part II: Physical validation. *Mon. Weather Rev.*, **131**, 1379–1395.
- Yaremchuk, M. and M. Carrier, 2012: On the renormalization of the Covariance operators. *Mon. Weather Rev.*, **140**, 637–649.
- Weaver, A. T. and D. L. T. Anderson, 1997: Variational assimilation of altimeter data in a multi-layer model of the tropical Pacific Ocean. J. Phys. Oceanogr., 27, 664-682.
- Weaver, A. T. and P. Courtier, 2001: Correlation modelling on the sphere using a generalized diffusion equation. Q. J. R. Meteorol. Soc., 127, 1815–1846.
- Weaver, A. T. and I. Mirouze, 2013: On the diffusion equation and its application to isotropic and anisotropic correlation modelling in variational assimilation. *Q. J. R. Meteorol Soc.*. **139**, 242–260.
- Weaver, A. T., Vialard, J., Anderson, D. L. T. and P. Delecluse, 2002: Threeand four-dimensional variational assimilation with a general circulation model of the tropical Pacific Ocean. ECMWF Tech. Memo. No. 365. Available at http://www.ecmwf.int/publications/.

- Weaver, A. T., Vialard, J. and D. L. T. Anderson, 2003: Three- and four-dimensional variational assimilation with an ocean general circulation model of the tropical Pacific Ocean. Part 1: formulation, internal diagnostics and consistency checks. *Mon. Weather Rev.*, **131**, 1360–1378.
- Weaver, A. T. and S. Ricci, 2004: 'Constructing a background-error correlation model using generalized diffusion operators'. In ECMWF proceedings of the seminar on Recent Developments in Data Assimilation for Atmosphere and Ocean, ECMWF, Reading, U. K., pp 327–339.
- Weaver, A. T., Deltel, C., Machu, E., Ricci, S. and N. Daget, 2005: A multivariate balance operator for variational ocean data assimilation. *Q. J. R. Meteorol. Soc.*, **131**, 3605–3625.

## Appendix A

## Selected articles

 Correlation modelling on the sphere using a generalized diffusion equation, A. T. Weaver and P. Courtier, *Quarterly Journal of the Royal Meteorological Society*, Volume 127, Issue 575.

Copyright © 2001 Royal Meteorological Society, first published by the Royal Meteorological Society.

2. Three- and four-dimensional variational assimilation with an ocean general circulation model of the tropical Pacific Ocean. Part 1: formulation, internal diagnostics and consistency checks, A. T. Weaver, J. Vialard and D. L. T. Anderson, *Monthly Weather Review*, Volume 131, Issue 7.

Copyright © 2003 American Meteorological Society, first published by the American Meteorological Society.

 A multivariate balance operator for variational ocean data assimilation, A. T. Weaver, C. Deltel, E. Machu, S. Ricci and N. Daget, *Quarterly Journal of the Royal Meteorological Society*, Volume 131, Issue 613.

Copyright © 2005 Royal Meteorological Society, first published by the Royal Meteorological Society.

4. Ensemble estimation of background-error variances in a three-dimensional variational data assimilation system for the global ocean, N. Daget, A. T. Weaver and M. A. Balmaseda, *Quarterly Journal of the Royal Meteorological Society*, Volume 135, Issue 641.

Copyright © 2009 Royal Meteorological Society, first published by John Wiley & Sons, Ltd.

5. On the diffusion equation and its application to isotropic and anisotropic correlation modelling in variational assimilation. A. T. Weaver and I. Mirouze, *Quarterly Journal of the Royal Meteorological Society*, Volume 139, Issue 670.

Copyright © 2013 Royal Meteorological Society, first published by John Wiley & Sons, Ltd.

### Correlation modelling on the sphere using a generalized diffusion equation

### By ANTHONY WEAVER<sup>1,2\*</sup> and PHILIPPE COURTIER<sup>1</sup>

<sup>1</sup>Laboratoire d'Océanographie Dynamique et de Climatologie, France <sup>2</sup>Centre Européen de Recherche et de Formation Avancée en Calcul Scientifique, France

(Received 27 June 2000; revised 8 February 2001)

#### SUMMARY

An important element of a data assimilation system is the statistical model used for representing the correlations of background error. This paper describes a practical algorithm that can be used to model a large class of two- and three-dimensional, univariate correlation functions on the sphere. Application of the algorithm involves a numerical integration of a generalized diffusion-type equation (GDE). The GDE is formed by replacing the Laplacian operator in the classical diffusion equation by a polynomial in the Laplacian. The integral solution of the GDE defines, after appropriate normalization, a correlation operator on the sphere. The kernel of the correlation operator is an isotropic correlation function. The free parameters controlling the shape and lengthscale of the correlation function are the products  $\kappa_p T$ ,  $p = 1, 2, \ldots$ , where  $(-1)^p \kappa_p$  is a weighting ('diffusion') coefficient ( $\kappa_p > 0$ ) attached to the Laplacian with exponent p, and T is the total integration 'time'. For the classical diffusion equation (a special case of the GDE with  $\kappa_p = 0$  for all p > 1) the correlation function is shown to be well approximated by a Gaussian with length-scale equal to  $(2\kappa_1 T)^{1/2}$ .

The Laplacian-based correlation model is particularly well suited for ocean models as it can be easily generalized to account for complex boundaries imposed by coastlines. Furthermore, a one-dimensional analogue of the GDE can be used to model a family of vertical correlation functions, which in combination with the twodimensional GDE forms the basis of a three-dimensional, (generally) non-separable correlation model. Generalizations to account for anisotropic correlations are also possible by stretching and/or rotating the computational coordinates via a 'diffusion' tensor. Examples are presented from a variational assimilation system currently under development for the OPA ocean general-circulation model of the Laboratoire d'Océanographie Dynamique et de Climatologie.

KEYWORDS: Correlation functions Diffusion equation Objective analysis Variational assimilation

### 1. INTRODUCTION

A central task in the development of a statistical (e.g. variational) data assimilation system is the estimation and representation of the background-state error covariances. In practice, these covariances must be modelled as approximations to the true covariances of background error. In general, this is done by parametrizing the covariances using smoothing functions or filters, and balance relationships which employ such simplifying assumptions as isotropy, homogeneity, geostrophy etc. (Daley 1991). The number of adjustable parameters in the covariance model, such as the standard deviations and correlation length-scales of the background error, is usually quite small compared to the number of individual covariances that actually need to be specified (typically of the order of  $10^{11}$  or greater for global applications with general-circulation models of the atmosphere or ocean). These parameters can be determined by fitting the covariance models to available statistical information on the background errors, obtained, for example, from statistics of the observed-minus-background field (Hollingsworth and Lönnberg 1986; Lönnberg and Hollingsworth 1986) or of some other proxy for background error.

Ideally, the covariance models should be made flexible enough to capture the main characteristics of the available estimates of the background-error covariances. This paper focusses on a particular aspect of covariance modelling, namely that of deriving efficient and flexible algorithms for representing two-dimensional (2D) and three-dimensional (3D), univariate correlations on the sphere, with a special emphasis

<sup>\*</sup> Corresponding author: CERFACS, 42 avenue Gaspard Coriolis, 31057 Toulouse Cedex, France.

e-mail: weaver@cerfacs.fr

<sup>©</sup> Royal Meteorological Society, 2001.

on the applicability of these algorithms for variational data assimilation systems with ocean models.

The background-error covariance models employed in variational assimilation systems with atmospheric models commonly rely on a representation of correlation functions in terms of a spherical-harmonic expansion. Such a formulation is used operationally at the European Centre for Medium-Range Weather Forecasts (ECMWF) (Courtier et al. 1998; Rabier et al. 2000), the National Centers for Environmental Prediction (NCEP) (Parrish and Derber 1992), the Canadian Meteorological Center (Gauthier et al. 1999) and the Met Office (Lorenc et al. 2000). The free parameters in the correlation model are the spectral variances of the correlation function, which are generally specified from statistics of forecast differences valid at the same time (Rabier et al. 1998). Spectrally-based correlation models, however, are not very practical for ocean models since the lateral boundary conditions imposed by the coastlines are difficult to handle within a global spectral basis function expansion. Correlation models formulated directly in physical (grid-point) space are generally easier to adapt to bounded domains and hence better suited for the ocean assimilation problem. For example, Lorenc (1992, 1997) and Parrish et al. (1997) have developed correlation models based on recursive grid-point filters, one variant of which is applied routinely in the Met Office real-time global ocean forecasting system (Bell et al. 2000). The recursive filter is very efficient and can accommodate geographical variations in correlation length-scale, but has limited flexibility in the shape of the correlation function and is difficult to make isotropic (Lorenc 1997).

Derber and Rosati (1989) proposed the use of an iterative Laplacian grid-point filter (smoother) in order to build a 2D isotropic correlation function that approximates a Gaussian. The Derber and Rosati scheme is used operationally at NCEP for producing ocean initial conditions for seasonal climate forecasting (Behringer *et al.* 1998). A close variant of the algorithm has been described in some detail by Egbert *et al.* (1994) and Bennett *et al.* (1997), who interpret the Laplacian filter in terms of a time-step integration of a diffusion equation. The key idea is that the integral solution of the diffusion equation defines a covariance operator. In this paper, we build on this idea to construct 2D and 3D univariate correlation models that are both numerically efficient and sufficiently general to support a number of desirable features: correlation functions with different shapes (not just Gaussian), geographically variable length-scales, horizontal/vertical non-separability, and 3D anisotropy.

The paper is organized as follows. Section 2 presents the variational assimilation problem. A general, multivariate formulation of the background-error covariance matrix is also described in order to isolate clearly the univariate, correlation component, the specification of which is the subject of this paper. In section 3, the main features of the 2D 'Gaussian' correlation model of Derber and Rosati are described. Particular attention is given to the interpretation of the model in terms of the spherical harmonics (ignoring boundaries). A straightforward extension of this correlation model is then proposed as a means of generating a general class of 2D isotropic correlation functions on the sphere. A one-dimensional (1D) version of the model is also introduced for representing a class of vertical correlation functions. Several issues related to the practical implementation of the algorithm are discussed in section 4. A 3D correlation model is also developed. In section 5, the correlation structures are illustrated in several examples from a variational assimilation system (Weaver and Vialard 2000) for the OPA Ocean General Circulation Model (OGCM) of the Laboratoire d'Océanographie Dynamique et de Climatologie (LODYC) (Madec et al. 1999). In particular, the impact of ocean boundaries and various extensions to account for 2D and 3D anisotropy are discussed and illustrated. Concluding remarks are made in section 6. An appendix provides some mathematical details and a numerical illustration of a matching procedure that can be used to relate the exact solution of the diffusion equation on the sphere to a Gaussian covariance operator.

The notation used in the paper closely follows that of Ide et al. (1997).

### 2. FORMULATION OF THE PROBLEM

### (a) Variational assimilation and its dual formulation

Following Courtier *et al.* (1998), the variational formulation of three-dimensional assimilation (3D-Var) is introduced. In its incremental formulation (Courtier *et al.* 1994), 3D-Var attempts to minimize the following cost function,

$$J(\delta \mathbf{x}) = \underbrace{\frac{1}{2} \delta \mathbf{x}^{\mathrm{T}} \mathbf{B}^{-1} \delta \mathbf{x}}_{J_{\mathrm{b}}} + \underbrace{\frac{1}{2} (\mathbf{H} \delta \mathbf{x} - \mathbf{d})^{\mathrm{T}} \mathbf{R}^{-1} (\mathbf{H} \delta \mathbf{x} - \mathbf{d})}_{J_{\mathrm{o}}}, \tag{1}$$

where  $\delta x$  defines an increment to the background state vector  $x^b$ . The observation vector,  $y^o$ , is contained within the innovation vector

$$\mathbf{d} = \mathbf{y}^{\mathbf{o}} - \mathbf{H}(\mathbf{x}^{\mathbf{b}}),\tag{2}$$

where H is the observation operator. The matrix H in (1) is defined as a linearization of H about  $x^b$ . The matrices B and R contain the assumed covariances of background and observation error, respectively. In practice, these matrices are defined as approximations to the true error covariance matrices

$$\mathbf{B} \equiv \mathbf{E}[(\boldsymbol{\epsilon}^{\mathrm{b}} - \mathbf{E}[\boldsymbol{\epsilon}^{\mathrm{b}}])(\boldsymbol{\epsilon}^{\mathrm{b}} - \mathbf{E}[\boldsymbol{\epsilon}^{\mathrm{b}}])^{\mathrm{T}}], \qquad (3)$$

$$\mathbf{R} \equiv \mathbf{E}[(\boldsymbol{\epsilon}^{0} - \mathbf{E}[\boldsymbol{\epsilon}^{0}])(\boldsymbol{\epsilon}^{0} - \mathbf{E}[\boldsymbol{\epsilon}^{0}])^{\mathrm{T}}], \qquad (4)$$

where  $\epsilon^{b} = \mathbf{x}^{b} - \mathbf{x}^{t}$  and  $\epsilon^{o} = \mathbf{y}^{o} - H(\mathbf{x}^{t})$ ,  $\mathbf{x}^{t}$  denoting the 'true' state vector and  $E[\cdot]$  mathematical expectation.

The observation term  $(J_0)$  measures the fit between the model increment in observation space and the innovation vector. The background term  $(J_b)$  penalizes the size of the increment vector (i.e. measures the fit to the background state). At the minimum, the resulting analysis increment  $\delta x^a$  is added to  $x^b$  in order to provide the analysis  $x^a$ :

$$\mathbf{x}^{\mathbf{a}} = \mathbf{x}^{\mathbf{b}} + \delta \mathbf{x}^{\mathbf{a}}.$$
 (5)

Following Courtier (1997) and Derber and Bouttier (1999), the cost function (1) can be rewritten in terms of a new variable  $\mathbf{v}$ , defined by

$$\mathbf{v} = \mathbf{B}^{-1/2} \delta \mathbf{x},\tag{6}$$

where  $\mathbf{B}^{1/2}$  is taken to be any 'square-root' matrix such that  $\mathbf{B} = \mathbf{B}^{1/2}\mathbf{B}^{T/2}$ . (The superscript T/2 will be used throughout this paper in order to denote transpose of a square-root factor.) This leads to the equivalent analysis problem of minimizing

$$J(\mathbf{v}) = \frac{1}{2}\mathbf{v}^{\mathrm{T}}\mathbf{v} + \frac{1}{2}(\mathbf{H}\delta\mathbf{x} - \mathbf{d})^{\mathrm{T}}\mathbf{R}^{-1}(\mathbf{H}\delta\mathbf{x} - \mathbf{d}),$$
(7)

with

$$\delta \mathbf{x} = \mathbf{B}^{1/2} \mathbf{v}.$$
 (8)

A gradient descent method is used to iterate to the minumum of J. On each iteration, the gradient of J with respect to v is required, which from (7) and (8) is given by

$$\nabla_{\mathbf{v}}J = \nabla_{\mathbf{v}}J_{\mathbf{b}} + \nabla_{\mathbf{v}}J_{\mathbf{o}} = \mathbf{v} + \mathbf{B}^{T/2}\nabla_{\delta\mathbf{x}}J_{\mathbf{o}},\tag{9}$$

where

$$\nabla_{\delta \mathbf{x}} J_{\mathbf{o}} = \mathbf{H}^{\mathrm{T}} \mathbf{R}^{-1} (\mathbf{H} \delta \mathbf{x} - \mathbf{d}).$$
(10)

At the start of minimization  $\mathbf{v} = \delta \mathbf{x} = 0$ , providing the initial guess is taken to be  $\mathbf{x} = \mathbf{x}^{b}$ . As a result, only the transformation (8) from  $\mathbf{v}$  to  $\delta \mathbf{x}$  (involving  $\mathbf{B}^{1/2}$ ), and the associated adjoint transformation in (9) from  $\nabla_{\delta \mathbf{x}} J_{o}$  to  $\nabla_{\mathbf{v}} J_{o}$  (involving  $\mathbf{B}^{T/2}$ ) are actually required on each iteration. This formulation thus allows us to circumvent the explicit specification of the inverse matrix  $\mathbf{B}^{-1}$ , while at the same time providing a generally efficient preconditioner for the minimization (Lorenc 1988).

The three-dimensional Physical-space Statistical Analysis System (3D-PSAS) (Cohn *et al.* 1997) provides an alternative algorithm for solving the incremental, variational analysis problem (Courtier 1997). 3D-PSAS involves the iterative minimization of a cost function of the form

$$F(\mathbf{w}) = \frac{1}{2}\mathbf{w}^{\mathrm{T}}(\mathbf{R} + \mathbf{H}\mathbf{B}\mathbf{H}^{\mathrm{T}})\mathbf{w} - \mathbf{w}^{\mathrm{T}}\mathbf{d},$$
(11)

with

$$\delta \mathbf{x} = \mathbf{B} \mathbf{H}^{\mathrm{T}} \mathbf{w},\tag{12}$$

**w** being an increment defined in the dual of observation space. As with 3D-Var, 3D-PSAS may be implemented using an algorithm for applying  $\mathbf{B}^{1/2}$  and  $\mathbf{B}^{T/2}$  to a given vector; the inverse operators are not required providing  $\mathbf{w} = 0$  ( $\mathbf{x} = \mathbf{x}^{b}$ ) is the first guess for minimization. Strictly speaking, only **B** is required for 3D-PSAS. Nevertheless, the definition of **B** through a square-root factorization is a practical way of ensuring both its symmetry and positive definiteness (Parrish and Derber 1992; Gaspari and Cohn 1999).

It is worth remarking that, although we focus the discussion in this paper on **B** and the associated 3D variational analysis problem, the basic techniques we develop will be equally applicable for representing model and/or boundary-forcing error covariances in weak-constraint four-dimensional variational assimilation (4D-Var) or its dual formulation (4D-PSAS) (Egbert *et al.* 1994; Courtier 1997).

### (b) General formulation of the background-error covariance matrix

Because of its size, **B** can be neither estimated completely nor stored explicitly. (Recall that **B** is an  $N \times N$  matrix where N, the dimension of the model state vector, is typically greater than 10<sup>6</sup>.) We are therefore forced to model **B** as an operator. Specifically, we wish to define **B** as a sequence of operators, from which a factorization of the form  $\mathbf{B}^{1/2}\mathbf{B}^{T/2}$  can be easily deduced. For convenience we continue to use matrix notation, although it is important to bear in mind that the matrices are never explicitly computed in practice. By definition, **B** is a symmetric and positive definite matrix. The diagonal and off-diagonal elements of **B** correspond to the background-state error variances and covariances respectively. The covariances can be split into two contributions: block-diagonal elements that represent the auto-covariances between grid points corresponding to a particular model variable (i.e. the univariate component of **B**), and the remaining elements that represent the cross-covariances between grid points corresponding to different model variables (i.e. the multivariate component of **B**).

Recently, Derber and Bouttier (1999) have proposed a compact and powerful formulation of the background-error covariance matrix which allows separating of the univariate and multivariate components of **B** into distinct operators. It is useful to outline their formulation in order to put the present work in the context of a fully multivariate covariance operator. In their formulation, the model variables are first partitioned into balanced and unbalanced components, except for one variable which is taken in totality (i.e. for this variable there is no distinction between balanced and unbalanced). This (total) variable is then used to establish linear balances with the remaining variables. This procedure can be represented through a linear balance operator,  $\mathbf{K}'_{b}$ , acting on the unbalanced variables. As the model increment is defined as the sum of the balanced and unbalanced components, the full balance operator is  $\mathbf{K}_{b} = \mathbf{K}'_{b} + \mathbf{I}$  where **I** is the identity matrix. In terms of the balance operator, **B** takes the form

$$\mathbf{B} = \mathbf{K}_{\mathbf{b}} \mathbf{B}_{\mathbf{u}} \mathbf{K}_{\mathbf{b}}^{\mathrm{T}},\tag{13}$$

where  $\mathbf{B}_{u}$  is the error covariance matrix for the unbalanced variables, which is taken to have a block-diagonal structure (i.e. the cross-covariances between the unbalanced variables are assumed to be negligible).

By definition,  $\mathbf{B}_{u}$  can be factored as

$$\mathbf{B}_{\mathbf{u}} = \mathbf{\Sigma} \mathbf{C} \mathbf{\Sigma},\tag{14}$$

where  $\Sigma$  is a diagonal matrix of background-error standard deviations and C a symmetric matrix of background-error correlations for the unbalanced variables. Thus, from (13) and (14), we can represent the transformation in (8) by the sequence of operators

- - - -

$$\delta \mathbf{x} = \mathbf{K}_{\mathbf{b}} \mathbf{\Sigma} \mathbf{C}^{1/2} \mathbf{v},\tag{15}$$

where  $\mathbf{C}^{1/2}$  is defined such that  $\mathbf{C} = \mathbf{C}^{1/2}\mathbf{C}^{T/2}$ . As **C** is block-diagonal, the operator  $\mathbf{C}^{1/2}\mathbf{v}$  in (15) can in turn be split into individual operators,  $\mathbf{C}_{\alpha}^{1/2}\mathbf{v}_{\alpha}$ , that act independently on the different variable components,  $\mathbf{v}_{\alpha}$ , of **v**.

The purpose of this paper is to focus on the design of univariate correlation operators  $(\mathbf{C}_{\alpha} = \mathbf{C}_{\alpha}^{1/2} \mathbf{C}_{\alpha}^{T/2})$  that are both efficient for large state-vector assimilation problems and well adapted to oceanographic applications in complex-boundary domains. Deriving an appropriate multivariate balance operator ( $\mathbf{K}_{b}$ ) and estimating the free statistical parameters in the covariance model (e.g.  $\Sigma$ ) are equally important issues but are beyond the scope of this paper.

### 3. CORRELATION MODELLING ON THE SPHERE

### (a) Modelling a Gaussian correlation function using the diffusion equation

The following 1D problem provides a simple framework for interpreting the basic procedure for constructing correlation operators on the sphere, which we develop subsequently in sections 3(b)-(c). The starting point is the 1D diffusion equation

$$\frac{\partial \eta}{\partial t} - \kappa \frac{\partial^2 \eta}{\partial z^2} = 0, \tag{16}$$

where  $\eta$  is an arbitrary scalar field (e.g. temperature) and  $\kappa$  is a constant diffusion coefficient. We consider solutions to (16) on the real line  $\mathbb{R}$  such that  $\eta(z, t)$  vanishes as  $z \to \pm \infty$ . It is convenient to represent the general solution symbolically by an integral

operator  $\mathcal{L}$ :

$$\eta(z,0) \xrightarrow{\pounds} \eta(z,T),$$
 (17)

where  $\eta(z, T)$  is the result of integrating (16) over a time interval  $0 \le t \le T$  with  $\eta(z, 0)$  as initial condition.

A simple interpretation of the above operator can be given by considering the explicit solution of (16). The Fourier transform of (16) yields

$$\frac{\partial \widehat{\eta}(\widehat{z},t)}{\partial t} = -\kappa \widehat{z}^2 \,\widehat{\eta}(\widehat{z},t),\tag{18}$$

where  $\hat{\eta}(\hat{z}, t)$  denotes the Fourier transform of  $\eta(z, t)$  and  $\hat{z}$  is the variable in Fourier space. Equation (18) is readily integrated in time to give

$$\widehat{\eta}(\widehat{z}, T) = \widehat{\eta}(\widehat{z}, 0) \exp(-\kappa T \widehat{z}^2).$$
(19)

In (19) we recognize the product of  $\hat{\eta}(\hat{z}, 0)$  with a Gaussian function. As the inverse Fourier transform of a Gaussian is also a Gaussian,  $\eta(z, T)$  is the result of the convolution of  $\eta(z, 0)$  with a Gaussian function:

$$\eta(z, T) = \frac{1}{(4\pi\kappa T)^{1/2}} \int_{z'} \exp\{-(z - z')^2 / 4\kappa T\} \eta(z', 0) \, \mathrm{d}z',$$
(20)

where the product  $2\kappa T$  can be interpreted as the square of the length-scale,  $L^2$ , of the Gaussian function.

The Gaussian in (20) is homogeneous and isotropic as it depends only on Euclidean distance r = |z - z'|. It is well known that such a function defines a valid (positive definite) covariance function on  $\mathbb{R}$ . This follows from the fact that its Fourier transform (the exponential in (19)) is everywhere positive (e.g. see Gaspari and Cohn 1999). The integral equation (20) thus defines a covariance operator, which provides an interpretation of the operator  $\mathcal{L}$ .

By definition, a correlation function has unit variance: i.e. its value at the origin, r = 0, equals one. The Gaussian covariance operator  $\mathcal{L}$  can be easily transformed into a Gaussian correlation operator  $\mathcal{C}_{\eta}$  by post-multiplying  $\eta(z, T)$  by the constant factor  $(4\pi\kappa T)^{1/2}$ , which can be represented symbolically by

$$\eta(z,0) \xrightarrow{\mathfrak{C}_{\eta}} (4\pi\kappa T)^{1/2} \eta(z,T).$$
(21)

On a discrete grid, the action of  $C_{\eta}$  may be effected in one of two ways: either by evaluating the convolution integral in (20) directly using numerical quadrature, or by the generally more efficient technique of iterating a discretized (e.g. finite difference) version of the differential equation (16) and normalizing the result as in (21). The latter is the essence of the Derber and Rosati scheme (see also Egbert *et al.* (1994)). In the next section, we present the theoretical basis of the scheme for application on the sphere.

### (b) Interpretation on the sphere

In 3D Euclidean space ( $\mathbb{R}^3$ ), the solution of the diffusion equation leads to a convolution integral with a Gaussian function that depends only on Euclidean distance  $r = \{(x - x')^2 + (y - y')^2 + (z - z')^2\}^{1/2}$  between points (x, y, z) and (x', y', z') in  $\mathbb{R}^3$ . This is an obvious extension of the 1D example in the previous section. Any

1820

homogeneous and isotropic correlation function in  $\mathbb{R}^3$ , such as the 3D Gaussian, can be readily transformed into a valid isotropic correlation function on the spherical space  $S^2$  if Euclidean distance is substituted by *chordal* distance

$$r = 2a\sin(\theta/2) = 2a(1 - \cos\theta)^{1/2},$$
(22)

where  $\theta$ ,  $0 \le \theta \le \pi$ , is the angular separation (great circle distance) between points on the sphere of radius *a* (Weber and Talkner 1993; Gaspari and Cohn 1999; Gneiting 1999)<sup>†</sup>. One might expect, therefore, that the solution of the diffusion equation on the sphere leads to an integral operator whose kernel is the Gaussian function  $e^{-r^2/2L^2}$ where *r* is given by (22) and  $L^2 = 2\kappa T$  as in the 1D example. This turns out not to be the case as shown below; the solution has a covariance interpretation but the actual covariance function is not Gaussian, even though numerically it can be very closely matched to a Gaussian as discussed by Hartman and Watson (1974). This feature is investigated in more detail in the appendix.

Consider the diffusion equation

$$\frac{\partial \eta}{\partial t} - \kappa \nabla^2 \eta = 0, \tag{23}$$

where  $\nabla^2$  is the Laplacian operator defined on the sphere  $\Sigma$ . When the domain over  $\Sigma$  contains boundaries (e.g. coastlines in the case of an ocean model), we must supplement (23) with boundary conditions. For the moment we consider  $\Sigma$  to be free of boundaries. Boundary related issues will be addressed in section 5(b).

A scalar field  $\eta(\lambda, \phi, t)$  may be expanded as

$$\eta(\lambda,\phi,t) = \sum_{n=0}^{\infty} \sum_{m=-n}^{n} \eta_n^m(t) Y_n^m(\lambda,\phi), \qquad (24)$$

where  $\lambda$  is longitude  $(0 \le \lambda \le 2\pi)$  and  $\phi$  is latitude  $(-\pi/2 \le \phi \le \pi/2)$ , *m* is the zonal wave number, *n* is the total wave number,  $Y_n^m(\lambda, \phi)$  are the spherical harmonics, and  $\eta_n^m(t)$  are spectral expansion coefficients. The  $Y_n^m(\lambda, \phi)$  are orthogonal and normalized following the usual convention in meteorology:

$$\frac{1}{4\pi a^2} \int_{\Sigma} Y_n^m(\lambda, \phi) Y_{n'}^{m'^*}(\lambda, \phi) \, \mathrm{d}\Sigma = \delta_{nn'} \delta_{mm'}, \tag{25}$$

where \* denotes complex conjugate,  $d\Sigma = a^2 \cos \phi \, d\lambda \, d\phi$ , and  $\delta_{nn'}$  is the Kroneker delta  $(\delta_{nn'} = 1, \text{ if } n = n'; \delta_{nn'} = 0, \text{ if } n \neq n')$ . Substituting (24) into (23) and noting that the spherical harmonics are the eigenvectors of the Laplacian operator on the sphere, with  $-n(n+1)/a^2$  the associated eigenvalues, gives

$$\frac{\mathrm{d}\eta_n^m}{\mathrm{d}t} = -\kappa \frac{n(n+1)}{a^2} \eta_n^m,\tag{26}$$

which over the interval [0, T] readily integrates as

$$\eta_n^m(T) = \eta_n^m(0) \exp\{-\kappa T n(n+1)/a^2\}.$$
(27)

<sup>†</sup> As pointed out by Gneiting, the Gaussian  $e^{-(a\theta)^2/2L^2}$ , whose argument  $a\theta$  follows from a first-order development of (22) for small  $\theta$ , is not a valid (positive definite) correlation function on the sphere even though it has often been used in meteorology and oceanography.
The spectral coefficients  $\eta_n^m(0)$  can be determined by multiplying (24) by  $Y_n^{m^*}(\lambda, \phi)$  and applying the orthonormality condition (25):

$$\eta_n^m(0) = \frac{1}{4\pi a^2} \int_{\Sigma} \eta(\lambda, \phi, 0) Y_n^{m^*}(\lambda, \phi) \,\mathrm{d}\Sigma.$$
(28)

Substituting (28) into (27), and (27) into (24), and applying the Addition Theorem for the spherical harmonics (Arfken 1966, p. 450), the solution to (23) can be written as

$$\eta(\lambda, \phi, T) = \frac{1}{4\pi a^2} \int_{\Sigma'} f(\theta; \kappa T) \eta(\lambda', \phi', 0) \, \mathrm{d}\Sigma', \tag{29}$$

where

$$f(\theta;\kappa T) = \sum_{n=0}^{\infty} f_n P_n^0(\cos\theta) = \sum_{n=0}^{\infty} (2n+1)^{1/2} \exp\{-\kappa T n(n+1)/a^2\} P_n^0(\cos\theta),$$
(30)

 $\theta$  being the angular separation between the points  $(\lambda, \phi)$  and  $(\lambda', \phi')$  on the sphere, and  $P_n^0$  the Legendre polynomials. As the coefficients  $f_n$  of the Legendre polynomials are positive, the integral kernel  $f(\theta; \kappa T)$  is the *representation* of an isotropic covariance function (Courtier *et al.* 1998). Equation (29) thus defines a covariance operator.

The value of the covariance function at the origin ( $\theta = 0$ ) gives the variance at any point:

$$f(0; \kappa T) = \sum_{n=0}^{\infty} f_n (2n+1)^{1/2} = \sum_{n=0}^{\infty} (2n+1) \exp\{-\kappa T n(n+1)/a^2\},$$
 (31)

where we have used the fact that  $P_n^0(1) = (2n+1)^{1/2}$  as implied by the normalization in (25). Each coefficient  $f_n(2n+1)^{1/2}$  gives the contribution to the total variance from a given total wave number *n*, and hence provides the variance power spectrum of  $f(\theta; \kappa T)$ . Furthermore,  $f(\theta; \kappa T)/f(0; \kappa T)$  defines the associated correlation function.

The length-scale of the correlation function can be defined following Daley (1991):

$$L^{2} = -2 \frac{f(0; \kappa T)}{\nabla^{2} f(0; \kappa T)}$$
  
=  $2a^{2} \frac{\sum_{n=0}^{\infty} (2n+1) \exp\{-\kappa T n(n+1)/a^{2}\}}{\sum_{n=0}^{\infty} n(n+1)(2n+1) \exp\{-\kappa T n(n+1)/a^{2}\}},$  (32)

which shows that, as in the 1D example,  $\kappa T$  is the parameter controlling the length-scale of the correlation function. Moreover, as shown in the appendix, for scales relevant in meteorology and oceanography ( $L \ll a$ ), matching (30) to a Gaussian leads to

$$L^2 \approx 2\kappa T,\tag{33}$$

in direct analogy with the 1D example.

# (c) Extension to a larger class of isotropic correlation functions on the sphere

The correlation model can be extended to represent a larger class of correlation functions by considering the solution of a more general partial differential equation,

$$\frac{\partial \eta}{\partial t} + \sum_{p=1}^{P} \kappa_p (-\nabla^2)^p \eta = 0, \qquad (34)$$

where  $\kappa_p$ ,  $p = 1, \ldots, P$ , are non-negative weighting coefficients. The classical diffusion equation (23) is a special case of (34) with  $\kappa_p = 0$  for all p > 1. Thus, we will refer to (34) hereafter as a generalized 'diffusion' equation (GDE).

The general solution of (34) is of the form (29),

$$\eta(\lambda, \phi, T) = \frac{1}{4\pi a^2} \int_{\Sigma'} f(\theta; \kappa_1 T, \dots, \kappa_P T) \eta(\lambda', \phi', 0) \, \mathrm{d}\Sigma', \qquad (35)$$

with integral kernel given by

$$f(\theta; \kappa_1 T, \dots, \kappa_P T) = \sum_{n=0}^{\infty} (2n+1)^{1/2} \exp\left[-\sum_{p=1}^{P} \kappa_p T \{n(n+1)/a^2\}^p\right] P_n^0(\cos\theta).$$
(36)

Equation (36) defines a family of isotropic covariance functions on the sphere; the associated correlation functions are  $f(\theta; \kappa_1 T, \ldots, \kappa_P T)/f(0; \kappa_1 T, \ldots, \kappa_P T)$ . The length-scale of the correlation function can be defined in the usual way:

$$L^{2} = -2 \frac{f(0; \kappa_{1}T, \dots, \kappa_{P}T)}{\nabla^{2} f(0; \kappa_{1}T, \dots, \kappa_{P}T)}$$
  
=  $2a^{2} \frac{\sum_{n=0}^{\infty} (2n+1) \exp\left[-\sum_{p=1}^{P} \kappa_{p}T\{n(n+1)/a^{2}\}^{p}\right]}{\sum_{n=0}^{\infty} n(n+1)(2n+1) \exp\left[-\sum_{p=1}^{P} \kappa_{p}T\{n(n+1)/a^{2}\}^{p}\right]}.$  (37)

The free parameters controlling the spectrum (shape) and length-scale of the correlation function are the sequence of products  $\kappa_p T$ , p = 1, ..., P, where P is defined such that  $\kappa_P \neq 0$  and  $\kappa_p = 0$  for all p > P. For the special case P = 1, (36) and (37) reduce to (30) and (32), respectively, with  $\kappa_1 = \kappa$ .

Figure 1(a) illustrates the variance power spectrum and Fig. 1(b) the grid-point representation of four different correlation functions generated numerically using (36) with a truncation at total wavenumber 106. The solid curves in Fig. 1 correspond to the near-Gaussian (P = 1) from Fig. A.1 in the appendix. The dashed and dashed-dotted curves correspond to P = 2 and P = 3 respectively, with  $\kappa_p = 0$  for all p < P in both cases. The dotted curves also correspond to P = 2 but with non-zero values defined for both  $\kappa_1$  and  $\kappa_2$ . For each correlation function the (non-zero) products  $\kappa_p T$  have been tuned according to (37) in order to give a length-scale of 500 km as in the Gaussian example in Fig. A.1.

The effect of including higher-order terms in the GDE (while holding the lengthscale fixed) is clear. From the spectrum, there is a progressive decrease in the variance at large scales (wave numbers 1 to 14), increase in the variance at intermediate scales (wave numbers 15 and 30), and sharper drop-off in the variance at small scales (wave numbers greater than 30). In grid-point space, the correlation functions are very similar over  $0 \le r < L$ . The higher-order terms in the GDE, however, increase the rate at which the correlation function drops off to zero over L < r < 2L, and increase the amplitude of the oscillations about the zero axis for r > 2L, with these oscillations tending to zero at large separation distances.

The possibility of representing oscillatory correlation functions is an attractive feature of the correlation model given that the background errors of certain geophysical fields (such as temperature) often display a tendency to oscillate about the zero axis (Hollingsworth and Lönnberg 1986; Carton *et al.* 2000). This feature is associated with the higher-order terms in the GDE and is reflected by the increase in spectral variance at



Figure 1. (a) The variance power spectrum and (b) grid-point values for four different correlations functions  $f(\theta; \kappa_1 T, \ldots, \kappa_P T)/f(0; \kappa_1 T, \ldots, \kappa_P T)$  generated from (36). For each correlation function, the diffusion parameters  $\kappa_P T$  have been tuned to give a length-scale of L = 500 km. The solid curves correspond to the near-Gaussian function in Fig. A.1 (P = 1 and  $\kappa_1 T = 3.08 \times 10^{-3}$ ). The other curves correspond to parameter choices P = 2,  $\kappa_1 = 0$  and  $\kappa_2 T = 3.02 \times 10^{-6}$  (dashed); P = 2,  $\kappa_1 T = 1.54 \times 10^{-3}$  and  $\kappa_2 T = 1.39 \times 10^{-6}$  (dotted); and P = 3,  $\kappa_1 = \kappa_2 = 0$  and  $\kappa_3 T = 3.78 \times 10^{-9}$  (dashed-dotted). See text for further explanation.

intermediate scales. On the other hand, the increase in the drop-off rate of the spectrum at high wave numbers is possibly a less desirable consequence of the higher-order terms as there is some evidence to suggest that the variance of background errors has a much broader spectral range: i.e. the tail of the spectrum tapers off much more gradually than in the curves in Fig 1(a) (Julian and Thiébaux 1975; Thiébaux *et al.* 1986; Rabier *et al.* 1998). The spectral characteristics of the correlation functions generated by the GDE are indeed similar to those of correlation functions implied by smoothing splines (Wahba and Wendelberger 1980; Bennett 1992). This is not surprising in view of the fact that both approaches advocate the use of high-order Laplacian operators; in the GDE they are used to model the correlation functions directly, whereas in spline smoothing they appear as penalty terms in the inner product defining the  $J_b$  term.

# (d) Modelling vertical correlation functions

In direct analogy to (34), a generalized 1D diffusion equation can be used as the basis for modelling a class of 1D correlation functions. Consider the equation

$$\frac{\partial \eta}{\partial t} + \sum_{r=1}^{R} \kappa_r \left( -\frac{\partial^2}{\partial z^2} \right)^r \eta = 0, \tag{38}$$

where R and  $\kappa_r$  are analogous to P and  $\kappa_p$  in (34). The general solution of (38) on  $\mathbb{R}$  can be easily obtained using the Fourier transform, following the example in section 3(a). In particular, the solution can be expressed as an integral operator, the kernel of which is a covariance function given by the inverse Fourier transform

$$h(z; \kappa_1 T, \dots, \kappa_R T) = \frac{1}{2\pi} \int_{\widehat{z}} \exp\left(-\sum_{r=1}^R \kappa_r T \widehat{z}^{2r}\right) \mathrm{e}^{(-\widehat{t}\widehat{z}\widehat{z})} \,\mathrm{d}\widehat{z},\tag{39}$$

where  $\hat{i} = (-1)^{1/2}$ . As with  $f(\theta; \kappa_1 T, \ldots, \kappa_P T)$ , the sequence of products of the diffusion coefficients  $\kappa_r$  with the integration time T controls the shape and length-scale of  $h(z; \kappa_1 T, \ldots, \kappa_R T)$ . For R = 1, (39) reduces analytically to the Gaussian in (20). For R > 1, with  $\kappa_R \neq 0$ , there is no obvious closed analytical form of the covariance functions; however, for given values of  $\kappa_r T$ , (39) can be solved numerically to yield a family of covariance functions similar to those in Fig. 1. Taking z as the vertical coordinate in the model, then (38), together with appropriate top and bottom boundary conditions, provides a statistical model that may be used for representing the vertical correlations in the autocorrelation matrix  $\mathbf{C}_{\alpha}$ , as discussed further in sections 4(c) and 5(d).

# 4. PRACTICAL IMPLEMENTATION

In the previous section, we showed that the integral solution of the diffusion equation on the sphere could be interpreted as a covariance operator. In particular, we showed that the kernel of the covariance operator could be expressed as an isotropic covariance function that could be closely matched to a Gaussian. The possibility of representing other isotropic covariance functions by solving a more general differential equation was also discussed. In this section, we discuss the practical implementation of the algorithm within the design of a univariate correlation operator. Some generalizations of the algorithm will be discussed and illustrated in section 5 with reference to specific examples.

# (a) Discrete representation of the generalized diffusion equation within a horizontal, univariate correlation operator

Consider the following semi-discrete representation of the GDE in (34), based on an explicit forward differencing of the temporal derivative:

$$\eta(t_m) = \eta(t_{m-1}) - \sum_{p=1}^{P} \kappa_p \Delta t \, (-\nabla^2)^p \eta(t_{m-1}), \tag{40}$$

where  $\Delta t = t_m/m$  is the step size and *m* is a positive integer. The total integration time is  $T = t_M$ , *M* being the total number of integration steps. The sequence of products  $\kappa_p T = \kappa_p M \Delta t$ ,  $p = 1, \ldots, P$ , determines the spectrum of the covariance function from (36), and the length-scale of the covariance function from (37). Once these parameters are specified,  $\kappa_p \Delta t$  in (40) is set to  $\kappa_p T/M$  with the value of *M* chosen large enough so as to maintain the numerical stability of the forward differencing scheme. For the diffusion equation (P = 1), the stability condition is roughly  $\kappa_1 \Delta t/e^2 < 1/4$ where *e* is the grid spacing and, from (33),  $\kappa_1 \Delta t \approx L^2/2M$ . For a given *L*, this leads to an approximate requirement that  $M > 2(L/e)^2$ .

In order to be consistent with the formulation of the ocean model used in the illustrative examples in section 5, it is convenient to express the Laplacian in terms of a set of general orthogonal curvilinear coordinates (i, j). Assuming that the geographical coordinates  $(\lambda, \phi)$  are continuous and differentiable functions of (i, j), the Laplacian in (40) takes the form

$$\nabla^2 \eta = \frac{1}{e_1 e_2} \left\{ \frac{\partial}{\partial i} \left( \frac{e_2}{e_1} \frac{\partial \eta}{\partial i} \right) + \frac{\partial}{\partial j} \left( \frac{e_1}{e_2} \frac{\partial \eta}{\partial j} \right) \right\},\tag{41}$$

where  $e_1 = e_1(\lambda(i, j), \phi(i, j)) \equiv \partial s_1/\partial i$  and  $e_2 = e_2(\lambda(i, j), \phi(i, j)) \equiv \partial s_2/\partial j$  are metric coefficients that define the curvilinear distance elements  $(ds_1, ds_2) = (e_1di, e_2dj)$  in the (i, j) coordinate system. Equation (41) is self-adjoint for the scalar product  $\langle \eta_1, \eta_2 \rangle = \int \eta_1 \eta_2 ds_1 ds_2$ : i.e.  $\langle \nabla^2 \eta_1, \eta_2 \rangle = \langle \eta_1, \nabla^2 \eta_2 \rangle$ . As shown in what follows, this fundamental property of the Laplacian can be exploited to identify a factorization of the correlation operator of the form  $\mathbf{C}^{1/2} \mathbf{C}^{T/2}$ .

The discrete operator,  $L_P$ , that defines the integral solution (35) can be represented by an application of (40) over M steps, with M assumed to be even for convenience:

$$\eta(t_M) = \mathbf{L}_P \eta(t_0), \tag{42}$$

where

$$\mathbf{L}_{P} = \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{D})^{p} \right\}^{M},$$
(43)

**D** denoting a self-adjoint matrix representation of the discretized Laplacian. Since **D** is self-adjoint, it may factored as  $\mathbf{D} = \mathbf{W}^{-1}\mathbf{D}^{T}\mathbf{W}$  where **W** is a diagonal matrix of local

area elements  $\Delta s_1 \Delta s_2$ . L<sub>P</sub> can then be factored as

$$\begin{aligned} \mathbf{L}_{P} &= \mathbf{L}_{P}^{1/2} \mathbf{L}_{P}^{1/2} \\ &= \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{D})^{p} \right\}^{M/2} \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{D})^{p} \right\}^{M/2} \\ &= \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{W}^{-1} \mathbf{D}^{\mathrm{T}} \mathbf{W})^{p} \right\}^{M/2} \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{D})^{p} \right\}^{M/2} \\ &= \mathbf{W}^{-1} \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{D}^{\mathrm{T}})^{p} \right\}^{M/2} \mathbf{W} \left\{ \mathbf{I} - \sum_{p=1}^{P} \kappa_{p} \Delta t (-\mathbf{D})^{p} \right\}^{M/2} \\ &= \mathbf{W}^{-1} \mathbf{L}_{P}^{T/2} \mathbf{W} \mathbf{L}_{P}^{1/2} \\ &= \mathbf{L}_{P}^{1/2} \mathbf{W}^{-1} \mathbf{L}_{P}^{T/2} \mathbf{W}. \end{aligned}$$
(44)

 $L_P$  is self-adjoint with respect to the metric W: i.e.  $L_P$  maps a field from the primal space (e.g. model space) into itself. The action of  $L_P$  thus conserves the physical units of the input field.  $L_P$  can be transformed into a symmetric covariance matrix by multiplying it either to the left by W or to the right by  $W^{-1}$ , as evident from the factorization (44). We may define the covariance matrix either way, but its interpretation is somewhat clearer by following the latter. Therefore, without loss of generality, we define the covariance matrix as

$$\mathbf{L}_{P}\mathbf{W}^{-1} = \mathbf{L}_{P}^{1/2}\mathbf{W}^{-1}\mathbf{L}_{P}^{T/2}.$$
(45)

The associated covariance operator then corresponds to an integration of the difference equation (40) from an initial condition scaled at each grid point by the inverse of the corresponding local-area element. This is consistent with the general definition of a covariance operator, which represents a mapping from the dual space into the primal space and hence does not conserve the physical units of the input field (Tarantola 1987).

A matrix,  $\Lambda$ , of normalization coefficients is required to convert (45) into a correlation matrix: i.e.  $\Lambda$  ensures that the standard-error deviations effectively used in  $\mathbf{B}_u$  (see (14)) are indeed those which are specified through  $\Sigma$ . In general, the elements of  $\Lambda$  will be spatially dependent (see section 5(b)) in which case  $\Lambda$ , like  $\Sigma$ , must be introduced on either side of the filter covariance matrix in order to maintain symmetry (cf. the scalar normalization in (21)). Thus, the auto-correlation matrix associated with an arbitrary scalar field  $\alpha$  has the form

$$\mathbf{C}_{\alpha} = \mathbf{\Lambda} \mathbf{L}_{p}^{1/2} \mathbf{W}^{-1} \mathbf{L}_{p}^{T/2} \mathbf{\Lambda}$$

$$(46)$$

$$(46)$$

$$= (\mathbf{A}\mathbf{L}_{\vec{p}} \ \mathbf{W}^{-1/2})(\mathbf{A}\mathbf{L}_{\vec{p}} \ \mathbf{W}^{-1/2})^{-1}$$
$$= \mathbf{C}_{\alpha}^{1/2}\mathbf{C}_{\alpha}^{T/2}.$$
(47)

The procedure for computing the normalization coefficients is discussed in some detail in the next section.

The presence of W in the general matrix expression for  $C_{\alpha}$  illustrates explicitly the grid-dependence of the correlation operator and hence of the scalar product of the background term  $J_b$  in (1) whose metric depends on  $C_{\alpha}^{-1}$ . This scalar product defines a weighted sum of increments over the model domain and thus naturally should depend on the model grid.

# A, WEAVER and P. COURTIER

Lorenc (1997) and Parrish *et al.* (1997) discuss a similar procedure for designing a 2D horizontal covariance operator using a 1D recursive filter. Their 2D covariance operator is formed from a matrix product of the 1D filter, applied successively along orthogonal grid directions (cf.  $\mathbf{L}_{P}^{1/2}$  in (44)), with its adjoint (cf.  $\mathbf{W}^{-1}\mathbf{L}_{P}^{T/2}\mathbf{W}$  in (44)). The resulting covariance operator, however, does not give an isotropic response as illustrated by Lorenc (1997). This contrasts the Laplacian-based filter ( $\mathbf{L}_{P}$ ) which is isotropic by construction. (A numerical validation of this result is provided by Fig. 2(a)). Furthermore, the Laplacian can be made invariant in any orthogonal curvilinear coordinate system, a feature that is particularly useful for models with distorted grids, such as global ocean models that have the north pole singularity rotated onto a land point (Madec and Imbard 1996).

# (b) The normalization matrix

As discussed in the previous section, a diagonal normalization matrix  $\Lambda$  is required to convert the covariance matrix  $\mathbf{L}_P \mathbf{W}^{-1}$  into a correlation matrix of the form (47)\*. The diagonal elements of  $\Lambda$  are the inverse of the square root of the corresponding diagonal elements of  $\mathbf{L}_P \mathbf{W}^{-1}$ : i.e. the inverse of the intrinsic standard deviations (square root of the variances) of  $\mathbf{L}_P \mathbf{W}^{-1}$ . Note that the variances of  $\mathbf{L}_P \mathbf{W}^{-1}$  have physical dimensions of inverse length squared, in contrast to the variances of  $\mathbf{B}_u$  which have physical dimensions equal to the square of the physical units of the background variables.

For the isotropic filter  $\mathbf{L}_P$ , applied in a boundary-free domain,  $\mathbf{\Lambda}$  is simply a constant multiple of the identity matrix where the constant can be determined analytically from the variance of the covariance functions in (36);  $\mathbf{\Lambda} = (1/\sqrt{t})\mathbf{I}$  where  $t = f(0; \kappa_1 T, \ldots, \kappa_P T)/4\pi a^2$ , the extra factor  $4\pi a^2$  coming from the normalization in (35). For the  $\mathbf{L}_1$ -filter,  $t \approx g(0; \gamma)/4\pi a^2 \approx 1/2\pi L^2$  (cf. (A.13) in the appendix).

In section 5, we discuss more general implementations of the filter (e.g. applications in bounded domains (section 5(b)) and extensions to account for anisotropic correlations (sections 5(c-d)) for which the true variances are no longer well approximated at all grid points by a constant. To compute the true variances at each grid point requires a specific algorithm. Two algorithms are proposed below, one of which computes the variances exactly, the other approximately.

Let  $t_l$  denote the *l*-th diagonal element of the matrix  $\mathbf{L}_P \mathbf{W}^{-1}$ , i.e. the intrinsic variance of the covariance filter at the *l*-th model grid point. Since, in practice,  $\mathbf{L}_P \mathbf{W}^{-1}$ is available as an operator, the element  $t_l$  can be computed by applying  $\mathbf{L}_P \mathbf{W}^{-1}$  to the vector  $\mathbf{e}_l = (0, \dots, 0, 1, 0, \dots, 0)^T$ , the entry equal to one being located at the *l*-th grid point. The value of the filtered field at the *l*-th grid point gives the filter variance at that point: i.e.  $t_l = \mathbf{e}_l^T \mathbf{L}_P \mathbf{W}^{-1} \mathbf{e}_l$  with  $1/\sqrt{t_l}$  defining the *l*-th diagonal element of  $\mathbf{A}$ . This algorithm is computationally expensive as it requires, in general, as many applications of the covariance filter as there are model grid points. The algorithm is, however, well suited for parallel computers as the filter variance at the different points can each be computed on a separate processor.

Rather than computing the variances exactly using the above algorithm, a cheaper alternative is to use a randomization method (Fisher and Courtier 1995; Andersson *et al.* 2000) to estimate the variances approximately. Consider the transformation  $\tilde{\mathbf{v}} = \mathbf{L}_P^{1/2} \mathbf{W}^{-1/2} \mathbf{v}$  where  $\mathbf{v}$  is a Gaussian random variable, drawn from a population

<sup>\*</sup>  $L_P$  is a scale-dependent filter, with the property that it conserves the spatial mean of the input field. For this reason, the diagonal elements of the covariance matrix  $L_P W^{-1}$  will not, in general, be equal to one.



Figure 2. The autocorrelation field generated by the  $L_1^h$ -filter (i.e. by applying (46) with P = 1 to the vector  $e_l = (0, \ldots, 0, 1, 0, \ldots, 0)^T$ ). The correlation length-scale is (a) 4° and (b) 1°. The zonal resolution is 1° and the meridional resolution is approximately 0.5°. The horizontal and vertical axes are in degrees longitude and latitude respectively.

#### A. WEAVER and P. COURTIER

having zero mean and unit variance: i.e.  $E[\mathbf{v}] = 0$  and  $E[\mathbf{v}\mathbf{v}^T] = \mathbf{I}$ . The transformed random variable verifies  $E[\mathbf{\tilde{v}}] = 0$  and  $E[\mathbf{\tilde{v}}\mathbf{\tilde{v}}^T] = \mathbf{L}_p^{1/2}\mathbf{W}^{-1}\mathbf{L}_p^{T/2} = \mathbf{L}_P\mathbf{W}^{-1}$ . Therefore, from an ensemble of Q random vectors  $\mathbf{\tilde{v}}_q$ ,  $q = 1, \ldots, Q$ , the variances of the filter covariance matrix can be estimated from the diagonal elements of

$$\mathbf{L}_{P}\mathbf{W}^{-1} \approx \frac{1}{Q} \sum_{q=1}^{Q} \widetilde{\mathbf{v}}_{q} \widetilde{\mathbf{v}}_{q}^{\mathrm{T}} = \frac{1}{Q} \sum_{q=1}^{Q} (\mathbf{L}_{P}^{1/2} \mathbf{W}^{-1/2} \mathbf{v}_{q}) (\mathbf{L}_{P}^{1/2} \mathbf{W}^{-1/2} \mathbf{v}_{q})^{\mathrm{T}}.$$
 (48)

The variances converge to the true variances when Q is large. For Gaussian statistics, it can be shown (Barlow 1989, p. 89) that the standard error in the estimated standard deviations (inverse of the normalization factors) is  $1/(2Q)^{1/2}$  (e.g. 10% for an ensemble of 50).

#### (c) A three-dimensional correlation model

In a multi-level ocean or atmosphere model, the univariate correlation matrix  $C_{\alpha}$  must account for error correlations in the vertical as well as the horizontal. Threedimensional correlation models are often assumed to be separable in that they are constructed as a product of a horizontal and a vertical correlation function, which depend only on horizontal and vertical separation, respectively (Daley 1991). Specifically, the function formed from the product of the covariance functions (36) and (39), normalized by their respective variances, corresponds to a special class of separable, 3D correlation functions, which may be modelled by a combined integration of the 1D and 2D GDEs.

Assuming that the vertical direction is normal to the geopotential surfaces on the sphere, and that the vertical distance z is a continuous function of a generalized vertical coordinate k, orthogonal to the curvilinear coordinate surfaces (i, j) defined in section 4(a), the vertical second-derivative operator in (38) can be expressed as

$$\frac{\partial^2 \eta}{\partial z^2} = \frac{1}{e_3} \left\{ \frac{\partial}{\partial k} \left( \frac{1}{e_3} \frac{\partial \eta}{\partial k} \right) \right\},\tag{49}$$

where  $e_3 = \partial s_3 / \partial k = \partial z / \partial k$  is the metric coefficient defining the vertical distance element  $ds_3 = e_3 dk$  in this coordinate system. In what follows, we will use the subscript/superscript v to distinguish the parameters and operators associated with the vertical correlation model from those associated with the horizontal correlation model of section 4(a), which for notational consistency will hereafter be denoted with a subscript/superscript h. Let  $\mathbf{D}^{v}$  be the matrix representation of a discretized version of (49). The matrix operator,  $\mathbf{L}_{R}^{v}$ , that integrates the 1D (vertical) GDE from t = 0 to  $t = T_{v} = M_{v} \Delta t_{v}$  can thus be expressed as

$$\mathbf{L}_{R}^{\mathsf{v}} = \left\{ \mathbf{I} - \sum_{r=1}^{R} \kappa_{r}^{\mathsf{v}} \Delta t_{\mathsf{v}} (-\mathbf{D}^{\mathsf{v}})^{r} \right\}^{M_{\mathsf{v}}}.$$
(50)

The symmetric vertical covariance matrix associated with (50) is then  $L_R^v W_v^{-1}$ , where  $W_v$  is a diagonal matrix of vertical grid elements  $\Delta s_3$  (cf. (45)).

The 3D covariance operator acting on an arbitrary scalar field can be represented by the following two-step procedure: (i) for each model level k, integrate the 2D GDE (34) from t = 0 to  $t = T_h$  with the scalar field taken as initial condition; (ii) for each horizontal grid point (i, j), integrate the vertical GDE (38) from t = 0 to  $t = T_v$  using the result from (i) as initial condition. Noting that the horizontal and vertical operators commute, we can express this 3D covariance operator as

$$\mathbf{L}_{R}^{v} \mathbf{W}_{v}^{-1} \mathbf{L}_{P}^{h} \mathbf{W}_{h}^{-1} = \mathbf{L}_{R}^{v}{}^{1/2} \mathbf{W}_{v}^{-1} \mathbf{L}_{R}^{v}{}^{T/2} \mathbf{L}_{P}^{h}{}^{1/2} \mathbf{W}_{h}^{-1} \mathbf{L}_{P}^{h}{}^{T/2}$$
$$= \mathbf{L}_{R}^{v}{}^{1/2} \mathbf{L}_{P}^{h}{}^{1/2} \mathbf{W}^{-1} \mathbf{L}_{P}^{h}{}^{T/2} \mathbf{L}_{R}^{v}{}^{T/2},$$
(51)

where  $\mathbf{W} = \mathbf{W}_{h}\mathbf{W}_{v}$  is a diagonal matrix of volume elements  $\Delta s_{1}\Delta s_{2}\Delta s_{3}$ . Thus, the 3D correlation matrix has the same form as (46) with  $\mathbf{L}_{P}$  replaced by the product  $\mathbf{L}_{R}^{v}\mathbf{L}_{P}^{h}$ . Note that the elements of the normalization matrix  $\mathbf{A}$  must now have physical dimensions of volume in order to render the correlation matrix dimensionally consistent.

In summary, the square-root factors of the three-dimensional correlation matrix take the form

$$\mathbf{C}_{\alpha}^{1/2} = \mathbf{\Lambda} \mathbf{L}_{R}^{\mathbf{v}\ 1/2} \mathbf{L}_{P}^{\mathbf{h}\ 1/2} \mathbf{W}^{-1/2},\tag{52}$$

$$\mathbf{C}_{\alpha}^{\mathrm{T/2}} = \mathbf{W}^{-1/2} \mathbf{L}_{P}^{\mathrm{h} T/2} \mathbf{L}_{R}^{\mathrm{v} T/2} \mathbf{\Lambda}.$$
 (53)

In words, the operator  $C_{\alpha}^{1/2}$ , which in 3D-Var is needed for transforming the increment from minimization (control) space into model space (15), involves the following sequence of operations:

(i) multiply each element of the input vector by the inverse of the square root of its associated volume element;

(ii) perform  $M_{\rm h}/2$  integration steps of the horizontal diffusion equation;

(iii) perform  $M_v/2$  integration steps of the vertical diffusion equation;

(iv) multiply each element of the filtered vector by its corresponding normalization factor.

Likewise, the 'adjoint' operator  $C_{\alpha}^{T/2}$ , which in 3D-Var is needed for defining  $\mathbf{B}^{T/2}$  in the gradient transformation (9), involves the sequence of operations 1 through 4 applied in reverse order, with the horizontal and vertical diffusion equations in operations (ii) and (iii) replaced by their respective adjoints.

# 5. ILLUSTRATIONS

The univariate  $C_{\alpha}$ -operator described in the previous section has been incorporated within a variational assimilation system (Weaver and Vialard 2000) being developed for the OPA OGCM (Madec *et al.* 1999). In this section, this system will be used to illustrate some of the basic features of the correlation model, including generalizations that account for anisotropic correlations. In all examples, the horizontal and/or vertical L<sub>1</sub>-filter is used.

The discrete differential operators used in the correlation model are based on similar operators used to solve the equations of the ocean model. These operators are formulated in orthogonal curvilinear coordinates and discretized on an Arakawa C-grid using second-order accurate, centred finite differences. The examples presented here have been conducted with a tropical Pacific version of the model. The configuration extends from  $120^{\circ}$ E to  $70^{\circ}$ W and from  $30^{\circ}$ S to  $30^{\circ}$ N. The entire model domain is enclosed by solid boundaries, with realistic coastlines imposed along the eastern and western boundaries, and artificial walls along the southern and northern boundaries. The zonal resolution is  $1^{\circ}$  and the meridional resolution varies smoothly from  $0.5^{\circ}$  at the equator to  $2^{\circ}$  at the northern and southern boundaries. The basic configuration has been adapted from Vialard and Delecluse (1998). The interested reader should refer to that paper and references therein for a more complete description of the model.

#### (a) Isotropic correlations

The auto-correlations implied by the  $L_1^h$ -filter are illustrated in Fig. 2 for a point (l) located on the equator at 160°W (i.e. in the central Pacific, far away from boundaries). The autocorrelation field has been obtained by applying (46) with P = 1 to the unit vector  $\mathbf{e}_l$  defined previously. Assuming that the variance of background errors is constant (i.e. assuming  $\Sigma = \sigma_b I$  with  $\sigma_b$  constant), then Fig. 2 corresponds to the analysis increment, up to a proportionality constant, that would be obtained from a single observation. This is easy to see from (12) noting that for a single observation, w becomes a scalar and  $\mathbf{H}^{\mathrm{T}}$  a column vector equal to  $\mathbf{e}_{l}$ , where the value of one in the *l*-th entry of  $e_l$  corresponds to the observation location. Figure 2(a) is a numerical validation of the isotropic response of the filter for a value of L large relative to the grid resolution; Fig. 2(b) illustrates how the degree of isotropy is degraded when Lis comparable to the grid resolution. In this latter example, L is equal to the zonal resolution and approximately twice the meridional resolution. Relative to the true correlations computed directly from the Gaussian function, those computed in Fig. 2(b) using the diffusion model are everywhere underestimated, particularly in the zonal direction where the resolution is coarsest.

# (b) Boundary considerations

In a bounded domain the Laplacian must be supplemented with boundary conditions. For the C-grid representation of the discrete  $L_1^h$ -filter, this requires specifying the value of the filtered variable's first derivatives normal to boundary points. For the higherorder  $L_P^h$ -filters (P > 1), additional boundary conditions are required; in general, the first, third, ..., and (2P - 1)th derivatives normal to the boundary must be specified. These boundary conditions can be viewed as free parameters of the correlation model.

Figure 3 illustrates the auto-correlations generated by the  $L_1^h$ -filter at three points adjacent to the coastlines, assuming vanishing first derivatives normal to boundary points (in the ocean model this boundary condition is analogous to the 'no-flux' condition on tracer fields or the 'free-slip' condition on the tangential component of the velocity field). The parameters of the correlation model correspond to those used for the  $L_1^h$ -filter in Fig. 2(a). The spatial structure of the correlations near boundary regions is determined primarily by the geometry of the coastlines; the exact nature of the boundary condition has a much weaker influence. For example, the correlations in Fig. 3 are relatively robust to changes in boundary condition from a no-flux to a dissipative-flux form.

The primary effect of changing the nature of the boundary condition is to modify the intrinsic variances  $(t_l)$  of the filter, and hence the normalization factors  $(1/t_l^{1/2})$ , in the vicinity of the boundary. Figure 4 shows the difference between the true normalization factors,  $1/t_l^{1/2}$ , computed using the exact procedure described in section 4(b), and the constant normalization factor,  $(2\pi)^{1/2}L$ , estimated from the analytical solution of the diffusion equation in a boundary-free domain with constant diffusion coefficient. The values of the normalization factors match the analytical estimate very closely except in regions closely confined to the boundaries where they are smaller by up to a factor of 2.5: i.e. the filter variances are greater by up to a factor of 6. However, changing from a no-flux boundary condition to a stronger, dissipative-flux condition results in a reduction of the filter variances, as one might expect intuitively. In terms of cost efficiency, the no-flux condition is preferable as it imposes a less stringent constraint on the size of the time step needed to maintain the numerical stability of the forward-differencing scheme.



Figure 3. The autocorrelation field generated by the  $L_1^h$ -filter (see text) at three grid-points adjacent to the continental boundaries. The domain of the ocean model covers the tropical Pacific basin and is closed along all boundaries. The contour interval is 0.1.



Figure 4. The relative difference  $(1 - (t_l/t)^{1/2})$  between the true normalization factors  $(1/t_l^{1/2})$  for the  $L_1^h$ -filter at tracer grid-points in the uppermost level in the model and the constant normalization factor  $(t = (2\pi)^{1/2}L)$  estimated from the analytical solution of the diffusion equation in a boundary-free domain with constant diffusion coefficient. The boundary conditions used for the  $L_1^h$ -filter assume vanishing normal derivatives ('no flux') at the boundaries. The contour interval is 0.1. See text for explanation.

# (c) Anisotropic correlations: coordinate stretching

There is substantial observational evidence to suggest that correlation scales near the equator are longer in the zonal direction than in the meridional direction (e.g. Meyers *et al.* 1991; Carton *et al.* 2000). Such anisotropy can be accounted for in the diffusion equation by stretching the zonal coordinates in the Laplacian operator, as discussed in Derber and Rosati (1989) and Behringer *et al.* (1998). With reference to the GDE in (34), this can be achieved by replacing  $\nabla^2$  with  $\nabla \cdot \mathbf{R} \nabla$  where  $\mathbf{R}$  is a diagonal, second-rank tensor of non-dimensional 'diffusion' coefficients  $r_i$  and  $r_j$ ,

$$\boldsymbol{R} = \begin{pmatrix} r_i & 0\\ 0 & r_j \end{pmatrix}, \tag{54}$$

 $\nabla$  is the gradient operator

$$\nabla = \left(\frac{1}{e_1}\frac{\partial}{\partial i}, \frac{1}{e_2}\frac{\partial}{\partial j}\right)^{\mathrm{T}},\tag{55}$$

and  $\nabla$ .c is the divergence operator

$$\nabla \cdot \mathbf{c} = \frac{1}{e_1 e_2} \left\{ \frac{\partial (e_2 c_1)}{\partial i} + \frac{\partial (e_1 c_2)}{\partial j} \right\},\tag{56}$$

with  $\mathbf{c} = (c_1, c_2)^{\mathrm{T}}$ . The tensor elements  $r_i$  and  $r_j$  act, respectively, on the *i* and *j* components of the gradient. For the isotropic case,  $r_i = r_j \equiv 1$ . Note that in the ocean model used in the examples here, there is a one-to-one correspondence between the computational coordinates *i* and *j* and the geographical coordinates  $\lambda$  and  $\phi$ : i.e.  $\lambda = \lambda(i)$  and  $\phi = \phi(j)$ .

Figure 5 illustrates the auto-correlations generated at three different latitudes by an anisotropic version of the  $L_1^h$ -filter (cf. Fig. 1 in Behringer *et al.* (1998)). The zonal and



Figure 5. The autocorrelation field generated at latitudes (a)  $0^{\circ}$ , (b)  $5^{\circ}N$  and (c)  $15^{\circ}N$ , by an anisotropic version of the  $L_1^{h}$ -filter (see text). The contour interval is 0.1.

meridional length-scales have been defined to be 8° and 2°, respectively, at the equator, and 4° in both directions at  $|\phi| \ge 20^{\circ}$ N/S, with a linear transition between these values over 20°S <  $\phi$  < 20°N. This has been achieved by setting  $\kappa_1 T$  to  $L^2/2$  (as in Fig. 2(a)), and varying the tensor elements linearly from  $r_i = 4$  and  $r_j = 1/4$  at the equator to  $r_i = r_j = 1$  at  $|\phi| \ge 20^{\circ}$ N/S. More generally, the  $r_i$  and  $r_j$  can be made inhomogeneous functions of both spatial coordinates (i, j) in order to allow for fully global variations in the length-scale.

# (d) Three-dimensional flow-dependent correlations: coordinate rotation

The surface over which the horizontal correlations are effectively computed depends on the choice of vertical coordinate in the correlation model. In section 4(c), the direction of the vertical coordinate was assumed to be normal to the geopotential (z-)surfaces, implying that the horizontal correlation function (36) is isotropic along these surfaces. Other choices are of course possible. For example, in the ECMWF 3D-Var (Courtier et al. 1998), the horizontal surfaces are defined with respect to the hybrid  $\eta$ -coordinate (a pressure-based, terrain-following vertical coordinate) used in the ECMWF atmospheric model. In general, the vertical coordinate of the correlation model can be independent of the vertical coordinate used in the actual atmosphere or ocean model. For example, Daley and Barker (1999) suggest the use of potential temperature (isentropic) coordinates as a means of introducing more flow dependency into the background-error correlation model of an atmospheric assimilation system. The analogy in an ocean assimilation system is to define the background-error correlations with respect to a potential density (isopycnal) coordinate (e.g. as proposed by Gavart and DeMey (1997) in the context of empirical orthogonal functions). In this section, we illustrate that a transformation from, for example, geopotential to isopycnal coordinates is a natural extension of the correlation model based on the GDE. For simplicity, we consider only the L<sub>1</sub>-filter, although in principle the analysis may be extended to accommodate the higher-order filters.

In the 3D L<sub>1</sub>-filter, the length-scale of the horizontal Gaussian function is determined from  $L = (2\kappa_1^h T_h)^{1/2}$  and the depth-scale of the vertical Gaussian function from  $D = (2\kappa_1^v T_v)^{1/2}$ . In other words, for a given correlation scale, we are free to adjust both the diffusion coefficient and the integration time, providing their product remains constant. (Note that for the higher-order filters, the product also controls the shape of the correlation function.) We can exploit this fact to derive an alternative form of the L<sub>1</sub>-filter, which will simplify the transformation to isopycnal coordinates.

First, for given scales D and L, let us define identical time integration parameters for the vertical and horizontal L<sub>1</sub>-filters: i.e. we set  $M_v = M_h$  and  $\Delta t_v = \Delta t_h$  (and hence  $T_v = T_h$ ), and redefine the vertical diffusion coefficient to be  $\tilde{\kappa}_1^v = \kappa_1^h D^2/L^2$  so that  $D = (2\tilde{\kappa}_1^v T_h)^{1/2}$ . From (43) and (50), we can express the 3D L<sub>1</sub>-filter as

$$\mathbf{L}_{1} = \mathbf{L}_{1}^{\mathbf{v}} \mathbf{L}_{1}^{\mathbf{n}}$$

$$= \{\mathbf{I} + \kappa_{1}^{\mathbf{v}} \Delta t_{\mathbf{v}} \mathbf{D}^{\mathbf{v}}\}^{M_{\mathbf{v}}} \{\mathbf{I} + \kappa_{1}^{\mathbf{h}} \Delta t_{\mathbf{h}} \mathbf{D}^{\mathbf{h}}\}^{M_{\mathbf{h}}}$$

$$= \{(\mathbf{I} + \widetilde{\kappa}_{1}^{\mathbf{v}} \Delta t_{\mathbf{h}} \mathbf{D}^{\mathbf{v}})(\mathbf{I} + \kappa_{1}^{\mathbf{h}} \Delta t_{\mathbf{h}} \mathbf{D}^{\mathbf{h}})\}^{M_{\mathbf{h}}}$$

$$\approx (\mathbf{I} + \kappa_{1}^{\mathbf{h}} \Delta t_{\mathbf{h}} \mathbf{D}^{\mathbf{h}} + \widetilde{\kappa}_{1}^{\mathbf{v}} \Delta t_{\mathbf{h}} \mathbf{D}^{\mathbf{v}})^{M_{\mathbf{h}}}$$

$$= \left\{\mathbf{I} + \kappa_{1}^{\mathbf{h}} \Delta t_{\mathbf{h}} \left(\mathbf{D}^{\mathbf{h}} + \frac{\widetilde{\kappa}_{1}^{\mathbf{v}}}{\kappa_{1}^{\mathbf{h}}} \mathbf{D}^{\mathbf{v}}\right)\right\}^{M_{\mathbf{h}}},$$
(58)

where we have neglected the term of  $O(\Delta t_h^2)$ . (This term can be made small by choosing  $M_h$  sufficiently large, which we are free to do but at the expense of increasing the cost of the filter.)

Equation (58) is the matrix representation of a 3D Laplacian diffusion operator acting on  $0 \le t \le T_h$ . The 3D correlation function associated with this operator is horizontally/vertically anisotropic by virtue of the relative scaling factor  $\tilde{\kappa}_1^v/\kappa_1^h = D^2/L^2$ between the vertical and horizontal second-derivative operators, where  $D \ll L$  for meteorological and oceanographic scales of interest. The relevance of this alternative form of the L<sub>1</sub>-filter is that we are now in a position to exploit a classical transformation from geopotential to isopycnal coordinates (Redi 1982), as commonly employed in z-coordinate OGCMs for parametrizing lateral mixing along isopycnal surfaces (e.g. Pacanowski 1996; Madec *et al.* 1999). First, as in the previous section, we can express the 3D Laplacian in a more general tensorial form,  $\nabla$ .  $\mathbf{R}\nabla$ , where

$$\boldsymbol{R} = \begin{pmatrix} r_i & 0 & 0\\ 0 & r_j & 0\\ 0 & 0 & r_k \end{pmatrix},$$
(59)

$$\nabla = \left(\frac{1}{e_1}\frac{\partial}{\partial i}, \frac{1}{e_2}\frac{\partial}{\partial j}, \frac{1}{e_3}\frac{\partial}{\partial k}\right)^{\mathrm{T}},\tag{60}$$

and

$$\nabla \cdot \mathbf{c} = \frac{1}{e_1 e_2} \left\{ \frac{\partial (e_2 c_1)}{\partial i} + \frac{\partial (e_1 c_2)}{\partial j} \right\} + \frac{1}{e_3} \frac{\partial c_3}{\partial k}, \tag{61}$$

with  $\mathbf{c} = (c_1, c_2, c_3)^{\mathrm{T}}$ . The tensor elements  $r_i$ ,  $r_j$  and  $r_k$  can accommodate coordinate stretching in the horizontal and vertical planes and can be varied with each grid point (i, j, k) in order to permit geographical variations of the length- and depth-scales. Defining the length-scales to be a function of depth or the depth-scale to be a function of horizontal position is one way of introducing non-separability into the correlation model. Equation (58) can be seen as a special case with  $r_i = r_j = 1$  and  $r_k = \tilde{\kappa}_1^{\mathrm{v}}/\kappa_1^{\mathrm{h}}$ .

Now, we make the assumption that **R** is diagonal in an isopycnal coordinate system, rather than in a z-coordinate system as considered above. This requires a simple reinterpretation of the tensor elements in **R**:  $r_i$  and  $r_j$  are assumed to be referenced to the 'horizontal' plane defined by a constant isopycnal surface, while  $r_k$  is assumed to be referenced to the 'vertical' plane normal to this surface. The representation of the isopycnal **R** in z-coordinates (our computational coordinate system) can then be obtained by a coordinate rotation, as described in Redi (1982) (see also Griffies *et al.* (1998)). For simplicity, we consider only the horizontally isotropic case  $r_i = r_j$ .

At the beginning of an assimilation cycle, we assume the availability of a statically stable background potential-density surface,  $\rho^b$ , with respect to which we can define the 'horizontal' surfaces of the correlation model. In a z-coordinate OGCM,  $\rho^b$  is computed diagnostically from the equation of state, given the background potentialtemperature field  $T^b_{\theta}(i, j, k)$  and background salinity field  $S^b(i, j, k)$  (these being the prognostic density variables in the OGCM), and the local ocean depth z(k) as input (i.e.  $\rho^b = \rho^b \{T^b_{\theta}(i, j, k), S^b(i, j, k), z(k)\}$ ). Implicitly, in defining  $\rho^b$  as our coordinate, we are assuming that the background errors are isotropic along the  $\rho^b$  surfaces. Such an hypothesis can only be verified objectively by computing the actual statistics of the observation-minus-background field (e.g. as in Hollingsworth and Lönnberg (1986)). Nevertheless, the use of  $\rho^b$  as a coordinate has a certain appeal physically, which makes it an attractive possibility providing the background density state is reasonably accurate. If not, then there may be little virtue in using a  $\rho^{b}$ -coordinate. In such a case, a z-coordinate may be a wiser choice.

The diagonal tensor  $\mathbf{R}$  in the isopycnal coordinate system transforms as  $\tilde{\mathbf{R}} = \mathbf{S}\mathbf{R}\mathbf{S}^{\mathrm{T}}$  in the z-coordinate system, where S is a rotation matrix. The elements of S depend on the components and magnitude of the isopycnal slope vector  $(-a_1, -a_2, 0)$  where

$$a_1 = \frac{e_3}{e_1} \left(\frac{\partial \rho^{\rm b}}{\partial i}\right) \left(\frac{\partial \rho^{\rm b}}{\partial k}\right)^{-1} \quad a_2 = \frac{e_3}{e_2} \left(\frac{\partial \rho^{\rm b}}{\partial j}\right) \left(\frac{\partial \rho^{\rm b}}{\partial k}\right)^{-1}.$$
 (62)

(For details, we refer the reader to Redi (1982) or Griffies *et al.* (1998).) Since the slopes  $a_1$  and  $a_2$  are generally less than  $10^{-2}$  in the ocean and since  $r_k \sim O(D^2/L^2) \ll 1$  while  $r_i, r_j \sim O(1)$ , the full expression for  $\tilde{R}$  can be simplified appreciably. This leads to the familiar 'small slope' approximation to  $\tilde{R}$  commonly used in an isopycnal mixing parametrization of a z-coordinate model (cf. (3) in Griffies *et al.* (1998)):

$$\widetilde{\mathbf{R}} \approx \begin{pmatrix} r_i & 0 & -r_i a_1 \\ 0 & r_i & -r_i a_2 \\ -r_i a_1 & -r_i a_2 & r_k + r_i (a_1^2 + a_2^2) \end{pmatrix}.$$
(63)

Figure 6 shows an example of the auto-correlation fields generated by a z-coordinate and a  $\rho^b$ -coordinate version of the 3D L<sub>1</sub>-filter. Figure 6(a) shows a vertical meridional section of a typical background temperature field through 110°W in the eastern Pacific basin. The corresponding potential-density profile is very similar (i.e. the isopycnals are only weakly modified by salinity). The background field actually corresponds to the analysis at the end of a one-month 4D-Var analysis of *in situ* temperature observations, as described in Weaver and Vialard (2000). The longitude 110°W coincides with a particular section of the Tropical Atmosphere–Ocean array of moored buoys, the temperature measurements from which were assimilated in the 4D-Var experiment. In other words, the background profile at 110°W is already quite close to observations so there may be good merit in building this information into the background-error correlation model.

Figure 6(b) shows the auto-correlation field generated by the separable, z-coordinate, 3D  $L_1$ -filter (57) for a point located at 8°N and at a depth of 55 m. The depth-scale and length-scale have been set to constant values of 20 m and 4°, respectively. As the separable filter acts along geopotential surfaces, the resulting auto-correlation field is insensitive to the background profile, in particular to the strong meridional gradient in the temperature field between 4°N and 9°N which is associated with the North Equatorial Counter Current. Figure 6(c) shows the auto-correlation field obtained using the alternative form of the 3D  $L_1$ -filter (58) in conjunction with the isopycnal tensor (63). As expected, the correlations are strongest along isopycnal surfaces and fall off rapidly across isopycnal surfaces. In a 3D analysis, these correlations would essentially define the surfaces over which to smooth an observation located at the correlation point. Thus, a correlation model based on Fig. 6(c) would clearly be less destructive to the background density profile than a correlation model based on Fig. 6(b).

The correlations in Fig. 6(b) were computed using ten iterations of the vertical filter  $(L_1^v)$  and 250 iterations of the horizontal filter  $(L_1^h)$ . The correlations in Fig. 6(c) were computed using 8000 iterations of the 3D  $L_1$ -filter, which was roughly the number required for numerical stability. This represents a substantial increase (roughly a factor of 30) in cost relative to the separable filter. Furthermore, as the isopycnal correlation



Figure 6. (a) A meridional section of a typical background potential-temperature field in the eastern tropical Pacific 110°W. (b) The auto-correlation field at a depth of 55 m and latitude 8°N generated by a separable form of the 3D L<sub>1</sub>-filter (57) defined with respect to the geopotential coordinate system. (c) The corresponding auto-correlation field obtained using the 3D L<sub>1</sub>-filter (58) defined with respect to an isopycnal coordinate system based on the background isopycnal surfaces associated with (a).

### A. WEAVER and P. COURTIER

model is dependent on the background state, in a cycling assimilation system the normalization factors would need to be recomputed at the start of each assimilation cycle, which could only be done feasibly using, for example, the randomization method described in section 4(b). In a z-coordinate correlation model, the normalization factors must be computed only at the beginning of the first assimilation cycle (providing the correlation parameters remain time-invariant). In Fig. 6, the normalization factors were computed using the randomization method with a 100-member ensemble, which explains why the variance at the correlation point is not exactly equal to one.

# 6. CONCLUDING REMARKS

A practical algorithm for modelling a large class of 2D and 3D univariate correlations functions on the sphere has been described. The theoretical basis of the algorithm lies in considering the solution of a generalized diffusion equation formed by replacing the Laplacian operator in the classical diffusion equation by a polynomial in the Laplacian. The integral solution of this equation can be interpreted as a covariance operator on the sphere; the kernel of the covariance operator is an isotropic covariance function and has an explicit representation in terms of the zonal spherical harmonics. The fundamental parameters of the correlation model are the products of the weighting ('diffusion') coefficient for each polynomial term with the total integration time of the generalized diffusion equation. The practical implementation of the algorithm is iterative and forms a class of Laplacian-based grid-point filters, one example of which is the basis of the well-known Derber and Rosati scheme (Derber and Rosati 1989).

The most important features of the algorithm are summarized as follows:

• The algorithm is easily generalized to account for complex boundary domains and hence is particularly well suited for modelling correlation functions in ocean data assimilation applications.

• The shape (spectrum) of the correlation function can be controlled by adjusting the relative weights of the different terms in the generalized diffusion equation. In general, the more weight given to the higher-order Laplacian terms, the more oscillatory the shape of the correlation function and the greater (lesser) the variance at intermediate (low and high) wave numbers. For the simplest case in which a zero-weight is given to all but the first-order Laplacian term, the correlation function is well approximated by a Gaussian (Derber and Rosati 1989).

• The algorithm can be used to represent a class of 1D correlation functions in the vertical as well as 2D correlation functions over the sphere. This provides the basis of a general 3D correlation operator on the sphere.

• Geographical variations in the length-scale and shape of the correlation function can be accounted for by defining the diffusion coefficients to be a function of the model's spatial coordinates.

• The correlation functions can be made anisotropic by stretching and/or rotating the computational coordinate system via a 'diffusion' tensor. The diffusion tensor is analogous to the 'metric tensor' introduced by Purser (1986) for modelling correlation functions in atmospheric data analysis. In an ocean model, the rotation of the 'horizontal' correlation surfaces from geopotential to isopycnal coordinates was one particularly attractive possibility illustrated in this paper.

• The numerical cost of the algorithm is primarily determined by the minimum number of iterations needed to maintain the filter stability, which in turn depends on the specific parameter choices of the correlation model. For the geopotential-coordinate correlation model, the minimum number of iterations is roughly proportional to the

square of the ratio of the length-scale to the grid spacing. For the isopycnal-coordinate correlation model, the stability criterion is much more stringent, particularly in regions of steep isopycnal slopes (e.g. see appendix C in Griffies *et al.* (1998) for a detailed discussion on the stability of isopycnal diffusion models). In this case, many iterations are required, thereby limiting the practicality of the algorithm. Nevertheless, there may be considerable scope for improving the efficiency of the algorithm using either an alternative (e.g. semi-implicit) time discretization scheme or a multigrid approach (Brandt 1977). These possibilities should be explored in the future.

#### ACKNOWLEDGEMENT

This work was initiated at LODYC. Much of the practical implementation was carried out by the first author during a year spent at ECMWF, which was organized as part of a collaboration between the French MERCATOR project and the seasonal forecasting project at ECMWF. In particular, he would like to thank David Anderson, Jean-Claude André, Philippe Courtier, Pascale Delecluse and Anthony Hollingsworth for making these arrangements possible. David Anderson, Jérôme Vialard and Andrea Piacentini provided many useful suggestions for improving an early draft of the manuscript. We are also grateful to two anonymous reviewers for their constructive and insightful remarks, particularly with regard to the class of correlation functions presented in the paper.

# APPENDIX

# Matching the solution of the heat equation on the sphere to a Gaussian covariance operator

Consider the function

$$g(\theta;\gamma) = \sum_{n=0}^{\infty} g_n P_n^0(\cos\theta) = \sum_{n=0}^{\infty} (2n+1)^{1/2} \frac{I_{n+1/2}(\gamma)}{I_{1/2}(\gamma)} P_n^0(\cos\theta)$$
(A.1)

where

$$\gamma = \frac{a^2}{L^2},\tag{A.2}$$

*L* being a length-scale to be interpreted shortly.  $I_{n+1/2}(\gamma)$  is the modified Bessel function of fractional order n + 1/2 and argument  $\gamma$ . Equation (A.1) is a covariance function by virtue of the positivity of the  $I_{n+1/2}(\gamma)$  for  $\gamma > 0$ . Using the expansion (Abramowitz and Stegun 1964, Formula 10.2.36)

$$e^{\gamma \cos \theta} = \sum_{n=0}^{\infty} (2n+1)^{1/2} \left(\frac{\pi}{2\gamma}\right)^{1/2} I_{n+1/2}(\gamma) P_n^0(\cos \theta),$$
(A.3)

(A.1) can be written in the alternative form

$$g(\theta; \gamma) = \left(\frac{2\gamma}{\pi}\right)^{1/2} \frac{1}{I_{1/2}(\gamma)} e^{\gamma \cos \theta} = \frac{\gamma}{\sinh \gamma} e^{\gamma \cos \theta}.$$
 (A.4)

The length-scale of  $g(\theta; \gamma)$  may be defined following Daley (1991):

$$L^{2} = -2\frac{g(0; \gamma)}{\nabla^{2}g(0; \gamma)}.$$
 (A.5)

### A. WEAVER and P. COURTIER

The denominator in (A.5) can be evaluated by applying the Laplacian in spherical coordinates,

$$\nabla^2 = \frac{1}{a^2 \cos \phi} \frac{\partial}{\partial \phi} \left( \cos \phi \frac{\partial}{\partial \phi} \right) + \frac{1}{a^2 \cos^2 \phi} \frac{\partial^2}{\partial \lambda^2}, \tag{A.6}$$

to (A.4) with  $\theta$ , the angular separation between points  $(\lambda, \phi)$  and  $(\lambda', \phi')$ , defined according to the great circle distance formula

$$\cos \theta = \cos \phi \cos \phi' \cos(\lambda - \lambda') + \sin \phi \sin \phi'. \tag{A.7}$$

This leads to

$$L^2 = \frac{a^2}{\gamma},\tag{A.8}$$

and thus the interpretation of L as the length-scale of the covariance function represented by  $g(\theta; \gamma)$ . In terms of chordal distance r (22), (A.4) can be written as

$$g(\theta(r); \gamma) = \left(\frac{2\gamma}{\pi}\right)^{1/2} \frac{e^{\gamma}}{I_{1/2}(\gamma)} e^{-r^2/2L^2} = \frac{\gamma e^{\gamma}}{\sinh \gamma} e^{-r^2/2L^2}.$$
 (A.9)

Hartman and Watson (1974) remark that (30) may be closely matched to (A.1) by equating their coefficients on the  $P_1^0$  Legendre polynomial,

$$e^{-2\kappa T/a^2} = \frac{I_{3/2}(\gamma)}{I_{1/2}(\gamma)} = \frac{\cosh \gamma}{\sinh \gamma} - \frac{1}{\gamma},$$
 (A.10)

the coefficients for n = 0 being already equal to one. For a length-scale small compared to the radius of the earth ( $\gamma \gg 1$ ), a first order development of the above equality leads to

$$\kappa T \approx \frac{L^2}{2} \tag{A.11}$$

as in the 1D example. Furthermore, from (A.9) the variance at any point is

$$g(0; \gamma) = \frac{\gamma e^{\gamma}}{\sinh \gamma} \approx 2\gamma.$$
 (A.12)

The covariance operator in (29) can thus be approximated by

$$\eta(\lambda, \phi, T) \approx \frac{1}{2\pi L^2} \int_{\Sigma'} e^{-r^2/2L^2} \eta(\lambda', \phi', 0) \, \mathrm{d}\Sigma', \qquad (A.13)$$

with  $L^2$  given by (A.11).

The  $g_n$  in (A.1) can be evaluated from the modified spherical functions of the first kind using, for example, the IMSL library. Figure A.1(a) depicts the variance power spectrum of the correlation functions  $f(\theta; \kappa T)/f(0; \kappa T)$  and  $g(\theta; \gamma)/g(0; \gamma)$ , as well as their difference. Figure A.1(b) depicts their grid-point values for a length-scale of 500 km. A truncation at wave number 106 has been used. The agreement is, as expected, excellent, particularly for the large scales. The length-scale of the truncated function  $g(\theta; \gamma)$  is 500.00 km while it is 500.26 km for  $f(\theta; \kappa T)$ . The maximum difference in grid-point space is  $10^{-4}$ .



Figure A.1. (a) The variance power spectrum and (b) grid-point values of the isotropic correlation functions  $g(\theta; \gamma)/g(0; \gamma)$  (solid curve) and  $f(\theta; \kappa T)/f(0; \kappa T)$  (dashed curve) (see text). The curves are indistinguishable; the dotted curve is their difference.

# A. WEAVER and P. COURTIER

#### REFERENCES

Abramowitz, M. and Stegun, I.	1964	Handbook of mathematical functions. Dover Publications, Inc., New York
Andersson, E., Fisher, M., Munro, R. and McNally, A.	2000	Diagnosis of background errors for radiances and other observ- able quantities in a variational data assimilation scheme and the explanation of a case of poor convergence. Q. J. R. Meteorol. Soc., 126, 1455–1472
Arfken G	1966	Mathematical methods for physicists. Academic Press, New York
Barlow, R. J.	1989	Statistics: a guide to the use of statistical methods in the physical sciences John Wiley Chichester UK
Dehringer D. L. M. and	1009	An improved counted model for ENSO prediction and implica-
Leetma, A.	1990	tions for ocean initialization. Part I: The ocean data assimi- lation system. <i>Mon. Weather Rev.</i> , <b>126</b> , 1013–1021
Bell, M. J., Forbes, R. M. and Hines, A.	2000	Assessment of the FOAM global data assimilation system for real- time ocean forecasting. J. Marine Systems, 25, 1–22
Bennett, A. F.	1992	Inverse methods in physical oceanography. Cambridge University Press
Bennett, A. F., Chua, B. S. and Leslie, L. M.	1997	Generalized inversion of a global numerical weather prediction model, II: Analysis and implementation. <i>Meteorol. Atmos.</i> <i>Phys.</i> , <b>62</b> , 129–140
Brandt, A.	1 <b>977</b>	Multilevel adaptive solutions to boundary-value problems. <i>Math. Comp.</i> , <b>31</b> , 333–390
Carton, J. A., Chepurin, G., Cao, X. and Giese, B.	2000	A simple ocean data assimilation analysis of the global upper ocean 1950–95. Part I: Methodology. J. Phys. Oceanogr., 30, 294–309
Cohn, S. E., Da Silva, A., Guo, J., Sienkiewicz, M. and Lamich, D.	1997	Assessing the effects of data selection with the DAO Physical- space Statistical Analysis System. <i>Mon. Weather Rev.</i> , <b>126</b> , 2913–2926
Courtier, P.	1997	Dual formulation of four-dimensional variational assimilation. O. J. R. Meteorol. Soc., 123, 2449–2462
Courtier, P., Thépaut, JN. and Hollingsworth, A.	1994	A strategy for operational implementation of 4D-Var, using an incremental approach. Q. J. R. Meteorol. Soc., <b>120</b> , 1367–1388
Courtier, P., Andersson, E., Heckley, W., Pailleux, J., Vasiljević, D., Hamrud, M., Hollingsworth, A., Rabier, F.	1998	The ECMWF implementation of three-dimensional variational assimilation (3D-Var). Part I: Formulation. Q. J. R. Meteorol. Soc., 124, 1783–1808
Daley R	1991	Atmospheric data analysis. Cambridge atmospheric and space
Daley, R.	1724	sciences series. Cambridge University Press
Daley, R. and Barker, E.	1999	The NAVDAS Source Book. Naval Research Laboratory, Monterey California
Derber, J. and Bouttier, F.	1999	A reformulation of the background error covariance in the ECMWF global data assimilation system. <i>Tellus</i> , <b>51A</b> , 195–221
Derber, J. and Rosati, A.	1989	A global oceanic data assimilation system. J. Phys. Oceanogr., 19, 1333-1347
Egbert, G. D., Bennett, A. F. and Foreman, M. G. G.	1994	Topex/Poseidon tides estimated using a global inverse model. J. Geophys. Res., 99, 24821-24852
Fisher, M. and Courtier, P.	1995	'Estimating the covariance matrices of analysis and forecast error in variational data assimilation'. ECMWF Technical Memo. No. 220. European Centre for Medium-Range Weather Fore- casts, Reading, UK
Gaspari, G. and Cohn, S.	1999	Construction of correlation functions in two and three dimensions. Q. J. R. Meteorol. Soc., <b>125</b> , 723–757
Gauthier, P., Charette, C., Fillion, L., Koclas, P. and Laroche, S.	1999	Implementation of a 3D variational data assimilation system at the Canadian Meteorological Centre. Part I: the global analysis. <i>AtmosOcean</i> , <b>37</b> , 103-156
Gavart, M. and DeMey, P.	1997	Isopycnal EOFs in the Azores Current Region: a statistical tool for dynamical analysis and data assimilation. J. Phys. Oceanogr., 27, 2146–2157
Gneiting, T.	1999	Correlation functions for atmospheric data analysis. Q. J. R.

Correlation functions for atmospheric data analysis. Q. J. R. 1999 Meteorol. Soc., 125, 2449-2464

Griffies, S. M., Gnanadesikan, A., Pacanowski, R. C., Larichev, V. D., Dukowicz, J. K. and	1998	Isoneutral diffusion in a z-coordinate ocean model. J. Phys. Oceanogr., 28, 805–830
Cmith D D		
Hartman, P. and Watson, G. S.	1974	'Normal' distribution functions on spheres and the modified Bessel functions. Annals of Probability 2, 503-607
Hollingsworth, A. and Lönnberg, P.	1986	The statistical structure of short-range forecast errors as deter- mined from radiosonde data. Part I: The wind field. <i>Tellus</i> , <b>38A</b> , 111–136
Ide, K., Courtier, P., Ghil, M. and Lorenc, A. C.	1997	Unified notation for data assimilation: operational, sequential and variational. J. Meteorol. Soc. Jpn., 75, 181-189
Julian, P. R. and Thiébaux, H. J.	1975	On some properties of correlation functions used in optimum internolation schemes Mon Weather Rev 103 605-616
Lönnberg, P. and Hollingsworth, A.	1986	The statistical structure of short-range forecast errors as deter- mined from radiosonde data. Part II: The covariance of beight and wind errors. <i>Tellus</i> <b>38A</b> 137-161
Lorenc, A. C.	1988	Optimal nonlinear objective analysis. Q. J. R. Meteorol. Soc., 114, 205-240
	1992	Iterative analysis using covariance functions and filters. Q. J. R. Meteorol Soc. 118 569-591
	1997	Development of an operational variational assimilation scheme.
Lorenc, A. C., Ballard, S. P., Bell, R. S., Ingleby, N. B., Andrews, P. L. F., Barker, D. M., Bray, J. R.	2000	The Met Office global three-dimensional variational data assimi- lation scheme. Q. J. R. Meteorol. Soc., <b>126</b> , 2991–3012
Clayton, A. M., Dalby, T., Li, D., Payne, T. J. and Saunders, F. W.		
Madec, G. and Imbard, M.	1996	A global ocean mesh to overcome the North Pole singularity. Clim. Dvn., 12, 381-388
Madec, G., Delecluse, P., Imbard, M. and Levy, C.	1999	'OPA, release 8.1, Ocean General Circulation Model reference manual'. Internal report, LODYC/IPSL, France
Meyers, G., Phillips, H., Smith, N. and Sprintall, J.	1991	Space and time scales for optimal interpolation—tropical Pacific Ocean. Prog. Oceanogr., 28, 189-218
Pacanowski, R. C.	1996	MOM 2 documentation, user's guide and reference manual. GFDL Ocean Technical Report No. 3.1, Geophysical Fluid Dynamics Laboratory/NOAA
Parrish, D. F. and Derber, J. C.	1992	The National Meteorological Center's spectral statistical interpo- lation analysis system. Mon. Weather Rev. 120, 1747-1763
Parrish, D. F., Derber, J. C., Purser, R. J., Wu, WS. and Pu, ZX.	1997	The NCEP global analysis system: Recent improvements and future plans. J. Meteorol. Soc. Jpn., 75, 359-365
Purser, R. J.	1986	'Bayesian optimal analysis for meteorological data'. Pp. 167–172 in Proceedings of the international symposium on variational methods in geosciences, 15–17 October, 1985, Norman, Ok- lahoma, Ed. Y. Sasaki, Elsevier Amsterdam
Redi, H. R.	1982	Oceanic isopycnal mixing by coordinate rotation. J. Phys. Oceanicr. 12, 1154–1158
Rabier, F., McNally, A., Andersson, E., Courtier, P., Undén, P., Eyre, J., Hollingsworth, A. and Bouttier, F.	1998	The ECMWF implementation of three-dimensional variational assimilation (3D-Var). Part II: Structure functions. Q. J. R. Meteorol. Soc., <b>124</b> , 1809–1829
Rabier, F., Järvinen, H., Klinker, E., Mahfouf, JF. and Simmons, A.	2000	The ECMWF operational implementation of four-dimensional variational assimilation. Part I: Experimental results with simplified physics <i>O L R Material</i> Soc. <b>126</b> 1142 1170
Tarantola, A.	1987	Inverse problem theory: Methods for data fitting and model par- ameter estimation Elsevier Amsterdam
Thiébaux, H. J., Mitchell, H. L. and Shantz, D. W.	1986	Horizontal structure of hemispheric forecast error correlations for geopotential and temperature. Mon. Weather Rev. 114

Vialard, J. and Delecluse, P.

- Honzontal structure of hemispheric forecast error correlations for geopotential and temperature. *Mon. Weather Rev.*, 114, 1048–1066
  An OGCM study for the TOGA decade. Part I: Role of salinity in the physics of the western Pacific fresh pool. *J. Phys. Oceanogr.*, 28, 1071–1088

# A. WEAVER and P. COURTIER

1980

Weaver, A. T. and Vialard, J.

2000 'Development of an ocean incremental 4D-Var scheme for seasonal prediction'. Pp. 191–194 in Proceedings of the third WMO international symposium on assimilation of observations in meteorology and oceanography, 7–11 June 1999, Québec City, Canada

Weber, R. O. and Talkner, P. 1993

Wahba, G. and Wendelberger, J.

- Some remarks on spatial correlation function models. Mon. Weather Rev., 121, 2611-2617
- Some new mathematical methods for variational objective analysis using splines and cross-validation. *Mon. Weather Rev.*, **108**, 36–57

# Three- and Four-Dimensional Variational Assimilation with a General Circulation Model of the Tropical Pacific Ocean. Part I: Formulation, Internal Diagnostics, and Consistency Checks

A. T. WEAVER

Centre Européen de Recherche et de Formation Avancée en Calcul Scientifique/SUC URA 1875, Toulouse, France

#### J. VIALARD

European Centre for Medium-Range Weather Forecasts, Reading, United Kingdom, and Laboratoire d'Oceanographie Dynamique et de Climatologie/CNRS/IRD/UPMC/MNHN, Paris, France

#### D. L. T. ANDERSON

European Centre for Medium-Range Weather Forecasts, Reading, United Kingdom

(Manuscript received 20 May 2002, in final form 27 November 2002)

#### ABSTRACT

Three- and four-dimensional variational assimilation (3DVAR and 4DVAR) systems have been developed for the Océan Parallélisé (OPA) ocean general circulation model (OGCM) of the Laboratoire d'Océanographie Dynamique et de Climatologie. An iterative incremental approach is used to minimize a cost function that measures the statistically weighted squared differences between the observational information and their model equivalent. The control variable of the minimization problem is an increment to the background estimate of the model initial conditions at the beginning of each assimilation window. In 3DVAR, the increment is transported between observation times within the window using a persistence model, while in 4DVAR a dynamical model derived from the tangent linear (TL) of the OGCM is used. Both the persistence and TL models are shown to provide reasonably good descriptions of the evolution of typical errors over the 10- and 30-day widths of the assimilation windows used in the authors' 3DVAR and 4DVAR experiments, respectively.

The present system relies on a univariate formulation of the background-error covariance matrix. In practice, the background-error covariances are specified implicitly within a change of control variable designed to improve the conditioning of the minimization problem. Horizontal and vertical correlation functions are modeled using a filter based on a numerical integration of a diffusion equation. The background-error variances are geographically dependent and specified from the model climatology. Single observation experiments are presented to illustrate how the TL dynamics act to modify these variances in a flow-dependent way by diminishing their values in the mixed layer and by displacing the maximum value of the variance to the level of the background thermocline.

The 3DVAR and 4DVAR systems have been applied to a tropical Pacific version of OPA and cycled over the period 1993–98 using in situ temperature observations from the Global Temperature and Salinity Pilot Programme. The overall effect of the data assimilation is to reduce a large bias in the thermal field, which was present in the control. The fit to the data in 4DVAR is better than in 3DVAR, and within the specified observationerror standard deviation. Intermittent updating of the linearization state of the TL model is shown to be an important feature of the incremental 4DVAR algorithm and contributes significantly to improving the fit to the data.

#### 1. Introduction

The El Niño–Southern Oscillation (ENSO) phenomenon is one of the main contributors to predictability on seasonal to interannual timescales (Barnett et al. 1993; Palmer and Anderson 1994). Accurate ENSO forecasts are thus a prerequisite to the development of a reliable dynamical seasonal forecasting system with a coupled ocean–atmosphere model (CGCM). Early studies (Cane et al. 1986; Latif and Flügel 1991; Balmaseda et al. 1994) showed that some ENSO forecasting skill could be obtained using only observed surface wind stress to initialize the ocean component of the forecasting system. More recent studies have shown that assimilating ocean observations to initialize CGCMs resulted in significantly better ENSO forecasts (Ji and

*Corresponding author address:* Dr. Anthony T. Weaver, CERFACS, 42 Ave. Gaspard Coriolis, 31057 Toulouse Cedex 1, France. E-mail: weaver@cerfacs.fr

<sup>© 2003</sup> American Meteorological Society

JULY 2003

Leetma 1997; Rosati et al. 1997; Segschneider et al. 2000, 2001; Alves et al. 2002). Subsurface temperature observations from expendable bathythermographs (XBTs) and the Tropical Atmosphere Ocean (TAO) array (McPhaden 1993) were shown to be especially beneficial to forecast skill (Ji and Leetma 1997; Alves et al. 2002; Segschneider et al. 2001). Combining available observations and an ocean model into a dynamically consistent picture of the ocean state can also help to provide better insight into the processes determining ocean variability at various timescales.

In meteorology, optimal interpolation (OI) was the standard technique used for many years for producing initial conditions for numerical weather prediction (NWP; Gandin 1965; Rutherford 1972; Lorenc 1981). More recently, however, most of the major NWP centers have replaced their OI systems by variational assimilation systems (Parrish and Derber 1992; Courtier et al. 1998; Gauthier et al. 1999; Rabier et al. 2000; Lorenc et al. 2000). In variational assimilation, the analysis problem is defined by the minimization of a cost function that measures the statistically weighted squared differences between observations (including a model background state) and their model counterpart. The cost function is minimized with respect to selected control variables and this is done iteratively using a gradient descent method. Variational assimilation overcomes many of the limitations of OI: it allows for greater flexibility for assimilating different observation types (possibly nonlinearly related to the model state); it eliminates the need to split the analysis domain into subsections so that all observations can, in principle, influence the analysis at every model grid point; it provides a more general framework for using more sophisticated background-error covariance models; and it provides a clearer development path toward advanced, four-dimensional assimilation techniques. These advantages are equally relevant for oceanographic data assimilation.

In this paper, we describe three- and four-dimensional variational assimilation (3DVAR and 4DVAR) systems that have been developed for the rigid-lid version of the Océan Parallélisé (OPA) ocean general circulation model (OGCM) of the Laboratoire d'Océanographie Dynamique et de Climatologie (LODYC; Madec et al. 1998). One of the main motivations for developing the system is to produce ocean analyses for seasonal climate forecasting. In the present study, the 3DVAR and 4DVAR systems are applied to produce a reanalysis of the tropical Pacific Ocean over the period 1993-98 using in situ temperature observations. Here, and in the companion paper by Vialard et al. (2003, hereafter referred to as Part II), we evaluate the analyses by focusing on their statistical and physical properties, and their comparison with independent datasets, rather than their impact on climate forecasts.

Both the 3DVAR and 4DVAR systems have been designed following the incremental approach (Courtier et al. 1994). The control variable of the minimization

problem is taken to be an increment to the background estimate of the model initial conditions at the beginning of a given assimilation window. In the cost function, the observations are compared to the sum of the background counterpart of the observations and an increment computed in observation space using a linear model. The fundamental difference between the 3DVAR and 4DVAR formulations lies in the level of sophistication of the linear model used to transport the state increment between observation times. In 3DVAR, a simple persistence model is used, whereas in 4DVAR a dynamical model based on the tangent linear (TL) of the OPA OGCM is used. The 4DVAR scheme involves substantially more development than does 3DVAR since an adjoint model must be derived for the linearized version of the OGCM in order to compute the gradient of the cost function with respect to the increment at initial time.

This particular incremental version of 3DVAR can be viewed as a limiting case of incremental 4DVAR in which the TL operator is replaced by the identity matrix. As in 4DVAR, observations can be assimilated at their appropriate measurement times since they are compared directly to the background state, which is propagated in time using the OGCM. For this reason, the scheme has been coined 3D-FGAT, for first guess at appropriate time (Fisher and Andersson 2001).<sup>1</sup> For example, the FGAT feature may be particularly important in the Tropics where an equatorial Kelvin wave can travel more than 2000 km in 10 days.<sup>2</sup>

In incremental 4DVAR, the dynamical model used to propagate the increment provides a time-dependent multivariate constraint on the analysis. Incremental 4DVAR is derived as an approximation to the complete 4DVAR problem in which the full nonlinear model (here an OGCM) is imposed as a constraint in the cost function with the model initial conditions taken as the control variables (Le Dimet and Talagrand 1986; Talagrand and Courtier 1987). In oceanography, most 4DVAR-related applications to date with OGCMs have concentrated on solving the complete problem directly, and in some cases using different or additional control variables (e.g., surface forcing fields; Tzipermann et al. 1992a,b; Greiner et al. 1998a,b, Greiner and Arnault 2000; Bonekamp et al. 2001). The incremental formulation was introduced in meteorology to overcome some important practical difficulties with solving the complete 4DVAR problem directly. In the latter, nonlinearities in the model constraint can significantly complicate the structure of the cost function and prevent its minimization at a

<sup>&</sup>lt;sup>1</sup>FGAT was initially introduced in OI in the mid-1980s by D. Vasiljevic at the European Centre for Medium Range Weather Forecasts (ECMWF). The 40-yr atmospheric reanalysis project (ERA-40) at ECMWF employs an FGAT version of 3DVAR.

<sup>&</sup>lt;sup>2</sup> Ten days is a typical window width used in OI-type ocean analysis systems such as the operational system at ECMWF (Alves et al. 2002). It is also the window width used for the 3DVAR experiments presented here.

reasonable computational cost using a gradient descent method. Furthermore, in order to compute a numerically accurate gradient, the adjoint of the *exact* TL of the full nonlinear model is required. When the nonlinear model contains discontinuous parameterizations or numerics, an accurate derivation of these models can be particularly difficult (Xu 1996).

The incremental algorithm should be viewed as a practical algorithm for approximately solving the complete problem. As constraints are linear, the cost function is quadratic and minimization with a gradient descent method is generally much more efficient. In the present study, our main objective is to reconstruct the largescale, low-frequency component of the tropical ocean circulation, which is well known to be largely governed by linear wave dynamics (e.g., see Philander 1989). Therefore, a priori the linearity assumption in incremental 4DVAR would not appear to be a very restrictive one. More generally, however, nonlinear effects can be partly accounted for by introducing a feedback (outer loop) in the algorithm to update the basic state of the linear model with increments generated during minimization (the inner loop) (Courtier et al. 1994; Laroche and Gauthier 1998). In terms of technical development, the incremental formulation has a distinct advantage over the complete formulation as the derivation of the linear and adjoint models can be greatly simplified by smoothing or neglecting discontinuous parameterizations (Mahfouf 1999). For example, this latter point has been exploited in the present study to neglect changes in vertical diffusion coefficients associated with perturbations in temperature, salinity, and velocity. Finally, the computational cost of 4DVAR may also be significantly reduced by computing the increments at lower resolution than that of the full model (Rabier et al. 2000), although this is not an issue in our present system, which employs a relatively low resolution version of the OGCM.

The purpose of this paper is to give a thorough description of the current 3DVAR and 4DVAR systems and to study certain algorithmic and statistical aspects of the two systems. Validation of physical aspects of the analyses is given in Part II. The organization of the paper is as follows. Section 2 describes the general formulation of the incremental 3DVAR and 4DVAR problems. In section 3, the different system components are described. In section 4, several diagnostics are presented to highlight some important properties of the two systems. A summary is given in section 5.

# 2. Formulation of the 3DVAR and 4DVAR problems

#### a. Incremental formulation

The notation used in this paper closely follows the recommendations of Ide et al. (1997). Let  $\mathbf{w}$  denote the

ocean state vector. The components of **w** consist of those model variables that are to be estimated from observations to produce the analysis  $\mathbf{w}^a$ . A model forecast, initiated from a previous analysis, provides a prior or background estimate,  $\mathbf{w}^b$ , of **w**. Since  $\mathbf{w}^b$  will already be close in some sense to the "true" state we wish to estimate, it is convenient to formulate the estimation problem in terms of an increment,  $\delta \mathbf{w}$ , where

$$\mathbf{w} = \mathbf{w}^b + \delta \mathbf{w}. \tag{1}$$

The state vector is propagated in time by the ocean model:

$$\mathbf{w}(t_i) = M(t_i, t_{i-1})[\mathbf{w}(t_{i-1})], \qquad (2)$$

where  $M = M(t_i, t_{i-1})$  represents the nonlinear model operator acting on  $\mathbf{w}(t_{i-1})$  between times  $t_{i-1}$  and  $t_i$ . Substituting (1) into (2) and expanding about  $\mathbf{w}^b(t_{i-1})$ gives, to first order,

$$\mathbf{w}(t_i) \approx M(t_i, t_{i-1})[\mathbf{w}^b(t_{i-1})] + \mathbf{M}(t_i, t_{i-1})\delta \mathbf{w}(t_{i-1}), \quad (3)$$

where  $\mathbf{M} = \mathbf{M}(t_i, t_{i-1})$  denotes a linear operator that acts on  $\delta \mathbf{w}(t_{i-1})$  between times  $t_{i-1}$  and  $t_i$ . We define the prognostic model for the increment as

$$\delta \mathbf{w}(t_i) = \mathbf{M}(t_i, t_{i-1}) \delta \mathbf{w}(t_{i-1}).$$
(4)

Now, let  $\mathbf{y}_i^o$  denote the observation vector at time  $t_i$ . Denoting  $H_i$  as the observation operator at  $t_i$ , then the model equivalent of  $\mathbf{y}_i^o$  can be written as

$$H_i[\mathbf{w}(t_i)] \approx H_i[\mathbf{w}^b(t_i)] + \mathbf{H}_i \delta \mathbf{w}(t_i), \tag{5}$$

where  $\mathbf{H}_i$  is a linear operator that acts on the increment at  $t_i$ . Assuming that a time sequence of observation vectors is available over an interval  $t_0 \leq t_i \leq t_n$ , then the model estimates of the observations within this interval can be directly related to the model initial conditions since  $\mathbf{w}(t_i) = M(t_i, t_0)[\mathbf{w}(t_0)]$ , where  $M(t_i, t_0) \equiv$  $M(t_i, t_{i-1})^\circ \cdots \circ M(t_1, t_0)$ . Thus,

$$H_i[\mathbf{w}(t_i)] = G_i[\mathbf{w}(t_0)] \approx G_i[\mathbf{w}^b(t_0)] + \mathbf{G}_i \delta \mathbf{w}(t_0), \quad (6)$$

where the combined operator  $G_i = H_i M(t_i, t_0)$  is a generalized observation operator and  $\mathbf{G}_i = \mathbf{H}_i \mathbf{M}(t_i, t_0)$  is its linearized counterpart with  $\mathbf{M}(t_i, t_0) \equiv \mathbf{M}(t_i, t_{i-1}) \cdots \mathbf{M}(t_1, t_0)$ .

In 4D variational assimilation, the analysis is defined as the state vector  $\mathbf{w}^a = \mathbf{w}^a(t_0)$  that simultaneously minimizes the "distance" to the background state  $\mathbf{w}^b =$  $\mathbf{w}^b(t_0)$  and to the time sequence of observations  $\mathbf{y}_i^o$  on  $t_0 \leq t_i \leq t_n$ . Distance is defined by an inner product (a cost function  $J^F$ ) whose weighting metric takes into account the statistical accuracy of the background and observational information. Expressed as a function of the increment  $\delta \mathbf{w} = \delta \mathbf{w}(t_0)$ , which constitutes the control vector of the minimization problem,  $J^F$  may be written as JULY 2003

$$J^{F}(\delta \mathbf{w}) = \frac{1}{2} \delta \mathbf{w}^{\mathrm{T}} \mathbf{B}^{-1} \delta \mathbf{w} + \frac{1}{2} \sum_{i=0}^{n} [G_{i}(\mathbf{w}^{b} + \delta \mathbf{w}) - \mathbf{y}_{i}^{o}]^{\mathrm{T}} \mathbf{R}_{i}^{-1} [G_{i}(\mathbf{w}^{b} + \delta \mathbf{w}) - \mathbf{y}_{i}^{o}],$$
(7)

where the matrices **B** and **R**<sub>i</sub> contain estimates of the covariances of background and observation error, respectively. In (7) the observation errors are assumed to be uncorrelated in time and uncorrelated with the background error. The observation term  $(J_o^F)$  measures the fit between the observations and their model equivalent. The background term  $(J_o^F)$  penalizes the size of the increment (i.e., measures the fit to the background state). The analysis is given by  $\mathbf{w}^a = \mathbf{w}^b + \delta \mathbf{w}^a$ , where  $\delta \mathbf{w}^a$  is the increment that minimizes  $J^F$ .

In the *incremental* formulation of variational assimilation (Courtier et al. 1994), (7) is approximated by a *quadratic* cost function J of  $\delta \mathbf{w}$  by replacing  $G_i(\mathbf{w}^b + \delta \mathbf{w})$  with its linearized counterpart (6). This results in an important practical simplification to the minimization problem, from one with potentially many minima due to the nonlinearity in  $G_i$ , to one with a unique minimum as guaranteed by the linearity of  $\mathbf{G}_i$ . The simplified cost function reads

$$J(\delta \mathbf{w}) = \frac{1}{2} \delta \mathbf{w}^{\mathrm{T}} \mathbf{B}^{-1} \delta \mathbf{w}$$
$$+ \frac{1}{2} \sum_{i=0}^{n} (\mathbf{G}_{i} \delta \mathbf{w} - \mathbf{d}_{i})^{\mathrm{T}} \mathbf{R}_{i}^{-1} (\mathbf{G}_{i} \delta \mathbf{w} - \mathbf{d}_{i}), \quad (8)$$

where the innovation vector

$$\mathbf{d}_i = \mathbf{y}_i^o - G_i(\mathbf{w}^b) = \mathbf{y}_i^o - H_i[\mathbf{w}^b(t_i)]$$
(9)

plays the role of an *effective* observation vector. Our 3DVAR (FGAT) and 4DVAR formulations differ principally in the choice of the linear operator **M** that is used in **G**<sub>i</sub> to propagate the increment  $\delta \mathbf{w}(t_i)$  in (4). In 3DVAR, the increment is evolved by a simple persistence model, which corresponds to setting  $\mathbf{M} = \mathbf{I}$ , the identity matrix. In 4DVAR, the increment is propagated by an approximate TL operator,  $\mathbf{M} \approx (\partial M/\partial \mathbf{w})|_{\mathbf{w}=\mathbf{w}^b}$ , the main approximation being introduced in the parametrization of vertical mixing as discussed in section 3b.

In practice, the cost function (8) is minimized approximately using an iterative gradient descent method. The increment is updated on each iteration using the gradient of the cost function with respect to the increment. The gradient of the  $J_o$  term with respect to the increment ( $\nabla_{\delta w} J_o$ ) can be obtained efficiently using the adjoint of the linear operator  $G_i$  (Le Dimet and Talagrand 1986; Thacker and Long 1988). A feedback between the linear and nonlinear models can then be in-

troduced by allowing the basic-state trajectory of the linear model to be regularly updated with the most recent estimate of the state trajectory obtained during minimization (Courtier et al. 1994). The updates are performed on an outer loop of the assimilation algorithm, while the iterations of the actual minimization are performed within an *inner* loop [see appendix A in Weaver et al. (2002) for details on how the inner-outer loop algorithm is implemented in our current system]. With frequent updates, the accuracy of the linear model should improve and the minimum of the incremental cost function should be closer to that of the original (nonincremental) cost function (7) involving the nonlinear model (Laroche and Gauthier 1998). This provides a practical way of accounting for nonlinearities in the assimilation algorithm while retaining the computational advantages of a quadratic minimization problem.

In order to improve the convergence properties of the minimization, a preconditioning transformation is employed by which the cost function is redefined in terms of a nondimensional variable:

$$\mathbf{v} = \mathbf{U}^{-\iota} \delta \mathbf{w},\tag{10}$$

where **U** is a rectangular matrix defined such that  $\mathbf{B} = \mathbf{U}\mathbf{U}^{\mathsf{T}}$  and  $\mathbf{B}^{-1} = (\mathbf{U}^{-1})^{\mathsf{T}}\mathbf{U}^{-1}$ , the superscript -I denoting generalized right inverse [i.e.,  $\mathbf{U}^{-1} = \mathbf{U}^{\mathsf{T}}(\mathbf{U}\mathbf{U}^{\mathsf{T}})^{-1}$ ]. The fact that **U** is rectangular, with dim $(\delta \mathbf{w}) < \dim(\mathbf{v})$ , is a feature that is particular to our current system and is related to the *rigid-lid* constraint used in the formulation of the ocean model. This point is discussed in more detail in section 3a and 3c. Introducing the transformation (10) directly into (8) leads to a simplified background term,  $J_b = \mathbf{v}^{\mathsf{T}}\mathbf{v}/2$ . The variational analysis problem is solved directly in **v** space, and then transformed back to model space using the generalized left inverse transformation:

$$\delta \mathbf{w} = \mathbf{U}\mathbf{v}.\tag{11}$$

The adjoint of (11) is used for computing the gradient of the  $J_o$  term in v space:

$$\boldsymbol{\nabla}_{\mathbf{v}} \boldsymbol{J}_o = \boldsymbol{\mathsf{U}}^{\mathrm{T}} \boldsymbol{\nabla}_{\delta \mathbf{w}} \boldsymbol{J}_o. \tag{12}$$

The conditioning of the  $J_b$  term in **v** space is optimal in the sense that the Hessian of  $J_b$  (the matrix of second derivatives of  $J_b$ ) is the identity matrix. For the special case of a single observation, the convergence of the minimization in **v** space is achieved in a single iteration using a gradient descent method with an exact line search.

The minimization routine used in this study is the

limited-memory quasi-Newton algorithm M1QN3 of Gilbert and Lemaréchal (1989). An exact line search has been employed with M1QN3 in order to improve the efficiency of the algorithm for quadratic minimization problems. The so-called warm start facility of M1QN3 is also employed when more than one outer iteration is performed in order to precondition the minimization using the information accumulated on the Hessian matrix during the preceding minimization.

#### b. The nonlinear analysis trajectory

The minimizing solution of the quadratic cost function (8) is a trajectory of increments  $\delta \mathbf{w}^a(t_i)$  satisfying exactly the equations of the linear model (4) on  $t_0 \leq t_i$  $\leq t_n$ . The corresponding model trajectory  $\mathbf{w}^b(t_i) + \delta \mathbf{w}^a(t_i)$  is used in the  $J_o$  term to compare with the observations [Eq. (8)]. It will be convenient to refer to this trajectory as the *linear* analysis in order to distinguish it from the *nonlinear* "analysis" trajectory  $\mathbf{w}^a(t_i)$ defined below. We will return to this point in section 4f and Part II of the paper.

There are several ways the analysis increment  $\delta w^a$  may be used to correct the trajectory of the nonlinear model. Two different approaches have been adopted in our 3DVAR and 4DVAR systems. Since our current assimilation system incorporates only temperature observations and relies on a univariate formulation of the background-error covariance matrix (section 3c), the 3DVAR produces an analysis increment for the temperature field only. A practical way of adjusting the nonanalyzed model fields while minimizing spurious adjustment processes is to apply the analysis increment gradually through a forcing term in the nonlinear model:

$$\mathbf{w}^{a}(t_{i}) = M(t_{i}, t_{i-1})[\mathbf{w}^{a}(t_{i-1})] + F(t_{i})\delta\mathbf{w}^{a}, \quad (13)$$

where  $\mathbf{w}^{a}(t_{0}) = \mathbf{w}^{b}$  and  $F(t_{i})$  is a weighting function defined such that  $\sum_{i=0}^{n} F(t_{i}) = 1$  so as to conserve the time-integrated value of the analysis increment  $\delta \mathbf{w}^{a}$ . The forcing term can be shown to behave as a low-pass time filter (Bloom et al. 1996). In this study, the temperature increment is applied uniformly over the time window via a constant forcing  $F(t_{i}) = 1/n$ . A similar procedure has been adopted in the ECMWF ocean analysis system (Alves et al. 2002) and is illustrated schematically in Fig. 1a.

In contrast to the 3DVAR, the 4DVAR produces a multivariate analysis increment since the TL model dynamics act to couple the different increment variables. This dynamical coupling thus allows us to generate increments in velocity and salinity even if only temperature data are assimilated. Since these increments will be in approximate dynamical balance, in 4DVAR we choose to initialize the nonlinear model directly using the analysis at  $t_0$ :

$$\mathbf{w}^{a}(t_{i}) = M(t_{i}, t_{0})[\mathbf{w}^{a}(t_{0})], \qquad (14)$$

where  $\mathbf{w}^{a}(t_{0}) = \mathbf{w}^{b} + \delta \mathbf{w}^{a}$  (Fig. 1b). Finally, when cy-

# a) Cycling of 3D-Var



FIG. 1. Schematic representation of the cycling procedures used for (a) 3DVAR and (b) 4DVAR. Two cycles are illustrated in each sketch. The dotted (solid) curves correspond to the background (analysis) trajectory; the cross symbols denote observations. The shaded square (circle) at the beginning of each cycle denotes the background (analysed) initial state. Note that in 4DVAR the analysis increment (represented by the difference between the shaded square and circle) is applied directly to the background initial state to produce the analysis trajectory, whereas in 3DVAR it is applied gradually as a forcing to the model equations. This explains why in 3DVAR the analysis and background start from the same point at the beginning of each cycle.

cling the 3DVAR or 4DVAR over an extended period, the analysis obtained from the trajectory at the end of the interval is taken to be the background state for a variational analysis performed on the following interval (Fig. 1).

#### 3. Components of the assimilation system

#### a. The ocean model

The ocean model used in this study is the OPA OGCM of the Laboratoire d'Océanographie Dynamique et de Climatologie (Madec et al. 1998). The model solves the primitive equations for horizontal currents  $\mathbf{u} = (u, v)$ , potential temperature  $T_{\theta}$ , and salinity S. The equations are formulated in orthogonal curvilinear z coordinates and discretized using finite differences on an Arakawa C grid. The basic configuration of the model is described in Vialard et al (2001). It covers the tropical Pacific Ocean from 30°N to 30°S, and 120°E to 70°W. The zonal resolution is 1° and the meridional resolution varies from  $0.5^{\circ}$  at the equator to  $2^{\circ}$  at the northern and southern boundaries. The model has 25 levels, with a resolution of 10 m in the upper 130 m, increasing to 1000 m in the bottom level. Realistic bathymetry is included using the Levitus (1982) land-sea mask.

At solid boundaries, conditions of no-slip and nonormal flux are applied on the velocity and tracer fields, respectively. At the ocean surface (z = 0), a rigid-lid and no-volume flux condition is applied (Roullet and Madec 2000). An obvious consequence of the rigid-lid condition is that, unlike in a free-surface model, there is no prognostic equation for sea surface height ( $\eta =$  JULY 2003

0). The no-volume flux condition is an additional requirement that the vertical velocity identically vanishes at the surface (w = 0 at z = 0) and, hence from the continuity equation, that the horizontal velocity field (u, v) is nondivergent (Bryan 1969). The latter condition is satisfied by defining the (u, v) field through a barotropic streamfunction  $\psi$  and a set of independent (subsurface) baroclinic velocities ( $\hat{u}$ ,  $\hat{v}$ ). This leads to some important technical difficulties in formulating the 4DVAR problem as detailed in appendix B of Weaver et al. (2002).

Surface fluxes of momentum, heat, and freshwater are prescribed at the ocean-atmosphere interface. The momentum flux is specified through weekly wind stress products from the European Remote Sensing (ERS) satellite's scatterometer (Grima et al. 1999). The heat and freshwater fluxes are specified as a daily climatology computed from the ECMWF (ERA-15) reanalysis (Gibson et al. 1997). The solar and nonsolar components of the heat flux are specified separately in order to allow penetration of the shortwave radiation in the upper ocean. A relaxation to weekly analyses of sea surface temperature (SST; Reynolds and Smith 1994) is applied through a Newtonian damping term added to the surface (nonsolar) heat flux. The relaxation coefficient is set to -40 W m<sup>-2</sup> K<sup>-1</sup>, which for a depth scale of 50 m corresponds to a restoring timescale of 1 month. No relaxation is applied to sea surface salinity. In the control integration only, in which no data are assimilated, a damping to Levitus climatological temperature and salinity is applied below the surface mixed layer outside the 10°S-10°N band.

#### b. The tangent-linear and adjoint models

The numerical codes of the TL and adjoint models have been derived directly from the numerical code of the nonlinear model by applying standard, hand-coding techniques (Talagrand 1991; Giering and Kaminski 1998). Some approximations have been introduced in the derivation of the TL and adjoint models, the most important one being in the parameterization of vertical diffusion, as described in more detail below. Another approximation relies on an intermittent storage of the basic-state trajectory to reduce computer memory requirements. In our experiments, the basic state has been stored once per day (every 16 time steps) and defined at intermediate times through linear interpolation. The impact of this approximation on the accuracy of the TL model was minor. Finally, to reduce the CPU cost of the adjoint integration, an approximation has been introduced in the adjoint of the elliptic solver that is applied on each time step to enforce the nondivergence constraint on the vertically integrated velocity. Details can be found in appendix B of Weaver et al. (2002).

#### SIMPLIFIED VERTICAL DIFFUSION

Let  $\alpha$  denote one of the prognostic state variables u, v,  $T_{\theta}$ , or S. In the complete TL model, the tendency of a perturbation  $\delta \alpha$  produced by vertical diffusion is described by the two-term partial differential equation:

$$\frac{\partial \delta \alpha}{\partial t} = \frac{\partial}{\partial z} \left( A_{\nu} \frac{\partial \delta \alpha}{\partial z} \right) + \frac{\partial}{\partial z} \left( \delta A_{\nu} \frac{\partial \alpha}{\partial z} \right), \quad (15)$$

where  $A_v = A_v(u, v, T_{\theta}, S)$  is the vertical diffusion coefficient and  $\delta A_{v} = (\partial A_{v}/\partial u)\delta u + (\partial A_{v}/\partial v)\delta v + (\partial A_{v}/\partial v)\delta v$  $\partial T_{\theta} \delta T_{\theta} + (\partial A_{\nu} \partial S) \delta S$  is the perturbation of  $A_{\nu}$  resulting from perturbations of u, v,  $T_{\theta}$ , and S. The first term on the right-hand side of (15) is a standard diffusion operator with coefficient  $A_v$ . In the incremental algorithm, A<sub>w</sub> is updated on each outer iteration and held constant only during the inner iterations. The second term on the right-hand side of (15) is more problematic because of both theoretical and practical difficulties in computing the perturbation  $\delta A_{v}$ . First, a direct computation of  $\delta A_{v}$ would involve linearizing the turbulent kinetic energy (TKE) and enhanced vertical diffusion parameterization schemes used in OPA. These schemes are highly nonlinear and discontinuous so any attempt to linearize them directly would have to be done with considerable caution (Xu 1996; Zou 1997). Discontinuities are a wellknown source of numerical noise in the TL model. One possibility for reducing such noise is to derive  $\delta A_{\mu}$  from a suitably smooth approximation to the original nonlinear scheme (Janiskova et al. 1999). Generally speaking, however, the problem is more complex than that of simply smoothing isolated discontinuities as the second term in (15) contains other generating mechanisms of spurious noise, which are present even in smoother parameterization schemes (Mahfouf 1999; Zhu and Kamachi 2000). The simplest way of avoiding these potential noise problems is to neglect the second term altogether by setting  $\delta A_{\mu} \equiv 0$ . This simplification, which has been used extensively in meteorology (e.g., Mahfouf 1999), is adopted here. Implementing a more elaborate linear physics parameterization was considered premature at this stage without first evaluating results using the simplified scheme.

#### c. The background-error covariance matrix

The background-error covariance matrix (**B**) plays an important role in determining the spatial structure of the analysis increment in both 3DVAR and 4DVAR. There are two basic difficulties in specifying **B**. First, given the sparsity of ocean observations, it is difficult, if not impossible, to obtain complete and accurate estimates of the covariances, even in the tropical Pacific, which is one of the best-observed ocean basins. Second, even if there were sufficient observations, the sheer size of **B** (roughly  $5 \times 10^{11}$  elements in our application) means that this matrix cannot be stored explicitly and so must be simplified. In practice, this is done by *modeling* the



FIG. 2. The horizontal correlation functions for the tracer fields at latitudes of (a) 0°, (b) 5°N, and (c) 15°N. The contour interval is 0.1.

covariances using analytical functions or filters that depend on a limited number of tunable parameters and that are numerically efficient for large-scale problems.

In our current system, **B** is univariate (three-block diagonal with respect to horizontal velocity, temperature, and salinity) and includes a relatively simple correlation model. Here, **B** is constructed as a symmetric product of several operators:

$$\mathbf{B} = \mathbf{S} \underline{\Sigma} \mathbf{\Lambda} \mathbf{L}^{1/2} \mathbf{W}^{-1} \mathbf{L}^{T/2} \mathbf{\Lambda} \underline{\Sigma} \mathbf{S}^{\mathrm{T}}$$
(16)

$$= (\mathbf{S} \boldsymbol{\Sigma} \boldsymbol{\Lambda} \mathbf{L}^{1/2} \mathbf{W}^{-1/2}) (\mathbf{W}^{-1/2} \mathbf{L}^{T/2} \boldsymbol{\Lambda} \boldsymbol{\Sigma} \mathbf{S}^{T}), \quad (17)$$

where **L** is a 3D filtering operator that is *self-adjoint* with respect to the scalar product whose metric is the diagonal matrix **W** of volume elements;  $\Lambda$  and  $\Sigma$  are diagonal matrices of normalization factors and background-error standard deviations, respectively; and **S** is a simplification operator that maps the horizontal components of total velocity (*u*, *v*) into *independent* components ( $\psi$ ,  $\hat{u}$ ,  $\hat{v}$ ) [for a detailed discussion see appendix B in Weaver et al. (2002)].

The underbraces in (17) highlight the preconditioning matrix  $\mathbf{U}$  and its adjoint  $\mathbf{U}^{T}$ , which are needed in (11) and (12), respectively. Here,  $\mathbf{U}$  is a rectangular matrix since the factor **S** is a mapping from a higher-dimensional space spanned by (u, v) into a lower-dimensional space spanned by  $(\psi, \hat{u}, \hat{v})$ . This explains why the generalized right inverse  $U^{-1}$  has been used in the transformation (10). The underbrace in (16) highlights that part of **B** corresponding to the correlation matrix **C**. The univariate correlations in **C** are assumed to be approximately Gaussian and are modeled implicitly with the filter L. The vertical (v) and horizontal (h) correlations are modeled separately using a 1D filter  $L_{v}$  and 2D filter  $L_{h}$ . The 3D correlation model is then constructed from the product  $\mathbf{L} = \mathbf{L}_{\nu} \mathbf{L}_{h}$ . The diagonal normalization matrix  $\Lambda$  is needed to ensure that the variances (diagonal elements) of **C** are equal to unity. Various filters exist for modeling correlation functions but some are better suited than others depending on the application. The complex boundaries associated with coastlines imposes

a particular constraint for oceanographic applications. For such applications, Laplacian- or diffusion-based filters are particularly well suited (Derber and Rosati 1989; Egbert et al. 1994; Weaver and Courtier 2001). Here, a diffusion-based filter has been used to model both the vertical and horizontal correlations;  $L_v$  is defined by an explicit time step integration of a 1D diffusion equation in the vertical direction, while  $L_h$  is defined by an explicit time step integration of a 2D diffusion equation over the sphere. The boundary conditions are chosen to be of Neumann type and are imposed directly within the finite-difference expression for the Laplacian using a land–ocean mask (Madec et al. 1998).

The correlation functions are made anisotropic and varied geographically by introducing a "diffusion" tensor in the Laplacian operator (Weaver and Courtier 2001). In particular, the tensor coefficients have been tuned to allow for longer correlation length scales near the equator in the zonal direction than in the meridional direction (Meyers et al. 1991; Kessler et al. 1996). Here, the horizontal length scales are taken to be a function of latitude and symmetric about the equator. The zonal and meridional length scales for the tracer fields have been defined to be  $8^{\circ}$  and  $2^{\circ}$ , respectively, at the equator, and 4° in both directions poleward of 20°N/S, with a linear transition between these values within the equatorial strip. The anisotropic ratio at the equator is consistent with the climatological observation statistics of Meyers et al. (1991), although the values of the actual length scales are somewhat smaller in our study. The values chosen here are broadly similar to those used in previous ocean data assimilation studies of the tropical Pacific (Smith et al. 1991; Behringer et al. 1998; Segschneider et al. 2001). The horizontal correlation functions for the tracer fields are illustrated in Fig. 2.

The horizontal correlations of the velocity field are modeled using a diffusion equation formed from the *vector* Laplacian rather than the scalar Laplacian used in the tracer diffusion model. The vector Laplacian has the desirable property of ensuring that the smoothing by the correlation function acts separately on the horizontal divergence and vorticity components of the velocity field (Madec et al. 1998). With the introduction of the anisotropic tensor, however, this is no longer



FIG. 3. The vertical correlation functions for the tracer fields at depths of (a) 96 (level 10), (b) 168 (level 15), and (c) 490 m (level 19). Note the different scales on the vertical axes.

strictly guaranteed. Geostrophy allows us to approximate the length scales of the vorticity correlations as  $L_h^{\xi} \approx 0.6L_h$ , where  $L_h$  is the specified length scale of the tracer correlations (see Weaver et al. 2002). These smaller length scales have also been used for the divergence correlations. It is worth remarking that, while geostrophy has been used as a constraint for defining the length scales of the tracer and velocity correlation functions, our covariance model does not at present include a geostrophic balance constraint on the covariances themselves. Note also that the correlation model for the velocity background errors is not required in our univariate 3DVAR, which assimilates temperature information only.

The vertical correlations of the velocity field are modeled using a diffusion equation that acts separately on the components u and v. The vertical length scales for both the tracers and the velocity components are taken to be a function of depth with a dependence on the model's vertical resolution to provide adequate smoothing between model levels. At each model level, the vertical length scale is set to twice the depth of that level. This results in rather sharp correlations above the thermocline where the resolution is highest and much broader correlations below the thermocline where the resolution is coarsest (Fig. 3).

The background-error standard deviations are allowed to vary with each grid point and have been computed with respect to the climatological model mean obtained from a control run without data assimilation. This specification is based on the assumption that background errors are likely to be largest in regions of strong ocean variability (e.g., in the thermocline). The same background-error standard deviations are used at the beginning of each assimilation cycle.

#### d. Observations

The assimilation dataset consists of in situ temperature observations from the Global Temperature and Salinity Pilot Project (GTSPP) of the National Oceanographic Data Centre (NODC). This includes data from mainly TAO moorings and XBTs, and from a limited number of conductivity-temperature-depth (CTD) casts and drifting buoys. A manual quality control procedure was used to remove suspect data (Alves et al. 2002). Observations falling within the surface level of the model (between 0 and 10 m) were also discarded to avoid potential redundancy with the observed SST (Reynolds and Smith 1994) used in the Newtonian damping term during outer iterations. The in situ temperatures retained for assimilation were then converted into potential temperature (the prognostic model variable) using a standard conversion formula [Eq. (A3.13) in Gill (1982)] with a reference salinity of 35.0 psu.

The observation-error covariances are assumed to be uncorrelated in space and time. The error variances are set to  $(0.5^{\circ}C)^2$  for TAO data and  $(1.0^{\circ}C)^2$  for all other data to ensure that the generally higher quality TAO data have more weight in the analysis. The observation operator *H* consists of horizontal bilinear interpolation and vertical linear interpolation. For TAO data, *H* also includes a time averaging since these data are available as daily averages.

# 4. Evaluation of 3DVAR and 4DVAR: Internal diagnostics and consistency checks

The purpose of this section is to evaluate certain algorithmic and statistical properties of the 3DVAR and 4DVAR systems. The analyses themselves will be discussed in detail in Part II. First, in section 4a, the assimilation experiments are introduced. In section 4b, the validity of the linear assumption, which underlies the incremental formulation, is investigated. Convergence and optimality properties of the assimilation systems are then examined in sections 4c and 4d. In section 4e, some of the flow-dependent characteristics of the backgrounderror statistics used implicitly in 4DVAR are examined in a simplified framework with single observations. Finally, in section 4f, the statistics of the innovation and residual vectors are examined to assess the fit of the background and of the analysis to the assimilated observations.

#### a. Experimental setup

Both the 3DVAR and 4DVAR systems have been cycled over the period 1 January 1993-30 December 1998 using a 10- and 30-day assimilation window, respectively. A total of 60 (inner) iterations of the minimization were performed per cycle, with an outer iteration performed every 10 inner iterations in 4DVAR. No outer iterations were performed in 3DVAR since  $H_i$ is linear and  $\mathbf{M} = \mathbf{I}$  is independent of the basic state. To provide a reference for evaluating 3DVAR and 4DVAR, an additional (control) experiment was performed in which no data were assimilated. In all experiments (hereafter referred to as EX3D, EX4D, and EXCL), the initial conditions on 1 January 1993 were generated from a 4-yr spinup of the model starting from rest and from Levitus (1982) climatological temperature and salinity. The wind stress forcing used for the first three years of the spinup was a climatology computed from the ERS wind stress products. The final year of the spinup was a transition year between ERS climatological and year 1992 products.

### b. Validity of the linear increment models

Incremental variational assimilation is founded on the linear approximation (6). In this section, we wish to examine the validity of this approximation by checking the accuracy of both the persistence model used in 3DVAR on a 10-day window and the TL model used in 4DVAR on a 30-day window. The accuracy of each model can be assessed by comparing the time evolution of an initial perturbation  $\delta \mathbf{w}$  in the nonlinear (*NL*) model with its evolution in the linear (*L*) model. The nonlinear perturbation can be computed from the finite difference,

$$\delta \mathbf{w}^{NL}(t_i) = M(t_i, t_0)(\mathbf{w}^b + \delta \mathbf{w}) - M(t_i, t_0)(\mathbf{w}^b), \quad (18)$$

while the linear perturbation is given by

$$\delta \mathbf{w}^{L}(t_{i}) = \mathbf{M}(t_{i}, t_{0}) \delta \mathbf{w}.$$
(19)

Ideally the initial perturbation  $\delta \mathbf{w}$  should have structure and amplitude typical of background errors. The actual background errors are not known, however; so in order to apply this test a suitable proxy must be defined. In meteorology, it is commonly assumed that the background errors can be roughly approximated from differences in model forecasts initiated from analyses at different time lags (Rabier et al. 1998). Here we consider a similar approach in which an initial perturbation is derived from the difference between two model states obtained from free integrations of the model having different initial states. One of the initial states is taken to be the *analysis* state ( $\mathbf{w}^a$ ) at the *end* of a 4DVAR assimilation interval [as defined by (14)], while the other initial state is taken to be the background state ( $\mathbf{w}^b$ ) valid at the same time (see Fig. 1b). In the first example presented below,  $\mathbf{w}^a$  and  $\mathbf{w}^b$  are taken from the end of the second 30-day cycle (on 1 March 1993) of EX4D, with  $\delta \mathbf{w} = \mathbf{w}^a - \mathbf{w}^b$  defining the initial perturbation in (18) and (19). The validity of the linear model was also investigated for other types of perturbations (analysis increments, differences between analyzed states at two different dates) and for other starting dates, and the results were qualitatively similar to those discussed below.

A meridional-vertical section of the perturbation temperature at 110°W is shown in Fig. 4a. The field is characterized by large perturbations of up to 4°C, appearing as a result of the temperature observations, which are assimilated during the second cycle. The perturbations after 10 and 30 days of integration in the nonlinear model are shown in Figs. 4b and 4c, respectively. Figure 4b allows us to check the assumption in 3DVAR that perturbations do not evolve significantly over 10 days outside the equatorial strip. The largeamplitude positive anomaly near 12°N and the smalleramplitude negative anomalies near 16°N and 7°S are indeed well approximated by the 10-day persistence model (cf. Figs. 4a and 4b). The persistence assumption breaks down nearer the equator where the oceanic dynamical response is much more rapid. In the persistence model, the positive anomaly at 4°N is largely overestimated, while the negative anomaly at 7°N is largely underestimated. In contrast, Figs. 4c and 4d illustrate that the TL model is able to provide a good description of the perturbation in the equatorial waveguide at least 30 days ahead. The similarity of Figs. 4c and 4d may come as little surprise since it is well known that the large spatial scale, low-frequency response of the tropical oceans involves predominantly linear wave dynamics (Philander 1989).

At smaller spatial and temporal scales (<1000 km, <1 month), energetic motions in the central and eastern equatorial Pacific are dominated by tropical instability waves (TIWs) (Legeckis 1977). Because of the importance of nonlinear processes in ultimately limiting the growth of unstable waves, it is interesting to see if some of the differences between perturbations in the TL and nonlinear models can be linked to TIWs. TIWs are most energetic during the autumn months. If TIW activity is indeed a limiting factor of the TL approximation, then one would expect this approximation to be less valid during this time. To check this, the above experiment was repeated using a starting date of 27 September 1993 and a temperature perturbation defined as the difference between the 4DVAR analysis and background state at that time (i.e., at the end of the ninth 30-day cycle of EX4D). A zonal-vertical section of the difference at 4°N between the nonlinear and TL model-predicted temperature perturbations after 30 days is shown in Fig. 5.



FIG. 4. (a) Meridional–vertical section at  $110^{\circ}$ W of a temperature perturbation defined as the difference between a 4DVAR analysis and the background state valid at the same time (1 Mar 1993). The perturbation after (b) 10 and (c) 30 days of evolution in the nonlinear model, and after (d) 30 days of evolution in the TL model. The contour interval is  $0.5^{\circ}$ C.



FIG. 5. Zonal-vertical section at 4°N of the difference between temperature perturbations in the nonlinear and TL models after 30 days. The initial perturbation in both experiments is defined as the difference between the 4DVAR analysis and the corresponding background state on 27 Sep 1993. The contour interval is 0.5°C.

West of 160°W, the TL approximation holds quite well. Most of the nonlinear behavior is located in the upper thermocline in the eastern and central Pacific, which is the area of largest thermal signal from TIWs. The differences between the nonlinear and linear perturbations are mostly associated with larger amplitudes of the linear perturbation. This confirms that nonlinear mechanisms play an important role in limiting the amplitude of the TIWs, and that the TL approximation is limited at their spatial scale. Over longer integration periods, the TIWs show up clearly as the dominant error signal (not shown), with the amplitude of the perturbations tending to be greatly overestimated in the TL model. This limitation is, however, only present at spatial scales of order 1000 km and during the active TIW season.

Strongly nonlinear processes associated with vertical mixing and convection are another factor contributing to the limitation of the TL approximation, particularly in the small vertical scales. For example, when a perturbation leads to a statically unstable water column in the nonlinear model, the water column will be mixed instantaneously as a result of enhanced vertical diffusion. This process will effectively reduce the amplitude of the perturbation in the nonlinear model. In the TL model, on the other hand, there is no linearized counterpart of this process and the perturbation will remain unchanged. This could explain the differences in the mixed layer seen in Fig. 5. Weaver et al. (2002) present additional diagnostics that provide further confirmation of these points.

In summary, the results from this section indicate that the 10-day persistence model used in 3DVAR is adequate for describing large-scale perturbations in offequatorial regions but has some limitations closer to the equator where the dynamical adjustment timescales are shorter. In comparison, the 30-day TL model provides a good description of large-scale perturbations in both off-equatorial and equatorial regions. Smaller-scale, regionally confined structures associated with TIWs, however, are more problematic and our experiments suggest that they are possibly the main source of error in the TL integration. Vertical mixing and convective processes also tend to limit the accuracy of the TL model in the small vertical scales. Further research is required to determine whether the accuracy of the TL model can be improved using a more sophisticated parametrization of vertical diffusion, in particular one that accounts for (nonzero) perturbations in the vertical mixing coefficients. However, it will be shown in section 4f and in Part II that these limitations on the validity of the TL model are much less severe than appears at first sight and that in fact a very good fit to the data is possible.

#### c. Convergence properties

Each 10-day 3DVAR cycle required roughly 25 min of CPU time on a Fujitsu VPP700, while each 30-day 4DVAR cycle required roughly 5 h of CPU time (i.e., a factor of 4 more costly than 3DVAR). The convergence of the minimization in 3DVAR was relatively quick, requiring on average less than 25 iterations to reach the effective minimum of *J*. The convergence of the minimization in 4DVAR was generally slower than in 3DVAR. After the final (60th) iteration, the reduction of the norm of  $\nabla_v J$  relative to its initial value was typically between three and four orders of magnitude compared to over six orders of magnitude in 3DVAR.

As a typical example, Fig. 6 (solid curve) shows the value of J from the second cycle of EX4D. The jumps in the solid curves in Fig. 6 occur after an outer iteration when the reference trajectory is reinitialized (here every 10 iterations). To illustrate the value of these outer iterations, the 4DVAR minimization on this particular cycle was repeated with only one outer iteration, that is, 60 minimization iterations with no trajectory updates (dashed curve in Fig. 6). Although it was found that convergence was more efficient without than with outer iterations, with a gain of approximately an order of magnitude in the reduction of the gradient norm, the reduction of the cost function in the experiment performed without outer iterations saturates at a level above the cost function with outer iterations. One has to be careful, however, in making direct comparisons between these two curves since, after the first outer iteration, the cost functions are no longer the same in the two experiments. It is more instructive to compare the final value of the full (nonapproximated) cost function  $J^F$  [Eq. (7)] for each experiment. These values are plotted in Fig. 6 with an asterisk (plus sign) for the experiment with (without) outer iterations. This figure provides a clear indication of the positive impact of the outer iterations. The final value of  $J^F$  with outer iterations is about half that without outer iterations. Furthermore, it is very close to the final value obtained with the approximated (quadratic) cost function (solid curve) and thus provides a good measure of the consistency of the incremental approach. In the absence of outer iterations, however, this consis-



FIG. 6. The value of the cost function (J) as a function of the minimization iteration on the second cycle of EX4D (solid curve). As a sensitivity test, the second cycle of the 4DVAR experiment has been repeated with only one outer iteration (dashed curve), compared to five outer iterations used in the reference experiment. The plus sign (asterisk) indicates the final value of the nonapproximated cost function (7) for the experiment with one (five) outer iteration(s). For clarity, these symbols have been displaced slightly to the left of the right border; *J* has been normalized by its respective value at the start of the minimization and plotted on a logarithmic vertical axis.

tency is lost as illustrated by the large discrepancy between the final value of  $J^F$  and that of the approximated cost function J (dashed curve).

In 3DVAR, 60 iterations was more than double the number actually needed to reach an acceptable level of convergence. Despite this potential to economize in 3DVAR, 60 iterations were retained simply to be consistent with the total number of iterations used in 4DVAR. In 4DVAR, convergence was slower but still reasonably efficient. Preconditioning techniques such as those described in Fisher and Andersson (2001) offer considerable scope for further improving the minimization efficiency. The outer iterations were clearly shown to be an essential feature of the algorithm. At present, however, it is not clear how many outer iterations are needed for the solution of the incremental problem to be a good approximation to the solution of the full problem. It is also not clear what combination of outer and inner iterations gives the best convergence rate. These issues are a matter for further research.

#### d. Optimality properties

The formulation of the variational assimilation problem relies on a number of hypotheses on the statistics of the background and observation errors. The validity of these hypotheses is an important factor in determining


FIG. 7. The value of  $\gamma = 2J_{\min}/p$  plotted as a function of the assimilation cycle during the period 1993–98 in (a) 3DVAR and (b) 4DVAR. On each cycle,  $J_{\min}$  represents the value of the cost function at the end of the minimization and p the total number of observations assimilated through the  $J_o$  term. A total of 219 (73) 10-day (30 day) cycles were performed in 3DVAR (4DVAR). The *expected* value of  $\gamma = 1$  (solid line) is plotted together with error bars (dashed lines) at  $1 \pm \sigma_{\gamma}$ , where  $\sigma_{\gamma} = \sqrt{2/p}$  is the *expected* standard deviation of  $\gamma$ . The dotted curves in (a) and (b) correspond to the *actual* values of  $\gamma$  computed on each cycle in the reference experiments EX3D and EX4D. The dashed–dotted curve in (a) corresponds to the *actual* values of  $\gamma$  from a 3DVAR experiment in which the observation-error variance for the TAO data is set to twice the value used in EX3D.

the optimality of the analysis. One particularly simple diagnostic that can be used to check whether the statistics of **B** and **R** are consistent with the innovation vector is the value of the cost function at its minimum  $(J_{\min})$ , which, for a linear system, should, on statistical average, be equal to p/2 where p is the total number of assimilated observations (Tarantola 1987; Bennett et al. 2000). If we assume further that the background and observation errors are Gaussian then it can be shown (e.g., see section 4.3.6 in Tarantola 1987) that the expected variance of  $J_{\min}$  is  $\sigma_{J_{\min}}^2 = E[(J_{\min} - E[J_{\min}])^2]$ = p/2, where E[] is the mathematical expectation operator. Therefore, provided **B** and **R** are correctly specified and that the system is quasi-linear, the value of  $J_{\min}$ on each assimilation cycle should be p/2 within a standard deviation of  $\sqrt{p/2}$ .

In order to compare the actual values of  $J_{\min}$  with its expected value, it is more convenient to consider the normalized quantity  $\gamma = 2J_{\min}/p$  for which  $E[\gamma] = 1$ and  $\sigma_{\gamma}^2 = E[(\gamma - E[\gamma])^2] = 2/p$ . Here,  $\sigma_{\gamma}^2$  scales with the inverse of p so will be small when a large number of observations are assimilated. Figures 7a and 7b (dotted curves) show the actual values of  $\gamma$  as a function of the assimilation cycle in EX3D and EX4D. The expected value of  $\gamma = 1$  (solid line) is also plotted together with the expected error bars (dashed lines) at  $1 \pm \sigma_{\gamma}$ . The average of the actual values of  $\gamma$  over all cycles in EX3D is close to its expected value of one ( $\overline{\gamma} = 0.90$ , where the overbar denotes cycle average) but in EX4D it is somewhat too small ( $\overline{\gamma} = 0.73$ ). On average, the actual values of  $\gamma$  exceed its expected value by  $9\sigma_{\gamma}$  in EX3D and  $29\sigma_{\gamma}$  in EX4D. The expected error bars are slightly larger in EX3D than in EX4D since there were fewer observations assimilated during a 10-day 3DVAR cycle than during a 30-day 4DVAR cycle (roughly 10 000 observations were available every 10 days). Even so, the actual values of  $\gamma$  nearly always greatly exceed the expected error bounds in both experiments.

Consider the expression for  $J_{\min}$  in terms of the innovation vector **d** (Tarantola 1987):

$$J_{\min} = \frac{1}{2} \mathbf{d}^{\mathrm{T}} (\mathbf{G} \mathbf{B} \mathbf{G}^{\mathrm{T}} + \mathbf{R})^{-1} \mathbf{d}, \qquad (20)$$

where the generalized quantities  $\mathbf{d} = (\dots, \mathbf{d}_i^{\mathrm{T}}, \dots)^{\mathrm{T}}, \mathbf{G}$ =  $(\ldots, \mathbf{G}_i^{\mathrm{T}}, \ldots)^{\mathrm{T}}$ , and  $\mathbf{R} = \operatorname{diag}(\ldots, \mathbf{R}_i, \ldots)$  have been introduced to simplify the notation. From (20) it is clear that the value of  $J_{\min}$ , or equivalently  $\gamma$ , may be decreased (increased) if the variances of either **B** or **R** are increased (decreased). The lower than expected value of  $\gamma$  in EX4D therefore could be a sign that either the background- or observation-error variances have been overestimated. Our prior estimate of the observation-error variances is very crude and in particular takes no account of possible representativeness error. For example, the value of 0.5°C used for the standard error of TAO data is only slightly larger than the documented estimate of TAO observation error, which accounts only for the instrumental component of that error (McCarty et al. 1997). The specified observation-error variances are thus probably underestimated and therefore acting to increase the value of  $\gamma$ . This effect is clearly demonstrated by the dashed-dotted curve in Fig. 7a, which shows the actual values of  $\gamma$  from a 3DVAR experiment in which the standard error for TAO data was increased to 1.0°C. This led to a factor of 2 decrease in the cycle average of  $\gamma(\overline{\gamma} = 0.44)$ .

It is more likely that the background-error variances, which have been specified from the climatology of a model integration performed without data assimilation, are substantially overestimated. Whereas this specification might yield a reasonable estimate of the background errors for the first cycle, it is probably too large an estimate in later cycles, particularly in well-observed regions, where the data assimilated in previous cycles have acted to reduce the innovation vector. Generally speaking, however, it is safer to overestimate than underestimate the background-error variances to prevent the model from drifting too far from observations. It should be noted that a misspecification of the background- and/or observation-error correlations could also change the value of  $\gamma$ .

Another interesting feature in Fig. 7 is the rather large variance in  $\gamma$  between cycles. This is an indication that, contrary to what has been assumed, the backgrounderror covariances are not stationary, in addition to being largely overestimated as already mentioned. From Fig. 7a, we see that the variability of  $\gamma$  is larger in EX3D than in the sensitivity experiment where the standard error used for TAO data was larger. This is understand-able since the smaller standard error for TAO in EX3D means that  $\gamma$  will be more sensitive to the true variations in the background-error covariances between cycles.

While the  $J_{\min}$  statistic gives some useful insight into the optimality properties of the system, it is not possible based on this information alone to correct unambiguously any misspecification of the background- and/or observation-error covariances. If the actual value of  $J_{\min}$ is not equal to p/2, then this can be rectified simply by muliplying **B** and **R** by the factor  $\gamma = 2J_{\min}/p$ . This procedure will change the absolute values of the background-, observation-, and analysis-error variances, but will have no influence on the analysis increment itself. What is important then is the *relative* magnitude of the variances in **B** and **R** (the absolute values can be obtained by postmultiplication by  $\gamma$ ). Adaptive procedures can be used to tune the variances in **B** and **R** using information on the mismatch between the expected and actual values of subparts of the cost function at its minimum (Desroziers and Ivanov 2001). To compute the expected minimum value of subparts of the cost function is more complicated, but practical methods do exist (Bennett et al. 2000; Desroziers and Ivanov 2001). We have not attempted to apply any of these methods in the present study but they do offer an interesting possibility for improving the variance estimates in the future.

#### e. Flow-dependent background-error variances

It is well known that, in the limit of a perfect, linear model, variational assimilation is equivalent to the Kalman filter in that, given identical input parameters, they produce the same analysis at the end of the assimilation window (e.g., see Courtier et al. 1994). Consider an assimilation window  $t_0 \le t_i \le t_n$  in which a background state  $\mathbf{w}^b$ , with error-covariance matrix **B**, is available at the beginning of the assimilation window and an observation vector  $\mathbf{y}_n^o$ , with error-covariance matrix  $\mathbf{R}_n$ , is available at the *end* of the window. Within the linear approximation, the error-covariance matrix  $\mathbf{P}^b(t_n)$  for the background state  $\mathbf{w}^b(t_n) = M(t_n, t_0)(\mathbf{w}^b)$  is obtained by evolving **B** using the linear model and its adjoint:

$$\mathbf{P}^{b}(t_{n}) = \mathbf{M}(t_{n}, t_{0})\mathbf{B}\mathbf{M}(t_{n}, t_{0})^{\mathrm{T}}, \qquad (21)$$

where  $\mathbf{P}^{b}(t_{0}) = \mathbf{B}$ . In an extended Kalman filter, (21) would be used explicitly to transport the covariances forward in time. In incremental variational assimilation, on the other hand, this propagation is implicit in the global minimization process. For the example above, the variational analysis increment at  $t_{0}$ , obtained by minimizing (8) exactly, can be written as

$$\delta \mathbf{w}^{a} = \mathbf{B} \mathbf{M}(t_{n}, t_{0})^{\mathrm{T}} \mathbf{H}_{n}^{\mathrm{T}} [\mathbf{H}_{n} \mathbf{M}(t_{n}, t_{0}) \mathbf{B} \mathbf{M}(t_{n}, t_{0})^{\mathrm{T}} \\ \times \mathbf{H}_{n}^{\mathrm{T}} + \mathbf{R}_{n}]^{-1} \mathbf{d}_{n}, \qquad (22)$$

where  $\mathbf{d}_n = \mathbf{y}_n^o - H_n[\mathbf{w}^b(t_n)]$ . The analysis increment at  $t_n$  can be obtained by applying  $\mathbf{M}(t_n, t_0)$  to both sides of (22) to yield, after inserting (21),

$$\delta \mathbf{w}^{a}(t_{n}) = \mathbf{P}^{b}(t_{n})\mathbf{H}_{n}^{\mathrm{T}}[\mathbf{H}_{n}\mathbf{P}^{b}(t_{n})\mathbf{H}_{n}^{\mathrm{T}} + \mathbf{R}_{n}]^{-1}\mathbf{d}_{n}, \quad (23)$$

where  $\delta \mathbf{w}^{a}(t_{n}) = \mathbf{M}(t_{n}, t_{0}) \delta \mathbf{w}^{a}$ . Equation (23) is the standard analysis step of an extended Kalman filter, which weights the background state at  $t_{n}$  using an error-covariance matrix predicted by (21).

Our FGAT version of 3DVAR can be viewed as a limiting case in which **M** is taken to be the identity matrix. This implies that  $\mathbf{P}^{b}(t_{n}) = \mathbf{B}$ , that is, that the background-error covariances at the end of the window (and at all intermediate times) are identical to those specified at the initial time. In incremental 4DVAR, on the other hand, **M** is the TL operator so that  $\mathbf{P}^{b}(t_{n})$  will be modified by dynamical processes acting within the window. In this section, we focus on how the TL dynamics act to modify the prior estimates of the background-error variances [the diagonal elements of  $\mathbf{P}^{b}(t_{n})$ ]. The TL dynamics will also modify the background-error correlations but a discussion of this feature is beyond the scope of the present paper.

Single-observation experiments provide a practical way of computing the effective background-error variances used implicitly in 4DVAR (Thépaut et al. 1993, 1996). For a single observation,  $d = \mathbf{d}_n$  and  $(\sigma^o)^2 = \mathbf{R}_n$  become scalars, and  $\mathbf{h}^T = \mathbf{H}_n$  becomes a vector of the same length as w. The error variance of the background equivalent of the observation is then given by the scalar product  $(\sigma_n^b)^2 = \mathbf{h}^T \mathbf{P}^b(t_n) \mathbf{h}$ . If we assume further that the observation is of one of the model state variables and that its location coincides with a model grid point, then  $\mathbf{h}^T = (0, \dots, 0, 1, 0, \dots 0)$ , where the





FIG. 8. (a) Vertical profiles of the background-error std devs (°C) used in 3DVAR and 4DVAR at the equator at 140°W. The dashed-dotted curves correspond to the standard deviations *specified* at the *beginning* of the assimilation window. In 3DVAR, these are also the *effective* standard deviations used at all future times within the window. The solid curves correspond to the *effective* standard deviations used in 4DVAR at the *end* of the 30-day window of the second cycle in EX4D. (b) The corresponding profile of  $|\partial T_{\theta}/\partial z|$  at this location, computed from the 30-day mean of the background temperature state on the second cycle. The values of  $|\partial T_{\theta}/\partial z|$  have been multiplied by a factor of 10 in order to be plotted with the same horizontal scale as in (a).

entry equal to one corresponds to that grid point, and  $(\sigma_n^b)^2$  becomes an element of the diagonal of  $\mathbf{P}^b(t_n)$ . It is straightforward to deduce  $(\sigma_n^b)^2$  by applying  $\mathbf{h}^T$  to both sides of (23):

$$\mathbf{h}^{\mathrm{T}} \boldsymbol{\delta} \mathbf{w}^{a}(t_{n}) = (\boldsymbol{\sigma}_{n}^{b})^{2} \left[ \frac{d}{(\boldsymbol{\sigma}_{n}^{b})^{2} + (\boldsymbol{\sigma}^{o})^{2}} \right], \qquad (24)$$

which can be rearranged to give

$$(\boldsymbol{\sigma}_n^b)^2 = \left[\frac{\mathbf{h}^{\mathrm{T}} \delta \mathbf{w}^a(t_n)}{d - \mathbf{h}^{\mathrm{T}} \delta \mathbf{w}^a(t_n)}\right] (\boldsymbol{\sigma}^o)^2.$$
(25)

The term  $\mathbf{h}^{\mathsf{T}} \delta \mathbf{w}^a(t_n)$  is the value of the analysis increment at the time and gridpoint location of the single observation. Equation (25) thus provides the basis of an algorithm for systematically diagnosing the exact diagonal elements of  $\mathbf{P}^b(t_n)$ . To compute all diagonal elements of  $\mathbf{P}^b(t_n)$  using (25) would be prohibitively expensive since (25) requires as many single observation experiments as the number of elements of  $\mathbf{w}$ . This algorithm is therefore only practical for computing a small subset of the  $\sigma_n^b$ , which is what is desired here. If an estimate of the complete diagonal of  $\mathbf{P}^b(t_n)$  were required, then a more efficient, but approximate, algorithm based on randomization would be more appropriate (e.g., Andersson et al. 2000).

Figure 8 shows vertical profiles of the backgrounderror standard deviations ( $\sigma_n^b$ ) for temperature at the equator in the central Pacific (140°W). The dasheddotted curve is the prior estimate of  $\sigma_n^b$ , which, as discussed earlier, was computed from a control experiment without data assimilation. In 3DVAR, the prior estimates are effectively used to weight the background state at all times within the assimilation window, while in 4DVAR they are used for weighting the background state at the beginning of the assimilation window. The prior estimate displays a maximum around the depth of the climatological thermocline where variability is greatest (around 170 m). The solid curve shows the profile of the effective  $\sigma_n^b$  used in 4DVAR at the end of the second 30-day cycle (3 March 1993) of EX4D (see Part II for a thorough description of this experiment). By repeating this experiment at different latitudes it was found that the flow dependency in the  $\sigma_n^b$  estimates is greatest near the equator, which is not surprising since dynamical effects have shorter timescales there. In 30 days, the ocean state can change significantly at the equator but much less so at higher latitudes. Figure 8a also shows that the TL dynamical processes tend to reduce the  $\sigma_n^b$  in the ocean mixed layer, particularly in the upper 70 m. Note that the Newtonian relaxation term for SST is contributing to this tendency by damping temperature perturbations in the top level of the TL model.

The TL dynamics also tend to move the level of maximum background-error variance to the level of the thermocline, as illustrated by comparing Fig. 8a with the profile of  $|\partial T_{\theta}/\partial z|$  computed from the 30-day-averaged background temperature state on the second cycle. This tendency is physically sensible since the level of maximum variability of the thermal field, and thus of maximum likely error in the background state, is located at the level of the thermocline. It can also be noted that the TL dynamics can substantially increase the maximum value of the background-error variance near the equator. This is related to the fact that the background trajectory in this experiment has already felt the influence of observations assimilated during the previous cycle through a tightening of the thermocline. This tighter, well-defined thermocline, relative to that of the control, leads to larger thermal signals and thus to an increase in the maximum value of  $\sigma_n^b$  by the TL dynamics. In some 3D assimilation systems (e.g., Behringer et al. 1998; Alves et al. 2002), the background-error variances have been parameterized by making them dependent on the vertical gradient of the background temperature field. Figure 8 suggests that relating  $\sigma^b$  to the background temperature gradient is dynamically sensible. Such a weighting could also be useful in 3DVAR and 4DVAR by introducing a flow dependence in the variances of **B** at the beginning of the window.

#### f. Fit to the observations

Statistics derived from the background minus observation vector {BmO =  $H_i[\mathbf{w}^b(t_i)] - \mathbf{y}_i^o$ } and analysis minus observation vector {AmO =  $H_i[\mathbf{w}^a(t_i)] - \mathbf{y}_i^o$ } can yield useful information about the internal consistency of the data assimilation system (Hollingsworth and Lönnberg 1989). Figures 9a and 9b show the 1993-98 averaged statistics of the BmO and AmO as a function of depth for all the assimilated TAO data in EX3D and EX4D. For reference, the average statistics of the difference between the control and TAO data are also shown, which for convenience will be referred to as the BmO of EXCL. In Fig. 9a the BmO curve of EXCL displays a large warm bias of about 2°C just below the thermocline. This bias is largely reduced in EX3D, with the maximum of the averaged BmO and AmO being about 0.3°C at 50 m. The bias is almost completely absent from EX4D: the average of AmO is very small, except at 250 m where the analysis is about 0.1°C too cold.

In Fig. 9b the rms of BmO of EXCL shows large differences both just below the thermocline, where there are large biases, and in the thermocline, where signals associated with the seasonal cycle and interannual variability are largest. These differences are substantially reduced in EX3D and EX4D. The rms of BmO in EX4D is similar to that of EX3D, with a maximum around 1.5°C at the level of the thermocline. On the other hand, the rms of AmO in EX4D is very much reduced, being less than 0.5°C over the entire water column, compared to the rms of AmO in EX3D, which is hardly smaller than that of BmO. The fit to the TAO data in EX4D is thus within the specified level of the observation error  $(0.5^{\circ}C)$ , which is not the case in EX3D. The larger AmO in EX3D is primarily an artifact of the gradual way the analysis increment is applied to the model. While minimizing nonphysical adjustment processes, this procedure ultimately degrades the fit to the data achieved by the linear analysis  $\mathbf{w}^{b}(t_{i}) + \delta \mathbf{w}^{a}(t_{i})$ . This is illustrated in Fig. 9c, which shows that for EX3D the rms of the  $J_o$  residual,  $\mathbf{H}_i \delta \mathbf{w}^a(t_i) - \mathbf{d}_i$ , is considerably less than the AmO and is comparable to the specified observation error (0.5°C). An additional experiment was performed



FIG. 9. Time-averaged statistics, plotted as a function of depth, of BmO and AmO for TAO data during the 1993–98 period. (a) The average of BmO (thin curves) and AmO (thick curves), (b) the rms of BmO (thin curves) and AmO (thick curves), and (c) the rms of AmO (thick curves) and of the  $J_o$  residual (thin curves). In (a) and (b), the dashed curve corresponds to EXCL; in all panels, the dashed-dotted curves correspond to EX3D, and the solid curves in (c)] are indistinguishable.

1375



FIG. 10. Spatially averaged rms of BmO (thin curves) and AmO (thick curves) for TAO data, plotted as a function of the assimilation cycle during the 1993–98 period. The dashed curve corresponds to EXCL, the dashed–dotted curves to EX3D, and the solid curves to EX4D. The statistics for EX3D are displayed as an average over three cycles (30 days) in order to be compared with those from EX4D.

in which the (temperature) analysis increment was added directly to the background initial conditions as in 4DVAR. The statistics of the AmO and  $J_o$  residual for that experiment were comparable to those of EX3D but the overall impact of the assimilation on the model fields was significantly worse and therefore justifies our approach of applying the analysis increment gradually. For EX4D, the residual and AmO are indistinguishable in Fig. 9c. This result is consistent with Fig. 6 of section 4c, which showed that, after several outer iterations, the full and incremental cost functions converged to similar values at the end of minimization. That result was demonstrated for a particular cycle (the second); Fig. 9c just confirms that this feature is consistent for all cycles of EX4D.

Figure 9 provided a time-averaged view of the BmO and AmO statistics for the assimilated TAO data. Figure 10 now provides a view of how these statistics change with time during the 1993–98 analysis period. On the first assimilation cycle, the background is provided by the control and as a result the BmO is very large (Fig. 10). After the first assimilation cycle, however, the properties observed in Fig. 9 begin to emerge: the average of AmO is small in both EX3D and EX4D, and the rms of AmO is equal to about 1°C in EX3D and 0.5°C in EX4D. It was found that, in both EX3D and EX4D, the statistics of BmO and AmO are approximately stationary, and in particular do not seem to exhibit any obvious dependence on either the number of assimilated observations or interannual variability.

The cycle average of the temporal evolution of BmO and AmO *within* the assimilation window is shown in Fig. 11. For the background, which is not constrained by observations, the fit to the data degrades with time. After 10 days in EX3D and 30 days in EX4D, the rms of BmO increases by  $0.2^{\circ}$  and  $0.5^{\circ}$ C from initial values of  $1.0^{\circ}$  and  $0.6^{\circ}$ C, respectively. The rms of BmO at the end of the window is  $1.2^{\circ}$ C in both experiments, compared to  $1.9^{\circ}$ C in EXCL. In both EX3D and EX4D, the fit of the analysis to the data is uniform over their respective assimilation windows, apart from a very slight increase of the AmO at the window boundaries in EX4D. Note that, since the analysis at the end of a given



FIG. 11. Rms of the 1993–98 cycle average of BmO (thin curves), AmO (thick curves), and  $J_o$  residual (thick curves) for TAO data plotted as a function of the day within the assimilation window. The dashed curve corresponds to EXCL, the dashed–dotted curves to EX3D, and the solid curves to EX4D. The upper (lower) of the two thick dashed–dotted curves corresponds to the AmO ( $J_o$  residual) of EX3D. The AmO and  $J_o$  residual (thick solid curves) of EX4D are indistinguishable.

cycle is used as the background at the beginning of the next cycle, the BmO at day 0 should be equal to the AmO at day 30, except for the first and last cycle for which there is no corresponding AmO (day 30) and BmO (day 0), respectively. This latter point explains the slight difference between these quantities in Fig. 11. When model error is important, and not explicitly accounted for in the assimilation method, it often manifests itself as a U shape in the fit to the data (Ménard and Daley 1996). Any signal of model error, however, is more likely to be visible in the  $J_o$  residual than in the AmO since the residual is the quantity that is minimized objectively. The residual will contain information about linearization errors in the nonlinear model itself.

The cycle average of the  $J_o$  residuals in EX3D and EX4D is also shown in Fig. 11. For EX3D (lower thick dashed–dotted curve), the U shape is evident, though quite weak, with a difference of about 0.2°C between the highest point at the window boundaries and the lowest point in the middle of the window. If this is indeed a sign of model error, then it would be consistent with the results of section 4b, which illustrated that the persistence model did have some limitations near the equator. In contrast, in EX4D, the residuals are nearly flat over the whole window width (and are very similar to the AmO). This result then suggests that model error is not a significant problem over the 30-day window used in 4DVAR.

#### 5. Summary

Three- and four-dimensional variational assimilation systems have been developed for the rigid-lid version of the OPA OGCM (Madec et al. 1998). The assimilation problem is defined by a cost function that penalizes departures from the data and from a background estimate that is the result of a previous assimilation. The control variable of the cost function is the initial condition at the start of each assimilation window. The cost function is minimized using an incremental approach by which the full minimization problem, involving the nonlinear OGCM as a constraint, is approximated by a sequence of quadratic minimization problems involving linear constraints. The control variable of each quadratic cost function is an increment to the initial conditions. In the incremental 3DVAR, the increment is transported forward in time by a persistence model. In the incremental 4DVAR, a linearized version of the full OGCM is used to propagate the increment. The adjoint of the linear model propagator is used to compute the gradient of the cost function with respect to the initial condition increment thus allowing the solution to be found iteratively using a gradient descent method. Once the analysis increment is found, it is applied at the beginning of the assimilation window to derive the 4DVAR analysis trajectory. In the case of 3DVAR, it is applied as a constant 3D forcing to the model equations over the assimilation window as a way of minimizing spurious adjustments.

The systems have been applied to assimilate in situ temperature data in the tropical Pacific Ocean. Three experiments were performed for the 1993-98 period: a control experiment without data assimilation, an experiment using 3DVAR cycled with a 10-day window, and an experiment using 4DVAR cycled with a 30-day window. The validity of the persistence model (3DVAR) and of the TL model (4DVAR) was investigated. It was shown that persistence was a reasonable assumption over 10 days, at least outside the equatorial waveguide, and that the TL model provides a good description of the large-scale oceanic state over at least 30 days. However, because of the nonlinear nature of convective and vertical mixing processes, the validity of the TL model could be degraded at small vertical scales. Tropical instability waves (TIWs) are nonlinear oscillations with a timescale of about 1 month. At their horizontal scale, the TL model is also less accurate. Nevertheless, in 4DVAR a very good fit to the observed TIWs was achieved (see Part II). This result points to the important role of the outer iterations in the incremental 4DVAR formulation in providing a feedback mechanism between the TL and nonlinear models so that the model can eventually achieve a very close fit to the data. With outer iterations, the final values of the incremental and nonincremental cost functions were shown to be very close, which provides a good measure of the consistency of the incremental approach for solving the original nonlinear minimization problem. Without outer iterations, the performance of the 4DVAR system is seriously degraded.

Single-observation experiments have been performed to illustrate the effect of the TL dynamics in modifying the prior estimates of the background-error variances. The TL dynamics were shown to modify the variances in a physically sensible way. The variances were diminished in the mixed layer, and the maximum value of the variance in the profile could be increased and displaced to the level of the background thermocline, where thermal variability (and background error) is greatest.

A detailed examination of the fit of the different analyses to the assimilated data has been made. The control experiment displays a large bias below the thermocline, which is strongly reduced in the 3DVAR analyses and almost entirely absent from the 4DVAR analyses. The rms difference between the analyses and observations is also very much reduced in the thermocline region. For example, for TAO data, whereas the rms difference is about 2.8°C in the control, it falls to 1.4°C in 3DVAR and to below 0.5°C in 4DVAR, which is less than the specified standard deviation of the observation error. Over the TAO region, the spatially averaged rms difference between the observations and the 3DVAR and 4DVAR analyses is stationary in time (equal to 1°C in 3DVAR and 0.4° in 4DVAR) and, in particular, does not exhibit any dependence on the number of observations or interannual variability.

In this paper, the 3DVAR and 4DVAR systems have been described and evaluated in terms of certain algorithmic and statistical diagnostics. In Part II, the analyses produced by the two systems are examined from a physical perspective. A discussion of the results of the two papers and possible avenues for future development are given at the end of Part II.

Acknowledgments. This work was initiated at LO-DYC with funding from an EC Human Capital and Mobility fellowship for the first author. Further funding was provided by the French MERCATOR and EC-FP5 EN-ACT projects. Discussions with Philippe Courtier, Pascale Delecluse, Eric Greiner, Gurvan Madec, and François Vandenberghe in the early stages of development of the assimilation system are gratefully acknowledged. The first author would like to thank ECMWF for hosting him for a year when much of this work was carried out. Erik Andersson, Magdalena Balmaseda, Mike Cullen, Mike Fisher, Andy Moore, Andrea Piacentini, and Philippe Rogel provided many useful suggestions for improving the manuscript. We are also grateful to three anonymous reviewers for their constructive and insightful remarks.

#### REFERENCES

- Alves, J. O., M. A. Balmaseda, D. L. T. Anderson, and T. N. Stockdale, 2002: Sensitivity of dynamical seasonal forecasts to ocean initial conditions. ECMWF Tech. Memo. 369, 24 pp. [Available online at http://www.ecmwf.int/publications/.]
- Andersson, E., M. Fisher, R. Munro, and A. McNally, 2000: Diagnosis of background errors for radiances and other observable quantities in a variational data assimilation system, and the explanation of a case of poor convergence. *Quart. J. Roy. Meteor. Soc.*, **126**, 1455–1472.
- Balmaseda, M., D. L. T. Anderson, and M. Davey, 1994: ENSO prediction using a dynamical ocean model coupled to statistical atmospheres. *Tellus*, 46A, 497–511.
- Barnett, T. P., M. Latif, N. Graham, M. Flugel, S. Pazan, and W. White, 1993: ENSO and ENSO-related predictability. Part I:

- Behringer, D., M. Ji, and A. Leetma, 1998: An improved coupled model for ENSO prediction and implications for ocean initialization. Part I: The ocean data assimilation system. *Mon. Wea. Rev.*, **126**, 1013–1021.
- Bennett, A. F., B. S. Chua, D. E. Harrison, and M. J. McPhaden, 2000: Generalized inversion of Tropical Atmosphere–Ocean (TAO) data using a coupled model of the tropical Pacific. J. Climate, 13, 2770–2785.
- Bloom, S. C., L. L. Takacs, A. M. Da Silva, and D. Ledvina, 1996: Data assimilation using incremental analysis updates. *Mon. Wea. Rev.*, **124**, 1256–1271.
- Bonekamp, H., G. J. van Oldenborgh, and G. Burgers, 2001: Variational assimilation of TAO and XBT data in the HOPE OGCM, adjusting the surface fluxes in the tropical ocean. J. Geophys. Res., 106, 16 693–16 709.
- Bryan, K., 1969: A numerical method for the study of the circulation of the world ocean. *J. Comput. Phys.*, **4**, 347–376.
- Cane, M. A., S. E. Zebiak, and S. C. Dolan, 1986: Experimental forecasts of El Niño. *Nature*, **321**, 827–832.
- Courtier, P., J.-N. Thépaut, and A. Hollingsworth, 1994: A strategy for operational implementation of 4DVAR, using an incremental approach. *Quart. J. Roy. Meteor. Soc.*, **120**, 1367–1388.
- —, and Coauthors, 1998: The ECMWF implementation of three dimensional variational assimilation (3DVAR). Part I: Formulation. *Quart. J. Roy. Meteor. Soc.*, **124**, 1783–1808.
- Derber, J. C., and A. Rosati, 1989: A global oceanic data assimilation system. J. Phys. Oceanogr., 19, 1333–1347.
- Desroziers, G., and S. Ivanov, 2001: Diagnosis and adaptive tuning of observation error parameters in a variational assimilation. *Quart. J. Roy. Meteor. Soc.*, **127**, 1433–1452.
- Egbert, G. D., A. F. Bennett, and M. G. G. Foreman, 1994: TOPEX/ Poseidon tides estimated using a global inverse model. J. Geophys. Res., 99, 24 821–24 852.
- Fisher, M., and E. Andersson, 2001: Developments in 4DVAR and Kalman filtering. ECMWF Tech. Memo. 347, 36 pp. [Available online at http://www.ecmwf.int/publications/.]
- Gandin, L. S., 1965: *Objective Analysis of Meteorological Fields*. Israeli Program for Scientific Translations, 242 pp.
- Gauthier, P., C. Charette, L. Fillion, P. Koclas, and S. Laroche, 1999: Implementation of a 3D variational data assimilation system at the Canadian Meteorological Centre. Part I: The global analysis. *Atmos.–Ocean*, 37, 103–156.
- Gibson, J., P. Kållberg, S. Uppala, A. Noumura, A. Hernandez, and E. Serrano, 1997: ERA description. ECMWF Re-Analysis Project Report Series, No. 1, 77 pp. [Available online at http:// www.ecmwf.int/research/era/ERA-15/Report\_Series/.]
- Giering, R., and T. Kaminski, 1998: Recipes for adjoint code construction. ACM Trans. Math. Software, 24, 437–474.
- Gilbert, J.-C., and C. Lemaréchal, 1989: Some numerical experiments with variable-storage quasi-Newton algorithms. *Math. Program.*, 45, 407–435.
- Gill, A. E., 1982: Atmosphere–Ocean Dynamics. Academic Press, 662 pp.
- Greiner, E., and S. Arnault, 2000: Comparing the results of a 4Dvariational assimilation of satellite and in situ data with WOCE CITHER hydrographic measurements in the tropical Atlantic. *Progress in Oceanography*, Vol. 47, Pergamon Press, 1–68.
  - —, —, and A. Morlière, 1998a: Twelve monthly experiments of 4D-variational assimilation in the tropical Atlantic during 1987: Part I: Method and statistical results. *Progress in Oceanography*, Vol. 41, Pergamon Press, 141–202.
  - , —, and —, 1998b: Twelve monthly experiments of 4D-variational assimilation in the tropical Atlantic during 1987: Part II: Oceanographic interpretation. *Progress in Oceanography*, Vol. 41, Pergamon Press, 203–247.
- Grima, N., A. Bentamy, K. Katsaros, Y. Quilfen, P. Delecluse, and C. Lévy, 1999: Sensitivity study of an oceanic general circulation

model forced by satellite wind-stress fields. J. Geophys. Res., 104, 7967–7989.

- Hollingsworth, A., and P. Lönnberg, 1989: The verification of objective analyses: Diagnostics of analysis system performance. *Meteor. Atmos.*, 40, 3–27.
- Ide, K., P. Courtier, M. Ghil, and A. C. Lorenc, 1997: Unified notation for data assimilation: Operational, sequential and variational. J. *Meteor. Soc. Japan*, **75**, 181–189.
- Janiskova, M., J.-N. Thépaut, and J.-F. Geleyn, 1999: Simplified and regular physical parametrizations for incremental four-dimensional variational assimilation. *Mon. Wea. Rev.*, 127, 26–45.
- Ji, M., and A. Leetma, 1997: Impact of data assimilation on ocean initialization and El Niño prediction. *Mon. Wea. Rev.*, 125, 742– 753.
- Kessler, W. S., M. C. Spillane, M. J. McPhaden, and D. E. Harrison, 1996: Scales of variability in the equatorial Pacific inferred from the Tropical Atmosphere–Ocean buoy array. J. Climate, 9, 2999– 3024.
- Laroche, S., and P. Gauthier, 1998: A validation of the incremental formulation of 4D variational data assimilation in a nonlinear barotropic flow. *Tellus*, **50A**, 557–572.
- Latif, M., and M. Flügel, 1991: An investigation of short-range climate predictability in the tropical Pacific. J. Geophys. Res., 96, 2661–2673.
- Le Dimet, F. X., and O. Talagrand, 1986: Variational algorithms for analysis and assimilation of meteorological observations: Theoretical aspects. *Tellus*, **38A**, 97–110.
- Legeckis, R., 1977: Long waves in the eastern equatorial Pacific; a view from a geostationary satellite. *Science*, **197**, 1177–1181.
- Levitus, S., 1982: *Climatological Atlas of the World Ocean*. NOAA Prof. Paper 13, 173 pp. and 17 microfiche.
- Lorenc, A. C., 1981: A global three-dimensional multivariate statistical interpolation scheme. *Mon. Wea. Rev.*, **109**, 701–721.
- —, and Coauthors, 2000: The Met. Office global three-dimensional variational data assimilation scheme. *Quart. J. Roy. Meteor. Soc.*, **126**, 2991–3012.
- Madec, G., P. Delecluse, M. Imbard, and C. Lévy, 1998: OPA 8.1 Ocean General Circulation Model reference manual. LODYC/ IPSL Technical Note 11, Paris, France, 91 pp. [Available online at http://www.lodyc.jussieu.fr/opa/.]
- Mahfouf, J.-F., 1999: Influence of physical processes on the tangentlinear approximation. *Tellus*, 51A, 147–166.
- McCarty, M. E., L. J. Mangum, and M. J. McPhaden, 1997: Temperature errors in TAO data induced by mooring motion. NOAA Tech. Memo. ERL PMEL-108, 14 pp.
- McPhaden, M. J., 1993: TOGA-TAO and the 1991-93 El Niño–Southern Oscillation event. *Oceanography*, **6**, 36–44.
- Ménard, R., and R. Daley, 1996: The application of Kalman smoother theory to the estimation of 4DVAR error statistics. *Tellus*, 48A, 221–237.
- Meyers, G., H. Phillips, N. Smith, and J. Sprintall, 1991: Space and time scales for optimal interpolation of temperature—Tropical Pacific Ocean. *Progress in Oceanography*, Vol. 28, Pergamon Press, 189–218.
- Palmer, T. N., and D. L. T. Anderson, 1994: Prospects for seasonal forecasting. *Quart. J. Roy. Meteor. Soc.*, **120**, 755–794.
- Parrish, D. F., and J. C. Derber, 1992: The National Meteorological Center's spectral statistical interpolation analysis system. *Mon. Wea. Rev.*, **120**, 1747–1763.
- Philander, S. G., 1989: El Niño, La Niña, and the Southern Oscillation. Academic Press, 293 pp.
- Rabier, F. A. McNally, E. Andersson, P. Courtier, P. Undén, J. Eyre, A. Hollingsworth, and F. Bouttier, 1998: The ECMWF implementation of three-dimensional variational assimilation (3DVAR). Part II: Structure functions. *Quart. J. Roy. Meteor. Soc.*, **124**, 1809–1829.
- —, H. Järvinen, E. Klinker, J. F. Mahfouf, and A. Simmons, 2000: The ECMWF operational implementation of four-dimensional variational assimilation. Part I: Experimental results with simplified physics. *Quart. J. Roy. Meteor. Soc.*, **126**, 1143–1170.

- Reynolds, R. W., and T. M. Smith, 1994: Improved global sea surface temperature analyses using optimal interpolation. J. Climate, 7, 929–948.
- Rosati, A., K. Miyakoda, and R. Gudgel, 1997: The impact of ocean initial conditions on ENSO forecasting with a coupled model. *Mon. Wea. Rev.*, **125**, 754–772.
- Roullet, G., and G. Madec, 2000: Salt conservation, free surface, and varying levels: A new formulation for ocean general circulation models. J. Geophys. Res., 105, 23 927–23 942.
- Rutherford, I., 1972: Data assimilation by statistical interpolation of forecast error fields. J. Atmos. Sci., 29, 809–815.
- Segschneider, J., D. L. T. Anderson, and T. N. Stockdale, 2000: Towards the use of altimetry for operational seasonal forecasting. *J. Climate*, **13**, 3116–3138.
- —, —, J. Vialard, M. A. Balmaseda, and T. N. Stockdale, 2001: Initialization of seasonal forecasts assimilating sea level and temperature observations. J. Climate, 14, 4292–4307.
- Smith, N. R., J. E. Blomley, and G. Meyers, 1991: A univariate statistical interpolation scheme for subsurface thermal analyses in the tropical oceans. *Progress in Oceanography*, Vol. 28, Pergamon Press, 219–256.
- Talagrand, O., 1991: The use of adjoint equations in numerical modeling of the atmospheric circulation. Automatic Differentiation of Algorithms: Theory, Implementation, and Application, A. Griewank and G. F. Corliss, Eds., SIAM, 169–180.
- —, and P. Courtier, 1987: Variational assimilation of meteorological observations with the adjoint vorticity equation. I: Theory. *Quart. J. Roy. Meteor. Soc.*, **113**, 1311–1328.
- Tarantola, A., 1987: Inverse Problem Theory: Methods for Data Fitting and Model Parameter Estimation. Elsevier, 613 pp.
- Thacker, W. C., and R. B. Long, 1988: Fitting dynamics to data. J. Geophys. Res., 93, 1227–1240.
- Thépaut, J.-N., R. N. Hoffman, and P. Courtier, 1993: Interactions of dynamics and observations in four-dimensional variational assimilation. *Mon. Wea. Rev.*, **121**, 3393–3414.

- —, P. Courtier, G. Belaud, and G. LeMaître, 1996: Dynamical structure functions in a four-dimensional variational assimilation: A case study. *Quart. J. Roy. Meteor. Soc.*, **122**, 535–561.
- Tzipermann, E., W. C. Thacker, R. B. Long, and S.-M. Hwang, 1992a: Oceanic data analysis using a general circulation model. Part I: Simulations. J. Phys. Oceanogr., 22, 1434–1457.
- —, —, —, , and S. R. Rintoul, 1992b: Oceanic data analysis using a general circulation model. Part II: A North Atlantic model. J. Phys. Oceanogr., 22, 1458–1485.
- Vialard, J., C. Menkes, J.-P. Boulanger, P. Delecluse, E. Guilyardi, M. J. McPhaden, and G. Madec, 2001: A model study of oceanic mechanisms affecting equatorial Pacific sea surface temperature during the 1997–98 El Niño. J. Phys. Oceanogr., 31, 1649–1675.
  —, A. T. Weaver, D. L. T. Anderson, and P. Delecluse, 2003: Three-
- and four-dimensional variational assimilation with a general circulation model of the tropical Pacific Ocean. Part II: Physical validation. *Mon. Wea. Rev.*, **131**, 1379–1395.
- Weaver, A. T., and P. Courtier, 2001: Correlation modelling on the sphere using a generalized diffusion equation. *Quart. J. Roy. Meteor. Soc.*, **127**, 1815–1846.
- —, J. Vialard, D. L. T. Anderson, and P. Delecluse, 2002: Threeand four-dimensional variational assimilation with a general circulation model of the tropical Pacific Ocean. ECMWF Tech. Memo. 365, 74 pp. [Available online at http://www.ecmwf.int/ publications/.]
- Xu, Q., 1996: Generalized adjoint for physical processes with parametrized discontinuities. Part I: Basic issues and heuristic examples. J. Atmos. Sci., 53, 1123–1142.
- Zhu, J., and M. Kamachi, 2000: The role of time step size in numerical stability of tangent linear models. *Mon. Wea. Rev.*, **128**, 1562– 1572.
- Zou, X., 1997: Tangent-linear and adjoint of 'on-off' processes and their feasibility for use in 4-dimensional variational data assimilation. *Tellus*, **49A**, 3–31.

# A multivariate balance operator for variational ocean data assimilation

By A. T. WEAVER<sup>1\*</sup>, C. DELTEL<sup>2</sup>, E. MACHU<sup>1</sup>, S. RICCI<sup>1</sup> and N. DAGET<sup>1</sup>

<sup>1</sup>Centre Européen de Recherche et de Formation Avancée en Calcul Scientifique/SUC URA 1875,

Toulouse, France

<sup>2</sup>Laboratoire d'Océanographie et du Climat-Expérimentation et Approches Numériques/IPSL, Unité Mixte de Recherche 7159 CNRS/IRD/UPMC/MNHN, Paris, France

(Received 6 June 2005; revised 19 December 2005)

#### SUMMARY

It is common in meteorological applications of variational assimilation to specify the error covariances of the model background state implicitly via a transformation from model space where variables are highly correlated to a control space where variables can be considered to be approximately uncorrelated. An important part of this transformation is a balance operator which effectively establishes the multivariate component of the error covariances. The use of this technique in ocean data assimilation is less common. This paper describes a balance operator that can be used in a variable transformation for oceanographic applications of three- and four-dimensional variational assimilation. The proposed balance operator has been implemented in an incremental variational data assimilation system for a global ocean general-circulation model. Evidence that the balance operator can explain a significant percentage of background-error variance is presented. The multivariate analysis structures implied by the balance operator are illustrated using single-observation experiments.

KEYWORDS: Background-error covariances Nonlinear balance Ocean analysis

# 1. INTRODUCTION

The importance of the background-error covariances for determining the quality of analyses and forecasts is well known (Daley 1991). Specifying appropriate background-error covariances is a complex research problem which requires careful consideration of physical, statistical and computational issues. One important problem is how best to define the multivariate component of the background-error covariances. The multivariate component is responsible for transferring observational information between model variables and thus is critical for extracting information about unobserved variables from directly observed quantities. The problem of defining multivariate covariances is also intimately related to that of producing balanced initial conditions for initializing forecasts. In particular, improvements in the specification of multivariate covariances will usually translate into better dynamically balanced analyses and therefore can reduce, or even eliminate, the need for a separate 'initialization' procedure.

In oceanography, various approaches have been developed to introduce multivariate constraints in data assimilation systems. In some systems, they take the form of dynamical or physical constraints (e.g. geostrophic or temperature–salinity (T-S) relations) that are applied a posteriori to a statistically generated univariate analysis (Burgers *et al.* 2002; Troccoli *et al.* 2002; Balmaseda 2004). While this generally leads to much better forecasts than if no constraints were applied at all, it does not make optimal use of multivariate information in defining the analysis itself and makes the assimilation of different data types more difficult.

A more effective way of incorporating multivariate constraints in the data assimilation system is through the background-error covariances. A popular method in oceanographic applications of sequential data assimilation schemes such as the Kalman filter is to compute the error covariances in a reduced-dimension subspace spanned by a limited number of three-dimensional (3D) empirical orthogonal functions (Testut *et al.* 2003) or a few members of an appropriately generated ensemble of ocean

\* Corresponding author: CERFACS, 42 Avenue Gaspard Coriolis, 31057 Toulouse Cedex 1, France.

e-mail: weaver@cerfacs.fr

© Royal Meteorological Society, 2005.

model states (Lermusiaux *et al.* 2000; Keppenne and Rienecker 2003). While reducedspace methods are capable of producing complex multivariate covariance structures, they have the disadvantage of restricting the analysis increment to lie only in the subspace spanned by the chosen basis vectors. Various localization techniques have been proposed to overcome this rank deficiency problem but unfortunately can be applied only at the expense of disrupting some of the attractive balance properties of the original covariances.

The specification of the multivariate component of the background-error covariances for variational ocean data assimilation has received much less attention. In ocean applications of four-dimensional variational assimilation (4D-Var), crossvariable correlations in the background errors are often neglected altogether (Bennett *et al.* 2000; Stammer *et al.* 2002; Weaver *et al.* 2003). This approximation is often justified by the fact that 4D-Var includes the ocean model (or a linearized version of the ocean model) as a constraint in the assimilation problem and so already contains a multivariate component. The validity of this approximation depends on several factors such as the length of the assimilation window, the choice of control variables, and the particular application. It is clearly a very poor approximation, however, in threedimensional variational assimilation (3D-Var) which does not include the ocean model as a constraint. In general, a well-tuned multivariate background-error covariance model is beneficial to 4D-Var as well as 3D-Var.

This paper describes a very general method for incorporating multivariate constraints in variational ocean data assimilation. It extends the work of Ricci *et al.* (2005) who proposed a technique for incorporating T-S constraints in a 3D-Var system. The fundamental idea is to simplify the specification of the background-error covariances by designing a transformation from model state space, where variables are highly correlated, to another (control) space where variables can be considered approximately mutually uncorrelated. The basic technique is commonly used in meteorological applications of variational assimilation (Derber and Bouttier 1999; Cullen 2003) but has seen limited use in oceanography. In effect, the specification of the background-error covariances in model state space is transformed into one of defining a more general observation operator. An obvious advantage with this approach is that observation operators can be nonlinear whereas constraints that are included in traditional covariance (matrix) formulations are necessarily linear.

The paper is organized as follows. An outline of the general approach for modelling background-error covariances is given in section 2. Special attention is paid to some important practical issues concerning the implementation of the technique in incremental versions of 3D-Var and 4D-Var. Section 3 describes a multivariate balance operator that can be used in a control variable transformation for 3D-Var and 4D-Var applications with ocean general-circulation models (OGCMs). The proposed balance operator has been implemented in a variational assimilation system for the OPA OGCM. Examples with this system are presented in section 4 to illustrate various properties of the balance operator. Conclusions are given in section 5. An appendix provides some mathematical details on the relationship between the balance operator and the multivariate component of the background-error covariance matrix.

## 2. AN IMPLICIT REPRESENTATION OF THE BACKGROUND-ERROR COVARIANCES

## (a) Formulation of the problem

The formulation of variational assimilation given by Derber and Wu (1998) provides a very general and convenient framework for representing background-error

covariances in model state space. In their formulation, the variational analysis is defined by the minimization of a cost function of the form

$$J(\mathbf{v}) = \frac{1}{2}(\mathbf{v} - \mathbf{v}^{\mathrm{b}})^{\mathrm{T}}(\mathbf{v} - \mathbf{v}^{\mathrm{b}}) + \frac{1}{2}(G(\mathbf{v}) - \mathbf{y}^{\mathrm{o}})^{\mathrm{T}}\mathbf{R}^{-1}(G(\mathbf{v}) - \mathbf{y}^{\mathrm{o}}),$$
(1)

where  $\mathbf{v}$  is the control (analysis) vector,  $\mathbf{v}^{b}$  is the background estimate of the control vector,  $\mathbf{y}^{o}$  is the vector of observations, **R** is an estimate of the observation-error covariance matrix, including contributions from measurement and representativeness error, and G is a nonlinear operator that maps the control vector onto the space of the observation vector\*. The background-error covariance matrix of the control vector is assumed to be the identity matrix  $(\mathbf{B}_{(\mathbf{v})} = \mathbf{I})$  as evident by the use of the canonical inner product for the background term in (1). In other words, background errors for  $\mathbf{v}^{b}$  are assumed to be uncorrelated and to have unit variance. Clearly the control vector must be constructed carefully for this to be a reasonable assumption; e.g. it would be a very poor assumption if  $\mathbf{v}$  were taken to be the model state vector. There are two advantages that result from this formulation where the background term takes on a very simple form. First, it generally improves the convergence properties of the minimization when the problem is solved with a conjugate gradient algorithm. For quadratic cost functions, this is often explained by a reduction in the condition number of the Hessian (Golub and Van Loan 1996). Second, all constraints in the assimilation problem are now imposed through the nonlinear observation operator G, including multivariate and smoothness constraints that are used in conventional model-space (matrix) formulations of the background-error covariances. In particular, this opens the way for incorporating potentially more realistic (nonlinear) multivariate balance relationships in the analysis problem.

The control vector  $\mathbf{v}$  is assumed to be related to the model state vector  $\mathbf{x}$  through a transformation of the form

$$\mathbf{v} = U^{-1}(\mathbf{x}),\tag{2}$$

where  $U^{-1}$  is a block-matrix operator, with possibly nonlinear blocks, which is assumed to be square and invertible in the following. There is no complication if U is rectangular (i.e. if there are fewer control variables than state variables) but in this case U would only be invertible in a generalized sense. If the observations are distributed over a time window  $t_0 \le t_i \le t_n$  then **x** can be interpreted, as in a conventional formulation of 4D-Var, as the initial state of the dynamical model used in G to propagate the model state forward to the observation times<sup>†</sup>.

Following Derber and Bouttier (1999), the operator  $U^{-1}$  can be split into three basic operators: a transformation  $K^{-1}$  that produces a set of approximately mutually uncorrelated variables by removing any known dynamical or physical balance relationships between model state variables; a diagonal matrix  $D^{-1}$  of normalization factors; and a roughening operator  $F^{-1}$  (the inverse of a smoothing operator) that acts separately on each of the uncorrelated variables. The change of variables (2) is needed to compute the background estimate,  $\mathbf{v}^{b}$ , of the control vector from the background estimate,  $\mathbf{x}^{b}$ , of the model state, while the inverse of the change of variables

$$\mathbf{x} = U(\mathbf{v}) = K(D(F(\mathbf{v}))) \tag{3}$$

<sup>\*</sup> Nonlinear and linear matrix operators will be distinguished throughout by italic and bold font, respectively.

 $<sup>\</sup>dagger$  By interpreting **x** to be the initial conditions, the model and external forcing fields are tacitly assumed to be perfect. This assumption can be relaxed in the above formulation by considering **x** to contain model-error or external forcing terms in addition to the initial conditions.

is needed to evaluate the term  $G(\mathbf{v})$  in the observation term. Equation (2) can be used to compute covariance statistics of the contrived control vector  $\mathbf{v}$  from estimates of background error for the state vector  $\mathbf{x}$ . In practice, only a few aspects of the covariances (e.g. average variances) can be estimated reliably. From these estimates, the assumption that  $\mathbf{B}_{(\mathbf{v})} \approx \mathbf{I}$  can be tested: if it is not well satisfied then either a new  $\mathbf{B}_{(\mathbf{v})} \neq \mathbf{I}$  could be used to weight the background term in (1) or the parameters in the operators F, D and K could be recalibrated so that the approximation is better satisfied.

## *(b) Incremental formulation*

The incremental formulation (Courtier *et al.* 1994) provides a practical algorithm for approximately minimizing (1). The incremental algorithm is defined by the iterative minimization of a sequence,  $k = 1, ..., K_0$ , of quadratic cost functions

$$J^{k}(\delta \mathbf{v}^{k}) = \frac{1}{2} (\delta \mathbf{v}^{k} - \mathbf{d}_{(\mathbf{v})}^{\mathbf{b},k})^{\mathrm{T}} (\delta \mathbf{v}^{k} - \mathbf{d}_{(\mathbf{v})}^{\mathbf{b},k}) + \frac{1}{2} (\mathbf{G}^{k-1} \delta \mathbf{v}^{k} - \mathbf{d}^{\mathbf{o},k})^{\mathrm{T}} \mathbf{R}^{-1} (\mathbf{G}^{k-1} \delta \mathbf{v}^{k} - \mathbf{d}^{\mathbf{o},k}),$$
(4)

where

$$\mathbf{d}_{(\mathbf{v})}^{\mathbf{b},k} = \mathbf{v}^{\mathbf{b}} - \mathbf{v}^{k-1},\tag{5}$$

$$\mathbf{d}^{\mathbf{o},k} = \mathbf{y}^{\mathbf{o}} - G(\mathbf{v}^{k-1}) \tag{6}$$

is the innovation vector,  $\mathbf{v}^{k-1}$  is a reference state,  $\delta \mathbf{v}^k$  is an increment defined by  $\mathbf{v}^k = \mathbf{v}^{k-1} + \delta \mathbf{v}^k$ , and  $\mathbf{G}^{k-1}$  is a linearized operator defined such that  $G(\mathbf{v}^{k-1} + \delta \mathbf{v}^k) \approx G(\mathbf{v}^{k-1}) + \mathbf{G}^{k-1}\delta \mathbf{v}^k$  (when this equation is satisfied exactly, (4) is identical to (1)). The superscript k - 1 indicates that  $\mathbf{G}^{k-1}$  is the result of linearizing G about  $\mathbf{v}^{k-1}$ . The sequence  $k = 1, \ldots, K_0$  is called outer iterations while the minimization iterations performed within each outer loop are called inner iterations. Equations (5) and (6) are the effective 'background' and 'observation' vectors for the inner-loop minimization. In practice, it is customary to set  $\mathbf{v}^0 = \mathbf{v}^b$  and to choose  $\mathbf{v}^{k-1}$ , for  $k = 2, \ldots, K_0$ , to be the solution obtained at the end of the previous outer loop. The minimum of (4) after the  $K_0$ th outer iteration defines the analysis increment,  $\delta \mathbf{v}^a = \delta \mathbf{v}^{K_0}$ . The analysis in model space is then given by  $\mathbf{x}^a = U(\mathbf{v}^a)$  where  $\mathbf{v}^a = \mathbf{v}^{K_0-1} + \delta \mathbf{v}^a$ .

The nonlinear transformation (3) is needed on each outer iteration to evaluate the term  $G(\mathbf{v}^{k-1})$  in (4). Through successive linearizations about  $\mathbf{v}^l$ , l = 0, ..., k - 2, this transformation can be approximated by

$$\mathbf{x}^{k-1} = U(\mathbf{v}^{k-1}) \approx U(\mathbf{v}^0) + \sum_{l=1}^{k-1} \mathbf{U}^{l-1} \delta \mathbf{v}^l.$$
(7)

By choosing  $\mathbf{v}^0 = \mathbf{v}^b$ , the first term on the right-hand side of (7) becomes

$$U(\mathbf{v}^0) = U(\mathbf{v}^b) \equiv \mathbf{x}^b.$$

Equation (7) then implies that  $\mathbf{x}^{k-1}$  can be approximated as the sum of the model-space background state and the model-space increments estimated using the inverse of the *linearized* change of variables. A further consequence of choosing  $\mathbf{v}^0 = \mathbf{v}^b$  is that the difference vector (5) can be written as minus the sum of the increments generated from previous outer iterations:

$$\mathbf{d}_{(\mathbf{v})}^{\mathbf{b},k} = \mathbf{v}^0 - \mathbf{v}^{k-1} = -\sum_{l=1}^{k-1} \delta \mathbf{v}^l.$$
(8)

Equation (8) together with the approximation (7) allow us to eliminate the explicit dependence of (4) on  $\mathbf{v}^{k-1}$  and thus to iterate the incremental minimization algorithm without the need to perform either the nonlinear transformation (2) or its inverse (3) (only the linearized transformations are required).

To complete the evaluation of  $G(\mathbf{v}^{k-1})$ ,  $\mathbf{x}^{k-1}$  must be propagated to the observation times using the model operator and then transformed to the observed quantities using the observation operator. The linearized counterpart of this operator is required to evaluate  $\mathbf{G}^{k-1}\delta\mathbf{v}^k$  in (4). As discussed in Weaver *et al.* (2003), 3D-Var and 4D-Var can be distinguished by the type of linear model that is used to evolve the increments between observation times. In 3D-Var the increments are persisted, whereas in 4D-Var they are evolved by a dynamical model that closely approximates the tangent-linear model. By distinguishing 3D-Var and 4D-Var at the incremental level, they can be viewed as two different algorithms for approximately solving the same 4D assimilation problem described by the non-quadratic cost function (1).

## (c) Diagnosing the effective background-error covariance matrix

Although the background-error covariance matrix in model space has not been defined explicitly, its effective form on a given outer iteration can be easily diagnosed by transforming the background term in (4) into model space using the linearized change of variables  $\delta \mathbf{v}^k = (\mathbf{U}^{k-1})^{-1} \delta \mathbf{x}^k$  and its inverse  $\delta \mathbf{x}^k = \mathbf{U}^{k-1} \delta \mathbf{v}^k$ . This yields

$$J_{b}^{k} = \frac{1}{2} (\delta \mathbf{x}^{k} - \mathbf{d}_{(\mathbf{x})}^{b,k})^{\mathrm{T}} (\mathbf{B}_{(\mathbf{x})}^{k})^{-1} (\delta \mathbf{x}^{k} - \mathbf{d}_{(\mathbf{x})}^{b,k}),$$
(9)

where

$$\mathbf{B}_{(\mathbf{x})}^{k} = \underbrace{\mathbf{K}^{k-1} \mathbf{D}_{(\widehat{\mathbf{x}})}^{k-1} \mathbf{F}^{k-1}}_{\mathbf{U}^{k-1}} \underbrace{\mathbf{F}^{k-1} \mathbf{D}_{(\widehat{\mathbf{x}})}^{k-1} \mathbf{K}^{k-1}}_{\mathbf{U}^{k-1}},$$
(10)

and  $\mathbf{d}_{(\mathbf{x})}^{\mathbf{b},k} = \mathbf{U}^{k-1} \mathbf{d}_{(\mathbf{v})}^{\mathbf{b},k}$ . Equation (10) corresponds to the model background-error covariance matrix on the *k*th outer iteration. Since  $\mathbf{B}_{(\mathbf{x})}^k$  depends, in general, on the linearization state  $\mathbf{x}^{k-1}$ , it may vary from one outer iteration to the next. In this way, the outer iterations provide an adaptive mechanism for modifying the background-error covariance model during the course of minimization. The background-error covariance matrix  $\mathbf{B}_{(\mathbf{x})}^{K_0}$  used on the final outer iteration  $K_0$  would be the effective covariance matrix used for the analysis. Note that in 4D-Var,  $\mathbf{B}_{(\mathbf{x})}^{K_0}$  would be evolved (implicitly) within the assimilation window through the action of the linearized dynamical model and its adjoint (Courtier *et al.* 1994). In 3D-Var, on the other hand,  $\mathbf{B}_{(\mathbf{x})}^{K_0}$  would be fixed within the assimilation window, although, as in 4D-Var, it may vary from one assimilation cycle to the next through its dependence on the background state.

Equation (10) provides a valuable statistical interpretation of the control variable transformation. The product  $\mathbf{F}^{k-1}(\mathbf{F}^{k-1})^{\mathrm{T}}$  of the linearized smoothing matrix and its transpose can be interpreted as a correlation matrix, provided that care has been taken to normalize the matrix so that the diagonal elements are all equal to one. The correlations in  $\mathbf{F}^{k-1}(\mathbf{F}^{k-1})^{\mathrm{T}}$  correspond to those of the errors of the transformed background variables  $\mathbf{\hat{x}}^{\mathrm{b}} = K^{-1}(\mathbf{x}^{\mathrm{b}})$ , not to the error correlations of  $\mathbf{x}^{\mathrm{b}}$  itself. By construction, cross-correlations between these variables are neglected so that  $\mathbf{F}^{k-1}(\mathbf{F}^{k-1})^{\mathrm{T}}$  is block-diagonal (univariate), where each block corresponds to the autocorrelation matrix for each variable in  $\mathbf{\hat{x}}^{\mathrm{b}}$ . While the cross-correlations will never be exactly zero in practice,

the intent is that, with an astutely chosen  $K^{-1}$  operator, they can be made sufficiently small so that neglecting them is an acceptable assumption.

The diagonal matrix  $\mathbf{D}_{(\widehat{\mathbf{x}})}^{k-1}$  in (10) contains estimates of the standard deviations of the errors in  $\widehat{\mathbf{x}}^{b}$ . In meteorology, it is typical to estimate the standard deviations (and parameters of the correlation model) from a suitably constructed ensemble of forecast differences (Parrish and Derber 1992; Buehner 2005). To estimate statistics of the control variables, the forecasts must first be transformed into  $\widehat{\mathbf{x}}$ -space using  $K^{-1}$ or, as an approximation, the forecast differences can be transformed directly using the linearized balance operator  $\mathbf{K}^{k-1}$ . In (10),  $\mathbf{K}^{k-1}$  couples the different model variables and thus establishes the multivariate component of the background-error covariances in  $\mathbf{x}$ -space (Derber and Bouttier 1999). The remainder of this article is devoted to the specification of a balance operator for ocean data assimilation. The problems of estimating background-error covariances and deriving efficient and general smoothing algorithms for representing background-error correlations are both very important, but a proper discussion of these issues goes beyond the scope of this paper.

## 3. A BALANCE OPERATOR FOR OCEAN STATE VARIABLES

# (a) General formulation

The variables comprising the model state vector are assumed to be potential temperature T, salinity S, sea surface height (SSH)  $\eta$ , and the components of the horizontal velocity vector  $\mathbf{u}^{h} = (u, v)^{T}$ . These variables correspond to the standard prognostic variables in a free-surface, hydrostatic OGCM. In this section, an operator  $K^{-1}$  is developed which can be used to transform  $\mathbf{x} = (T, S, \eta, \mathbf{u}^{h})^{T}$  into a vector  $\hat{\mathbf{x}} = (T, S_{U}, \eta_{U}, \mathbf{u}_{U}^{h})^{T}$  whose elements  $T, S_{U}, \eta_{U}$  and  $\mathbf{u}_{U}^{h} = (u_{U}, v_{U})^{T}$  can be considered to be approximately mutually uncorrelated. This can be achieved by separating the state variables into *unbalanced* and *balanced* components (Derber and Bouttier 1999), except for one variable, taken here to be T, which is treated in totality and used as the starting point to establish the balanced part of the other variables. The other elements  $S_{U}, \eta_{U}$  etc. of  $\hat{\mathbf{x}}$  represent the unbalanced part of that particular variable.

The balance relationships used to define  $\mathbf{x} = K(\hat{\mathbf{x}})$  are described in detail in the next section. Symbolically, the balance operator can be summarized by the sequence of equations

$$T = T, 
S = K_{ST}(T) + S_{U} = S_{B} + S_{U}, 
\eta = K_{\eta\rho}(\rho) + \eta_{U} = \eta_{B} + \eta_{U}, 
u = K_{up}(p) + u_{U} = u_{B} + u_{U}, 
v = K_{vp}(p) + v_{U} = v_{B} + v_{U},$$
(11)

where

$$\rho = K_{\rho TS}(T, S),$$
  

$$p = K_{p\rho}(\rho) + K_{p\eta}(\eta)$$
(12)

are diagnostic quantities corresponding to density and pressure, respectively, and  $K_{xy}$  represents the transformation from the variable(s) y to x. The variables with a subscript B on the right-hand side of (11) represent the balanced component of those variables. The lower block-triangular structure of the balance operator (11) implies that a balanced variable can be a function of the variables preceding it in the sequence but

will be independent of the variables following it in the sequence. It also allows the inverse balance operator  $K^{-1}$  to be obtained trivially from the sequence of equations

$$T = T,$$
  

$$S_{U} = S - S_{B},$$
  

$$\eta_{U} = \eta - \eta_{B},$$
  

$$u_{U} = u - u_{B},$$
  

$$v_{U} = v - v_{B}.$$
  
(13)

A linearized version of  $\mathbf{x} = K(\hat{\mathbf{x}})$  is required for the incremental formulation. According to (7), the linearized balance equation can be approximated as

$$\mathbf{x}^{k-1} \approx \mathbf{x}^{\mathbf{b}} + \sum_{l=1}^{k-1} \delta \mathbf{x}^{l}, \tag{14}$$

where  $\delta \mathbf{x}^l = \mathbf{K}^{l-1} \delta \mathbf{\hat{x}}^l$ . This approximation is very convenient since it eliminates the need to specify the nonlinear version of the balance operator (it is implicit in  $\mathbf{x}^b$ ). It is used in the rest of this section and in the illustrations presented in section 4. It can be expected to be a good approximation when the balance operator is weakly nonlinear.

From (11) and (12), the linear balance equations for the increment can be written in the general form

$$\delta T^{k} = \delta T^{k},$$

$$\delta S^{k} = \mathbf{K}_{ST}^{k-1} \delta T^{k} + \delta S_{U}^{k} = \delta S_{B}^{k} + \delta S_{U}^{k},$$

$$\delta \eta^{k} = \mathbf{K}_{\eta\rho} \delta \rho^{k} + \delta \eta_{U}^{k} = \delta \eta_{B}^{k} + \delta \eta_{U}^{k},$$

$$\delta u^{k} = \mathbf{K}_{up} \delta p^{k} + \delta u_{U}^{k} = \delta u_{B}^{k} + \delta u_{U}^{k},$$

$$\delta v^{k} = \mathbf{K}_{vp} \delta p^{k} + \delta v_{U}^{k} = \delta v_{B}^{k} + \delta v_{U}^{k},$$
(15)

where

$$\delta \rho^{k} = \mathbf{K}_{\rho T}^{k-1} \delta T^{k} + \mathbf{K}_{\rho S}^{k-1} \delta S^{k}, \delta p^{k} = \mathbf{K}_{p\rho} \delta \rho^{k} + \mathbf{K}_{p\eta} \delta \eta^{k}.$$
(16)

As described later, nonlinear operators are used for the salinity balance  $K_{ST}$  and the density balance  $K_{\rho TS}$ . All the other balance operators are linear and thus independent of the linearization state  $\mathbf{x}^{k-1}$ . This has been made clear in (15) by omitting the superscript k - 1 from those matrix operators.

# (b) A set of linearized balance relationships

Temperature plays an important role in the balance formulation since is used to compute all, or a significant part of, the balanced component of the other variables. The relationship between temperature T and salinity S is complex and traditionally determined empirically from scatter plots of historical T and S data. Han *et al.* (2004) propose fitting a high-order polynomial function to T-S diagrams in order to determine an explicit S(T) relationship. This procedure works reasonably well in some data-rich regions such as the western tropical Pacific Ocean, as illustrated by Han *et al.* (2004). Using a somewhat different formulation to the one presented here, Han *et al.* (2004) then go on to show how such an S(T) relationship, together with an estimate of the uncertainty in this relationship, can be used to correct salinity from temperature data within a variational assimilation framework.

Troccoli and Haines (1999) propose an alternative and simpler method for adjusting salinity when only temperature information is available. Their approach is designed to preserve the T-S properties of the background state by making vertical displacements of the local background salinity field in response to changes to the local background temperature field produced by the assimilation of temperature data. The attractive features of their method are that it can be applied in a global system, it allows for state-dependency in the T-S relation, and it does not require any prior statistical analysis of an observational database.

Ricci *et al.* (2005) describe a simple variant of the Troccoli and Haines (1999) scheme for implementation within a linear balance operator. In their study, balanced salinity increments are defined by

$$\delta S_{\rm B}^{k} = \gamma^{k-1} \frac{\partial S}{\partial z} \bigg|_{S=S^{k-1}} \frac{\partial z}{\partial T} \bigg|_{T=T^{k-1}} \delta T^{k}, \tag{17}$$

where  $\gamma = \gamma^{k-1}(T^{k-1}, S^{k-1}, u^{k-1}, v^{k-1})$  is a coefficient that is set to either zero or one, depending on various conditions in the reference state. For example, to take into account the weak correlation between temperature and salinity in well-mixed regions,  $\gamma^{k-1}$  is set to zero at grid points where the reference vertical mixing coefficient is large, such as in the ocean mixed layer. When  $\gamma^{k-1} = 0$ ,  $\delta S^k$  is entirely described by its unbalanced component  $\delta S_{U}$ . To avoid a discontinuity in the balance at the base of the mixed layer,  $\delta S_{\rm B}^{k}$  is smoothly reduced to zero at the surface from its non-zero value just below the mixed layer. The vertical derivatives in (17) are used to estimate the local derivative of the background T-S relation and can be computed using finite differences or a cubic spline. In practice, it has been found desirable to apply a horizontal smoothing operator (e.g. the one used in  $\mathbf{F}^{k-1}$ ) to the balance coefficient in (17) in order to avoid generating noisy salinity increments. The impact of the T-S balance has been evaluated in detail by Ricci et al. (2005) in a multi-annual cycled 3D-Var experiment for the tropical Pacific Ocean. When assimilating temperature data alone, they showed that the constraint can have a significant positive impact on velocity as well as salinity compared to a 3D-Var analysis in which no T-S constraint is applied. Notice that, as the T-S constraint is dependent on the reference state, it can evolve both during the course of minimization (via outer iterations) and from one assimilation cycle to the next.

Density can be computed from potential temperature and salinity using a nonlinear equation of state (e.g. McDougall *et al.* 2003). The (balanced) density increment can be defined by linearizing the equation of state about the reference state:

$$\delta \rho^k = \rho_0(-\alpha^{k-1}\delta T^k + \beta^{k-1}\delta S^k), \tag{18}$$

where

$$\alpha^{k-1} = (1/\rho_0) \partial \rho / \partial T |_{S=S^{k-1}, T=T^{k-1}},$$
  
$$\beta^{k-1} = (1/\rho_0) \partial \rho / \partial S |_{S=S^{k-1}, T=T^{k-1}}$$

are thermal and saline expansion coefficients, respectively, and  $\rho_0$  is a constant reference density.

SSH can be computed diagnostically as a function of the state variables T, S and  $\mathbf{u}^{h}$  by filtering out non-stationary contributions to SSH (e.g. from high-frequency gravity waves) using the rigid-lid approximation (Fukumori *et al.* 1998). Furthermore, for flow regimes where the Rossby number is weak (regimes close to geostrophic balance), contributions to SSH from advection, dissipation, and surface forcing can be neglected.

The resulting equation approximates SSH as the sum of two terms: a baroclinic term that depends on density and a barotropic term that depends on the depth-integrated transport. For the global model used in the illustrations in the next section, Ferry (2003) has demonstrated that SSH variability is indeed dominated by its baroclinic and barotropic components, except in coastal regions where the contribution from surface forcing can be important. In the following, the baroclinic and barotropic contributions to SSH are taken to be the balanced and unbalanced parts of SSH, respectively.

The balanced (baroclinic) component can be estimated by computing the dynamic height at the surface z = 0 relative to a reference depth  $z_{ref}$ :

$$\delta \eta_{\rm B}^k = -\int_{z'=z_{\rm ref}}^0 (\delta \rho^k(z')/\rho_0) \, \mathrm{d}z' \tag{19}$$

 $(z_{ref} = 1500 \text{ m in the examples in section 4})$ . Equation (19) is only an approximation of the baroclinic part of (the increment of) SSH. The complete expression involves the solution of an elliptic equation (Fukumori *et al.* 1998):

$$\nabla \cdot H \nabla \delta \eta_{\rm B}^k = -\nabla \cdot \int_{z=-H}^0 \int_{z'=z}^0 (\nabla \delta \rho^k(z')/\rho_0) \, \mathrm{d}z' \, \mathrm{d}z \tag{20}$$

where  $z = -H(\lambda, \phi)$  is the total ocean depth,  $\lambda$  is longitude,  $\phi$  is latitude, and  $\nabla$  and  $\nabla$ · are the horizontal gradient and divergence operators, respectively. Equation (20) takes into account variations in topography and is independent of a reference depth, and therefore would be more accurate than (19) in regions where bathymetry is important or where the ocean is shallow. For this study, however, the simpler equation (19) has been adopted.

The balanced pressure increment at any depth z can be computed by integrating the hydrostatic equation from z to the surface:

$$\delta p^{k}(z) = \int_{z'=z}^{0} \delta \rho^{k}(z') g \, \mathrm{d}z' + \rho_{0} g (\delta \eta_{\mathrm{B}}^{k} + \delta \eta_{\mathrm{U}}^{k}), \tag{21}$$

where g is the acceleration due to gravity, and the second term on the right-hand side of (21) is the pressure exerted by the surface elevation,  $\delta p^k(0) = \rho_0 g \delta \eta^k$ , with  $\delta \eta^k$  given by (15). Substituting (19) in (21) and reversing the order of integration of the first term on the right-hand side of (21) leads to

$$\delta p^k(z) = -\int_{z'=z_{\text{ref}}}^z \delta \rho^k(z') g \, \mathrm{d}z' + \rho_0 g \delta \eta_{\mathrm{U}}^k. \tag{22}$$

Away from the equator, the balanced part of the horizontal velocity components  $(\delta u_B^k, \delta v_B^k)$  is assumed to be in geostrophic balance; i.e. proportional to the horizontal gradient of (22) divided by the Coriolis parameter f. The horizontal gradient of the first term in (22) is associated with a *baroclinic* geostrophic velocity, while that of the second term is associated with a *barotropic* geostrophic velocity. The ageostrophic components of the velocity increment are assumed to be associated with the unbalanced components  $(\delta u_{II}^k, \delta v_{II}^k)$ .

Special treatment of the geostrophic velocity balance is required near the equator where  $f \rightarrow 0$ . There, the zonal component  $\delta u_{\rm B}^k$  is taken to be geostrophically balanced while the meridional component  $\delta v_{\rm B}^k$  is reduced to zero. Geostrophic balance for  $\delta u_{\rm B}^k$ is computed near the equator using a  $\beta$ -plane geostrophic approximation (Lagerloef *et al.* 1999), which involves the meridional derivative of the geostrophic equation. For this balance to exist, the meridional pressure gradient must exactly vanish at the equator so that the standard (undifferentiated) form of the geostrophic equation is satisfied when f = 0. Picaut and Tournier (1991) suggest adding a latitudinally dependent correction term to the pressure field in order to force a zero meridional gradient at the equator while leaving the meridional curvature of the original pressure field, and hence the estimate of the zonal geostrophic current via the  $\beta$ -plane approximation, unaltered. A similar technique is adopted here. The correction term effectively filters out all flows with anti-symmetric pressure structures about the equator. An important exception is an equatorial Kelvin wave, which is associated with a strictly zonal current in geostrophic balance and is thus described by the proposed velocity balance on the equator.

To allow for a smooth transition between the equatorial ( $\beta$ -plane) geostrophic velocity and the standard (f-plane) geostrophic velocity away from the equator, weighting functions  $W_{\beta} = \exp(-\phi^2/2L_{\beta}^2)$  and  $W_f = 1 - W_{\beta}$  are introduced, where  $L_{\beta}$  is a length-scale whose size is of the order of the equatorial Rossby radius of deformation (Lagerloef *et al.* 1999). At the equator,  $W_{\beta} = 1$  and  $W_f = 0$ , while far away from the equator,  $W_{\beta} \approx 0$  and  $W_f \approx 1$ . Experimental evidence is given by Lagerloef *et al.* (1999) to justify the Gaussian form for the weighting function. The complete expression for the increments of the balanced velocity components in spherical coordinates is then given by

$$\delta u_{\rm B}^{k} = -\frac{1}{\rho_0} \left( \frac{W_f}{f} + \frac{W_{\beta}}{\beta} \frac{1}{a} \frac{\partial}{\partial \phi} \right) \frac{1}{a} \frac{\partial \delta \widetilde{p}^{k}}{\partial \phi}, \tag{23}$$

$$\delta v_{\rm B}^k = \frac{1}{\rho_0} \frac{W_f}{f} \frac{1}{a \cos \phi} \frac{\partial \delta \widetilde{p}^k}{\partial \lambda},\tag{24}$$

where  $\beta = \partial f / \partial (a\phi)$ , and *a* is the radius of the earth. To simplify the  $\beta$ -plane approximation, the differentiated term involving the product  $f \partial \delta u_{\rm B}^k / \partial (a\phi)$  has been neglected in (23). This term can be expected to be relatively small near the equator where  $f \approx 0$ . Following Picaut and Tournier (1991), the modified pressure increment in (23) and (24) is defined by

$$\delta \widetilde{p}^{k} = \delta p^{k} - \phi \left(\frac{\partial \delta p^{k}}{\partial \phi}\right)_{\phi=0} \exp\left(-\phi^{2}/2L_{p}^{2}\right), \tag{25}$$

where the second term on the right-hand side (25) corresponds to the pressure correction factor. The correction term does not affect (24), is negligible far from the equator, and satisfies both the  $\beta$ -plane constraint

$$(\partial^2 \delta \widetilde{p}^k / \partial \phi^2)_{\phi=0} = (\partial^2 \delta p^k / \partial \phi^2)_{\phi=0}$$

and the necessary condition for geostrophic balance at the equator,

$$(\partial \delta \widetilde{p}^k / \partial \phi)_{\phi=0} = 0.$$

The length-scales  $L_p$  in the correction term and  $L_\beta$  in the weighting functions are taken to be equal and set to 1.55° as in Lagerloef *et al.* (1999).

## 4. Illustrations

The balance operator described in the previous section has been implemented in a 3D-Var/4D-Var system (Weaver *et al.* 2003) for a global, free-surface version of the OPA OGCM (Madec *et al.* 1998; Roullet and Madec 2000). The system has been exploited within the framework of the European ENACT project (see http://www.ecmwf.int/research/EU\_projects/ENACT) to produce global ocean re-analyses using historical temperature and salinity data (Ingleby and Huddleston 2006) and surface forcing fields from the ERA40 atmospheric re-analysis (Uppala *et al.* 2005). It is beyond the scope of this paper to provide a thorough description of the system and assessment of the re-analyses. Only certain aspects of the system concerning the balance operator and background-error covariance formulation are discussed and illustrated here.

## (a) Evidence of balance in ocean background errors

The balance operator can be considered effective if the variance of background error of the balanced variables explains a substantial part of the variance of background error of the full variables. If this is not the case then the balance operator would provide little useful information for the analysis. Since actual background error is unknown, a suitable proxy must be defined in order to estimate its statistical properties. Here, background error is approximated as the difference between the background state ( $\mathbf{x}^{b}(t_{n})$ ) and the reference state ( $\mathbf{x}^{K}(t_{n})$ ) at the end of an assimilation window. The two states will differ since the background state over the window  $t_{0} \le t_{i} \le t_{n}$  is obtained by forcing the model with the atmospheric fluxes only, whereas the reference state is obtained by assimilating data over  $t_{0} \le t_{i} \le t_{n}$  in addition to applying the surface forcing. This approach is analogous to the so-called NMC method used in meteorology to estimate backgrounderror statistics (Parrish and Derber 1992). Berre *et al.* (2006) discuss the conditions for which the NMC method is a good approximation to true forecast error.

A set of 328 background-minus-reference state differences has been obtained by cycling the 3D-Var system over the nine-year period 1993–2001 using a ten-day window. The inverse of the linearized balance operator was then applied to each of these difference fields in order to retrieve the unbalanced components. From the average variance,  $\sigma_x^2$ , of the full fields and the average variance,  $\sigma_{xU}^2$ , of the unbalanced fields, the percentage ratio of explained variance  $r = (1 - \sigma_{xU}^2 / \sigma_x^2) \times 100\%$  was computed. For the global average, r is 37% for salinity, 94% for SSH, and 70% and 44% for the zonal and meridional components of velocity, respectively. Errors computed using the NMC method were artificially small in regions poleward of 65°N/S and below 1000 m since no data were assimilated there. Those regions were thus excluded from the global average. In so far as the NMC method provides a reasonable representation of background errors, these results suggest that the proposed balance operator can explain a substantial amount of actual background-error variance.

The percentage variance ratio has also been computed as a function of depth from the horizontally averaged variances in each model level, and as a function of latitude from the zonally and depth-averaged variances. The results are displayed in Figs. 1(a) and (b). For salinity, r is largest (up to 80%) below the level of the mean thermocline (below 200 m), and reduces gradually to zero between 200 m and the surface (Fig. 1(a)). The small value of r close to the surface is understandable since the salinity balance is deliberately reduced in regions of strong mixing such as the surface mixed layer. Figure 1(b) suggests that the salinity balance is most effective in the subtropical gyre regions (between 10°N/S and 30°N/S). For the u-component of velocity, r is relatively uniform with depth, with values between 60% and 70%. The explained variance is about 20% to 40% smaller for v than u, and decreases more rapidly with depth. The velocity balance for the u-component is effective at all latitudes, even at the equator where it



Figure 1. The percentage ratio r of background-error variance explained by the balanced part of salinity (solid curve), SSH (dashed-dotted curve), the *u*-component of velocity (dashed curve) and the *v*-component of velocity (dotted curve). (a) r computed from horizontally averaged variances and plotted as a function of depth;
(b) r computed from zonally and depth-averaged variances and plotted as a function of latitude. Background errors have been estimated from a set of 329 background-minus-reference state differences.

explains about 50% of the variance. The value of r for the v-component is, as expected, small near the equator where the weight given to the geostrophic equation for v is reduced to zero (see (24)), but is comparable to the value of r for the u-component poleward of about 10°N/S. The SSH balance is particularly effective and explains over 90% of the variance within 40° of the equator. At mid- and high latitudes, the barotropic (unbalanced) component is known to be important, which probably explains the reduction in r in this region.

It is worth noting that if the balanced and unbalanced fields were truly independent then  $\sigma_x^2 = \sigma_{x_B}^2 + \sigma_{x_U}^2$ , where  $\sigma_{x_B}^2$  denotes the variance of the balanced field. The percentage variance ratio r would thus be equivalent to  $\hat{r} = (\sigma_{x_B}^2 / \sigma_x^2) \times 100\%$ . Comparing  $\hat{r}$  with r would then provide a measure of the validity of the assumption that the two fields are approximately uncorrelated. In particular, comparing Figs. 1(a) and (b) with the equivalent figures for  $\hat{r}$  (not shown) illustrates that r and  $\hat{r}$  do have similar structure and amplitude for all fields, except for u and v for which there is a tendency for  $\hat{r}$  to increase with depth rather than decrease with depth as in Fig. 1(a). The reasons for this discrepancy are not known at present but will need to be explored in future work.

## (b) Single-observation experiments with 3D-Var and 4D-Var

The multivariate properties of the background-error formulation are most clearly illustrated using single-observation experiments. A mathematical demonstration of this point is given in the appendix. For simplicity, the unbalanced components of salinity, SSH and velocity are ignored (they are assumed to have zero error variance) so that only the univariate T covariances need to be specified. In other words, the balance operator is applied as a strong constraint (Lorenc 2003). This is sufficient to illustrate basic properties of the balance operator which is the objective here. For practical applications, however, it would be better to apply the balance operator as a weak constraint by prescribing a non-zero covariance to the unbalanced components, provided reasonable estimates of these covariances can be computed (e.g. using ensemble methods).

The univariate 3D smoothing operator for *T* is defined as the product of a 1D and 2D anisotropic diffusion operator (Weaver and Courtier 2001). The resulting correlation structures are approximately Gaussian. The parameters of the 3D diffusion operator are the same as those used for the T-T correlations in the study of Weaver *et al.* (2003), except for the vertical correlation scales which have been slightly reduced here. The error variances,  $(\sigma_T^k)^2$ , for *T* have been made dependent on the vertical gradient of the reference *T* field in order to focus the largest errors at the level of the thermocline where thermal variability is greatest. Weaver *et al.* (2003) illustrate how this simple parametrization of the background *T* errors can account for some of the dynamical effects implicit in a Kalman filter. A similar parametrization is used in the operational ocean data assimilation systems at the National Centers for Environmental Prediction (Behringer *et al.* 2004). To avoid prescribing unrealistically small variances in the mixed layer and deep ocean where vertical *T* gradients are small, the parametrization is modified so that

$$\sigma_T^k = \begin{cases} \max(\tilde{\sigma}_T^k, \sigma_T^{\text{ml}}), & \text{in the mixed layer,} \\ \max(\tilde{\sigma}_T^k, \sigma_T^{\text{do}}), & \text{below the mixed layer,} \end{cases}$$
(26)

where

$$\tilde{\sigma}_T^k = \min\{|(\partial T/\partial z|_{T=T^{k-1}})\delta z|, \, \sigma_T^{\max}\},\tag{27}$$

 $\sigma_T^{\text{max}}$  being the maximum-allowed value of  $\sigma_T^k$ ,  $\delta_Z$  a vertical scale, and  $\sigma_T^{\text{ml}}$  and  $\sigma_T^{\text{do}}$  lower bounds in the mixed layer and deep ocean, respectively. The specification of  $\sigma_T^k$  is thus transformed into one of choosing appropriate values for these parameters. For the examples presented here,  $\sigma_T^{\text{max}} = 1.5 \text{ K}$ ,  $\delta_Z = 10 \text{ m}$ ,  $\sigma_T^{\text{ml}} = 0.5 \text{ K}$ , and  $\sigma_T^{\text{do}} = 0.07 \text{ K}$ . In the first example, the impact of a single *T* observation in 3D-Var is considered.

For this special case, the analysis of the T field depends entirely on the univariate Tcovariances (it is independent of the balance operator) and the analysis increments for the other variables are independent of the univariate covariances of their unbalanced component. Those increments could be obtained a posteriori by applying the linearized balance operator directly to the analysed T increment. This point is clarified in the appendix. Figure 2 shows the 3D-Var analysis increment for a single T observation chosen to be 1 K higher than the background T value, and located in the thermocline (100 m) on the equator in the central Pacific (160°W). The observation-error variance has been set to  $(1.0 \text{ K})^2$ . These increments are proportional to the implicitly defined background-error covariances with T at the observation point ( $K_0 = 1$  in all experiments). The structures are physically sensible. The positive T anomaly in the subsurface (Fig. 2(a)) is associated with an elevated SSH (Fig. 2(e)) and a geostrophic current at the surface, with an eastward zonal component that is symmetric about the equator (Figs. 2(c)) and a meridional component that is asymmetric about the equator (Figs. 2(d)). The dependence of the T-S balance and the T error variances on the reference state can lead to an anisotropic response in the T and S increments. To avoid generating noisy increments, both the T-S balance coefficients and  $\sigma_T^{b}$  were smoothed in each level using the horizontal diffusion operator in  $\mathbf{F}^{k-1}$ .

The previous example does not illustrate the full potential of the balance operator for exploiting different observation types in the assimilation process. When information about state variables other than T is assimilated, the analysis results from a generally complex interaction between the balance operator, its adjoint and the covariance statistics of the uncorrelated variables. For example, a SSH observation would provide direct



Figure 2. Horizontal sections of the analysis increments for (a) temperature, (b) salinity, (c) zonal velocity, (d) meridional velocity, and (e) SSH generated by the 3D-Var assimilation of a single-temperature observation (positive innovation) located at a depth of 100 m on the equator in the central Pacific. The contour interval is 2.0 K in (a), 0.2 psu in (b), 0.1 m s<sup>-1</sup> in (c), 0.01 m s<sup>-1</sup> in (d), and 0.02 m in (e). The fields have been multiplied by a factor 100. Solid (dashed) contours indicate positive (negative) values.

information on SSH as well as indirect information on T and S via the dynamic height relation (19). In this case, the covariances for the unbalanced components of S and SSH, as well as those for T, would influence the analysis, and the adjoint of the balance operator would be required in the minimization process to map gradient information from SSH into gradient information for the other fields (the appendix).

Figure 3 shows a zonal-vertical section at the equator of the T increment (Fig. 3(a)) and S increment (Fig. 3(b)) generated by the 3D-Var assimilation of a single SSH observation, chosen to be 5 cm higher than the background SSH, on the equator in the eastern Pacific (110°W). The observation-error variance has been set to  $(0.5 \text{ cm})^2$ . To fit the SSH observation, 3D-Var produces T and S increments with largest amplitude at the level of the thermocline. The vertical structures are noticeably anisotropic. The increments display a pronounced upward tilt from west to east commensurate with the tilt of the background isotherms in this region. This anisotropic response is produced by the gradient-dependent T variances. The S increment has a dipole-like structure where the transition from negative to positive values occurs at the level of the salinity maximum in the background state. Above this level, the vertical derivative of the background salinity is negative (salinity increases with depth), whereas below this level, the vertical derivative is positive (salinity decreases with depth). Since the vertical derivative of the background temperature is everywhere negative (temperature decreases with depth), there is a change in sign in the derivative  $\partial S/\partial T|_{T=T^b, S=S^b}$  in (17) which gives rise to the dipole in Fig. 3(b).



Figure 3. Vertical cross-section at the equator of the analysis increments for (a) temperature and (b) salinity generated by the 3D-Var assimilation of a single SSH observation (positive innovation) located on the equator in the central Pacific. The contour interval is 2.0 K in (a), and 0.1 psu in (b). The fields have been multiplied by a factor 100. Solid (dashed) contours indicate positive (negative) values.



Figure 4. Horizontal section of the SSH analysis increments generated by the 4D-Var assimilation of a single-temperature observation (positive innovation) located ten days into an assimilation window at the same geographical location as in the example in Fig. 2. The increments are displayed on day 10 for a 4D-Var experiment (a) without and (b) with the balance operator activated. The fields have been multiplied by a factor 100 and the same contour interval has been used here as in Fig. 2(e). Solid (dashed) contours indicate positive (negative) values.

The previous examples illustrate the fundamental importance of the balance operator in establishing a physically sensible (multivariate) response in 3D-Var. The balance operator also plays an important role in 4D-Var. This is illustrated in Fig. 4 which shows the SSH increments produced from two 4D-Var single T observation experiments performed without and with the balance operator activated (Figs. 4(a) and (b), respectively). The geographical location of the single T observation is the same as in the example in Fig. 2. In these experiments, the control variables are a function of the model initial conditions which are taken to be ten days before the observation time. For the experiment without the balance operator, the background-error covariances must be specified for the full fields at initial time. The correlation models for S and velocity are taken to be identical to those used by Weaver et al. (2003) for a rigid-lid version of OPA, while the correlation model for SSH is taken to be identical to the horizontal correlation model for T and S. The variances are set to values typical of the climatological variability of these fields:  $(0.08 \text{ m})^2$  for SSH, and surface values of  $(0.25 \text{ psu})^2$  for S,  $(0.4 \text{ m s}^{-1})^2$ for u, and  $(0.1 \text{ m s}^{-1})^2$  for v. The variances for S, u and v are gradually reduced below the surface. For the experiment with the balance operator, the unbalanced variances are set to zero as in the previous example, while the variances of the balanced components are defined implicitly via interactions between the balance operator and univariate Tcovariances.

The increments shown in Figs. 4(a) and (b) are those produced at the observation time (day 10) and have been computed by using the tangent-linear model to propagate forward the analysis increment at initial time. The SSH increment in the first 4D-Var

experiment has localized structure similar to that obtained by 3D-Var with the balance operator (cf. Figs. 4(a) and 2(e)). In terms of the analysis of SSH, nothing much appears to have been gained by using 4D-Var. When the balance operator is included, however, the temperature observation projects much more effectively onto large-scale equatorial wave-modes as clearly illustrated in Fig. 4(b) by the presence of a westward-propagating baroclinic Rossby wave to the west of the observation location. Contrary to the 4D-Var experiment without the balance operator, the observation is able to have a much wider impact than in 3D-Var.

## 5. SUMMARY AND CONCLUSIONS

The background-error covariances are often cited as a critical component of a statistical data assimilation system. An arguably more fundamental component is the operator that is needed to compute the model counterpart of the assimilated observations. In this paper, it was shown how linear balance and smoothness constraints, that are traditionally used to model multivariate covariances of background error, could be cast within the more general, nonlinear, framework of an observation operator. The key aspect of this procedure is the design of a transformation, possibly nonlinear, from the space of highly correlated model state variables to a space of non-dimensional control variables that are approximately mutually uncorrelated. In the space of the transformed variables, the background-error covariance matrix is assumed to be the identity matrix. The inverse of the transformation, or its generalized inverse if the dimension of the control space is smaller than that of model space, is also needed so not all transformations are suitable.

This paper outlined a control variable transformation for application to variational ocean data assimilation. The focus was on the balance operator, the inverse of which is designed to decorrelate the model state variables of temperature, salinity, SSH and velocity. In the proposed formulation, the inverse of the sequence of balance relationships left temperature unaltered but removed parts from salinity that could be related to temperature, parts from SSH that could be related to temperature and salinity, and parts from velocity that could be related to temperature, salinity and SSH. Both linear constraints (geostrophy, hydrostatic, dynamic height) and nonlinear constraints (T-S)relationship, equation of state) were employed. In incremental variational assimilation, nonlinear constraints are linearized about a reference state as part of the minimization process. Furthermore, by linearizing the control variable transformation within the definition of the reference state itself, the minimization problem can be solved without the need to perform either the nonlinear transformation or its inverse. This is a convenient approximation but may break down when the increments are large. Further research is needed to quantify the impact of this approximation for the nonlinear balance operator proposed here.

Evidence that the proposed ocean balance operator can explain a substantial amount of actual background-error variance was provided by considering the statistical balance properties of a large set of differences between model forecasts verifying at the same time. Single-observation experiments were performed to illustrate the multivariate analysis structures implied by the balance operator. One example illustrated how the balance operator could be used as an effective way to project SSH (altimeter) data onto the subsurface density field in 3D-Var. Another example illustrated the potential benefits of the balance operator for equatorial analysis with 4D-Var. To obtain full benefit from the balance formulation in realistic implementations will require careful specification of the error covariance statistics of the transformed (uncorrelated) state variables. Ensemble methods could be very promising for this purpose.

## **ACKNOWLEDGEMENTS**

Support for this work was provided by the European ENACT project (contract No. EVK2-CT2001-00117) and the Groupe Mission MERCATOR/CORIOLIS. We are grateful to M. Balmaseda and two anonymous referees for their thoughtful comments on the paper.

#### APPENDIX

# Matrix representation of the background-error covariance model

In this appendix, an explicit form of the background-error covariance matrix is derived to illustrate how, on a given outer iteration of the incremental variational algorithm (10), the components of the balance operator combine with the univariate blocks of the covariance matrix of the uncorrelated variables to produce a full-rank multivariate covariance matrix for the model variables. For clarity of notation, the superscript k - 1 on linearized operators will be omitted. From (10), the background-error covariance matrix of the model state **x** is related to the background-error covariance matrix of the uncorrelated state variables  $\hat{\mathbf{x}}$  by

$$\mathbf{B}_{(\mathbf{x})} = \mathbf{K} \mathbf{B}_{(\widehat{\mathbf{x}})} \mathbf{K}^{\mathrm{T}},\tag{A.1}$$

where  $\mathbf{B}_{(\widehat{\mathbf{x}})} = \mathbf{D}_{(\widehat{\mathbf{x}})} \mathbf{F} \mathbf{F}^{\mathrm{T}} \mathbf{D}_{(\widehat{\mathbf{x}})}^{\mathrm{T}}$  is a block matrix of the form

$$\mathbf{B}_{(\widehat{\mathbf{x}})} = \begin{pmatrix} \mathbf{B}_{TT} & 0 & 0 & 0 & 0 \\ 0 & \mathbf{B}_{S_{\mathrm{U}}S_{\mathrm{U}}} & 0 & 0 & 0 \\ 0 & 0 & \mathbf{B}_{\eta_{\mathrm{U}}\eta_{\mathrm{U}}} & 0 & 0 \\ 0 & 0 & 0 & \mathbf{B}_{u_{\mathrm{U}}u_{\mathrm{U}}} & \mathbf{B}_{v_{\mathrm{U}}u_{\mathrm{U}}}^{\mathrm{T}} \\ 0 & 0 & 0 & \mathbf{B}_{v_{\mathrm{U}}u_{\mathrm{U}}} & \mathbf{B}_{v_{\mathrm{U}}u_{\mathrm{U}}}^{\mathrm{T}} \end{pmatrix},$$
(A.2)

with  $\mathbf{B}_{\hat{x}\hat{x}} = \mathbf{D}_{\hat{x}}\mathbf{F}_{\hat{x}\hat{x}}\mathbf{D}_{\hat{x}}^{\mathrm{T}}, \hat{x} = T, S_{\mathrm{U}}, \eta_{\mathrm{U}} \text{ and } \mathbf{u}_{\mathrm{U}}^{\mathrm{h}}$ . The four-block submatrix in the lower corner of (A.2) corresponds to  $\mathbf{B}_{\mathbf{u}_{\mathrm{U}}^{\mathrm{h}}\mathbf{u}_{\mathrm{U}}^{\mathrm{h}}}$ . A non-zero cross-covariance between  $u_{\mathrm{U}}$  and  $v_{\mathrm{U}}$  arises since the smoothing operator  $\mathbf{F}_{\mathbf{u}_{\mathrm{U}}^{\mathrm{h}}\mathbf{u}_{\mathrm{U}}^{\mathrm{h}}}$  employed involves a vector Laplacian operator which smooths separately horizontal divergence and relative vorticity (Weaver *et al.* 2003).

The balance operator is a lower diagonal matrix of the form

$$\mathbf{K} = \begin{pmatrix} \mathbf{I} & 0 & 0 & 0 & 0 \\ \mathbf{K}_{ST} & \mathbf{I} & 0 & 0 & 0 \\ \mathbf{K}_{\eta T} & \mathbf{K}_{\eta S} & \mathbf{I} & 0 & 0 \\ \mathbf{K}_{uT} & \mathbf{K}_{uS} & \mathbf{K}_{u\eta} & \mathbf{I} & 0 \\ \mathbf{K}_{vT} & \mathbf{K}_{vS} & \mathbf{K}_{v\eta} & 0 & \mathbf{I} \end{pmatrix},$$
(A.3)

where, from (15) and (16),

 $\mathbf{K}_{\eta T} = \mathbf{K}_{\eta \rho} \ \mathbf{K}_{\rho T},$   $\mathbf{K}_{uT} = \mathbf{K}_{up} \ \mathbf{K}_{p\rho} \ \mathbf{K}_{\rho T},$   $\mathbf{K}_{vT} = \mathbf{K}_{vp} \ \mathbf{K}_{p\rho} \ \mathbf{K}_{\rho T},$   $\mathbf{K}_{\eta S} = \mathbf{K}_{\eta \rho} \ \mathbf{K}_{\rho S},$  $\mathbf{K}_{uS} = \mathbf{K}_{up} \ \mathbf{K}_{p\rho} \ \mathbf{K}_{\rho S},$ 

$$\mathbf{K}_{vS} = \mathbf{K}_{vp} \ \mathbf{K}_{p\rho} \ \mathbf{K}_{\rho S},$$
$$\mathbf{K}_{u\eta} = \mathbf{K}_{up} \ \mathbf{K}_{p\eta},$$
$$\mathbf{K}_{v\eta} = \mathbf{K}_{vp} \ \mathbf{K}_{p\eta}.$$

Substituting (A.2) and (A.3) in (A.1) and carrying out the matrix multiplication gives

$$\mathbf{B}_{(\mathbf{x})} = \begin{pmatrix} \mathbf{B}_{TT} & \mathbf{B}_{ST}^{\mathrm{T}} & \mathbf{B}_{\eta T}^{\mathrm{T}} & \mathbf{B}_{u T}^{\mathrm{T}} & \mathbf{B}_{v T}^{\mathrm{T}} \\ \mathbf{B}_{ST} & \mathbf{B}_{SS} & \mathbf{B}_{\eta S}^{\mathrm{T}} & \mathbf{B}_{u S}^{\mathrm{T}} & \mathbf{B}_{v S}^{\mathrm{T}} \\ \mathbf{B}_{\eta T} & \mathbf{B}_{\eta S} & \mathbf{B}_{\eta \eta} & \mathbf{B}_{u \eta}^{\mathrm{T}} & \mathbf{B}_{v \eta}^{\mathrm{T}} \\ \mathbf{B}_{u T} & \mathbf{B}_{u S} & \mathbf{B}_{u \eta} & \mathbf{B}_{u u} & \mathbf{B}_{v u}^{\mathrm{T}} \\ \mathbf{B}_{v T} & \mathbf{B}_{v S} & \mathbf{B}_{v \eta} & \mathbf{B}_{v u} & \mathbf{B}_{v v} \end{pmatrix},$$
(A.4)

where

$$\begin{aligned} \mathbf{B}_{ST} &= \mathbf{K}_{ST} \mathbf{B}_{TT}, \\ \mathbf{B}_{\eta T} &= \mathbf{K}_{\eta T} \mathbf{B}_{TT}, \\ \mathbf{B}_{uT} &= \mathbf{K}_{uT} \mathbf{B}_{TT}, \\ \mathbf{B}_{uT} &= \mathbf{K}_{uT} \mathbf{B}_{TT}, \\ \mathbf{B}_{vT} &= \mathbf{K}_{vT} \mathbf{B}_{TT}, \\ \mathbf{B}_{SS} &= \mathbf{K}_{ST} \mathbf{B}_{TT} \mathbf{K}_{ST}^{T} + \mathbf{B}_{S_U S_U}, \\ \mathbf{B}_{\eta S} &= \mathbf{K}_{\eta T} \mathbf{B}_{TT} \mathbf{K}_{ST}^{T} + \mathbf{K}_{\eta S} \mathbf{B}_{SU S_U}, \\ \mathbf{B}_{uS} &= \mathbf{K}_{uT} \mathbf{B}_{TT} \mathbf{K}_{ST}^{T} + \mathbf{K}_{uS} \mathbf{B}_{SU S_U}, \\ \mathbf{B}_{vS} &= \mathbf{K}_{vT} \mathbf{B}_{TT} \mathbf{K}_{ST}^{T} + \mathbf{K}_{vS} \mathbf{B}_{SU S_U}, \\ \mathbf{B}_{\eta \eta} &= \mathbf{K}_{\eta T} \mathbf{B}_{TT} \mathbf{K}_{\eta T}^{T} + \mathbf{K}_{\eta S} \mathbf{B}_{SU S_U} \mathbf{K}_{\eta S}^{T} + \mathbf{B}_{\eta U \eta U}, \\ \mathbf{B}_{u\eta} &= \mathbf{K}_{uT} \mathbf{B}_{TT} \mathbf{K}_{\eta T}^{T} + \mathbf{K}_{uS} \mathbf{B}_{SU S_U} \mathbf{K}_{\eta S}^{T} + \mathbf{K}_{u\eta} \mathbf{B}_{\eta U \eta U}, \\ \mathbf{B}_{v\eta} &= \mathbf{K}_{vT} \mathbf{B}_{TT} \mathbf{K}_{\eta T}^{T} + \mathbf{K}_{vS} \mathbf{B}_{SU S_U} \mathbf{K}_{\eta S}^{T} + \mathbf{K}_{v\eta} \mathbf{B}_{\eta U \eta U}, \\ \mathbf{B}_{uu} &= \mathbf{K}_{uT} \mathbf{B}_{TT} \mathbf{K}_{uT}^{T} + \mathbf{K}_{uS} \mathbf{B}_{SU S_U} \mathbf{K}_{uS}^{T} + \mathbf{K}_{u\eta} \mathbf{B}_{\eta U \eta U}, \\ \mathbf{B}_{uu} &= \mathbf{K}_{uT} \mathbf{B}_{TT} \mathbf{K}_{uT}^{T} + \mathbf{K}_{vS} \mathbf{B}_{SU S_U} \mathbf{K}_{uS}^{T} + \mathbf{K}_{u\eta} \mathbf{B}_{\eta U \eta U} \mathbf{K}_{u\eta}^{T} + \mathbf{B}_{u U u}, \\ \mathbf{B}_{vv} &= \mathbf{K}_{vT} \mathbf{B}_{TT} \mathbf{K}_{uT}^{T} + \mathbf{K}_{vS} \mathbf{B}_{SU S_U} \mathbf{K}_{uS}^{T} + \mathbf{K}_{v\eta} \mathbf{B}_{\eta U \eta U} \mathbf{K}_{u\eta}^{T} + \mathbf{B}_{v U u}, \\ \mathbf{B}_{vv} &= \mathbf{K}_{vT} \mathbf{B}_{TT} \mathbf{K}_{vT}^{T} + \mathbf{K}_{vS} \mathbf{B}_{SU S_U} \mathbf{K}_{uS}^{T} + \mathbf{K}_{v\eta} \mathbf{B}_{\eta U \eta U} \mathbf{K}_{u\eta}^{T} + \mathbf{B}_{v U u}. \end{aligned}$$

To interpret the results of the single-observation experiments in section 4, it is helpful to illustrate how the algebraic structure of (A.4) determines the expression for the increment  $\delta \mathbf{x}^k$  on each outer iteration. Consider the exact minimizing solution of (4), which is found by setting the gradient of (4) to zero and solving for  $\delta \mathbf{v}^k$  (e.g. see Daley 1991):

$$\delta \mathbf{v}^k = \mathbf{G}^{\mathrm{T}} (\mathbf{G} \mathbf{G}^{\mathrm{T}} + \mathbf{R})^{-1} \mathbf{d}^{\mathrm{o},k}.$$
(A.5)

For a single observation,  $\mathbf{d}^{o,k} = d^{o,k}$ ,  $\mathbf{R} = (\sigma^o)^2$  and  $\mathbf{G}\mathbf{G}^T = (\sigma^k)^2$  are scalars, where the latter two quantities correspond to, respectively, the observation-error variance and the effective background-error variance for the observation on the *k*th outer iteration. If the observation is situated at the end of the assimilation window  $(t = t_n)$  then  $\mathbf{G} = \mathbf{H}_n \mathbf{M}(t_n, t_0) \mathbf{U}$  where  $\mathbf{H}_n = \mathbf{h}^T$  is the observation operator, which for a single observation is a vector of the same length as  $\delta \mathbf{x}$ , and  $\mathbf{M}(t_n, t_0)$  is the linearized forward propagator which is the identity matrix in 3D-Var and the tangent-linear operator in 4D-Var. Substituting these expressions into (A.5) and transforming the increment into model space gives

$$\delta \mathbf{x}^{k} = \mathbf{U} \delta \mathbf{v}^{k} = c \mathbf{B}_{(\mathbf{x})} \mathbf{M}(t_{n}, t_{0})^{\mathrm{T}} \mathbf{h}, \qquad (A.6)$$

where  $c = d^{0,k} \{(\sigma^k)^2 + (\sigma^0)^2\}^{-1}$  and, from (10),  $\mathbf{B}_{(\mathbf{x})} = \mathbf{U}\mathbf{U}^T$ . From (A.6) it is clear that  $\delta \mathbf{x}^k$  will be proportional to the columns of the matrix  $\mathbf{B}_{(\mathbf{x})} \mathbf{M}(t_n, t_0)^T$ , or simply the columns of  $\mathbf{B}_{(\mathbf{x})}$  in the case of 3D-Var. For example, for a temperature observation,  $\mathbf{h} = (\mathbf{e}^T, 0, 0, 0, 0)^T$ , where  $\mathbf{e}$  is a vector corresponding to the temperature components of  $\delta \mathbf{x}^k$ , and the other elements of  $\mathbf{h}$  are zero vectors corresponding to the other variable components. (If the temperature observation is located exactly at a model grid point then  $\mathbf{e} = (0, \dots, 0, 1, 0, \dots, 0)^T$  where the non-zero entry is associated with that gridpoint.) In this case, it is easy to see that the 3D-Var increment will be proportional to the first block-column of  $\mathbf{B}_{(\mathbf{x})}$ , in particular dependent on  $\mathbf{B}_{TT}$  and the forward balance operators only. Likewise, the 3D-Var increment will be proportional to the second blockcolumn of  $\mathbf{B}_{(\mathbf{x})}$  for a salinity observation, proportional to the third block-column of  $\mathbf{B}_{(\mathbf{x})}$  for a vector velocity observation. Notice that in 4D-Var, regardless of what type of observation is assimilated,  $\delta \mathbf{x}^k$  will be a non-trivial linear combination of all blockcolumns of  $\mathbf{B}_{(\mathbf{x})}$  since the action of the adjoint operator  $\mathbf{M}(t_n, t_0)^T$  will result in a transfer of information from the observed quantity to all model variables.

### REFERENCES

Alves, O., Balmaseda, M. A., 2004 Sensitivity to dynamical seasonal forecasts to ocean initial Anderson, D. L. T. and conditions. Q. J. R. Meteorol. Soc., 130, 647-668 Stockdale, T. 'Ocean data assimilation for seasonal forecasts'. Pp. 301-325 Balmaseda, M. A. 2004 in Proceedings of the ECMWF Seminar on recent developments in data assimilation for atmosphere and ocean, 8–12 September 2003, Reading, UK Behringer, D., Ji, M. and 1998 An improved coupled model for ENSO prediction and implica-Leetma, A. tions for ocean initialization. Part I: The ocean data assimilation system. Mon. Weather Rev., 126, 1013–1021 Bennett, A. F., Chua, B. S., 2000 Generalized inversion of Tropical Atmosphere-Ocean (TAO) data Harrison, D. E. and using a coupled model of the tropical Pacific. J. Climate, 13, McPhaden, M. J. 2770-2785 Berre, L., Ştefănescu, S. E. and 2006 The representation of the analysis effect in three error simulation Pereira, M. B. techniques. Tellus, 58, 196-209 Ensemble-derived stationary and flow-dependent background-Buehner, M. 2005 error covariances: Evaluation in a quasi-operational NWP setting. Q. J. R. Meteorol. Soc., 131, 1013-1044 Burgers, G., Balmaseda, M. A., 2002 Balanced ocean-data assimilation near the equator. J. Phys. Vossepoel, F. C., Oceanogr., 32, 2509-2519 van Oldenborgh, G. J. and van Leeuwen, P. J. Courtier, P., Thépaut, J.-N. and 1994 A strategy for operational implementation of 4D-Var, using an incremental approach. Q. J. R. Meteorol. Soc., 120, Hollingsworth, A. 1367-1388 Cullen, M. J. P. 2003 Four-dimensional variational data assimilation: A new formulation of the background-error covariance matrix based on a potential-vorticity representation. Q. J. R. Meteorol. Soc., **129,** 2777–2796 1991 Atmospheric data analysis. Cambridge atmospheric and space Daley, R. sciences series, Cambridge University Press Derber, J. and Bouttier, F. 1999 A reformulation of the background error covariance in the ECMWF global data assimilation system. Tellus, 51A, 195-221 Derber, J. and Wu, W.-S. 1998 The use of TOVS cloud-cleared radiances in the NCEP SSI analysis system. Mon. Weather Rev., 126, 2287-2299

1998

Ferry, N.

- Fukumori, I., Raghunath, R. and Fu, L.
- Golub, G. and Van Loan, C. F.
- Han, G., Zhu, J. and Zhou, G.
- Ingleby, B. and Huddleston, M.
- Keppenne, C. L. and Rienecker, M. M.
- Lagerloef, G. S. E., Mitchum, G. T., 1999 Lukas, R. B. and Niiler, P. P. Lermusiaux, P. F. J., 2000 Anderson, D. G. M. and Lozano, C. J.
- Lorenc, A. C.
- McDougall, T. J., Jackett, D. R., Wright, D. G. and Feistel, R.
- Madec, G., Delecluse, P., Imbard, M. and Levy, C. Parrish, D. F. and Derber, J. C.
- Picaut, J. and Tournier, R.
- Ricci, S., Weaver, A. T., Vialard, J. 2005 and Rogel, P.
- Roullet, G. and Madec, G.
- Stammer, D., Wunsch, C., Giering, R., Eckert, C., Heimbach, P., Marotzke, J., Adcroft, A., Hill, C. N. and Marshall, J.
- Testut, C.-E., Brasseur, P., Brankart, J. M. and Verron, J.
- Troccoli, A. and Haines, K.
- Troccoli, A., Balmaseda, M. A., Segschneider, J., Vialard, J., Anderson, D. L. T., Stockdale, T., Haines, K. and Fox, A. D.

- 2003 Assimilation et surface libre dans les modèles océaniques MERCATOR. Rapport interne Projet MERCATOR, Référence MOO-ST-410-231-MER
  - Nature of global large-scale sea level variability in relation to atmospheric forcing: a modeling study. *J. Geophys. Res.*, **103**, 5493–5512
- 1996 Matrix computations. The Johns Hopkins University Press, London
- 2004 Salinity estimation using the *T*-*S* relation in the context of variational data assimilation. *J. Geophys. Res.*, **109**, C03018, doi: 10.1029/2003JC001781
- 2006 Quality control of ocean temperature and salinity profiles historical and real-time data. J. Mar. Sys. (in press)
- 2003 Assimilation of temperature into an isopycnal ocean general circulation model using a parallel ensemble Kalman filter. *J. Mar. Sys.*, **40–41**, 363–380
  - Tropical Pacific near-surface currents estimated from altimeter, wind, and drifter data. J. Geophys. Res., **104**, 23313–23326
  - On the mapping of multivariate geophysical fields: Error and variability subspace estimates. Q. J. R. Meteorol. Soc., 126, 1387–1430
- 2003 Modelling of error covariances by 4D-Var data assimilation. *Q. J. R. Meteorol. Soc.*, **129**, 3167–3182
- Accurate and computationally efficient algorithms for potential temperature and density of seawater. J. Atmos. Ocean. Technol., 20, 730–741
   OPA 8.1 Ocean General Circulation Model reference manual.
- 1998 OPA 8.1 Ocean General Circulation Model reference manual. Technical note no. 11, LODYC/IPSL, Paris, France
   1992 The National Meteorological Center's spectral statistical
  - The National Meteorological Center's spectral statistical interpolation analysis system. *Mon. Weather Rev.*, **120**, 1747–1763
- 1991 Monitoring the 1979–1985 equatorial Pacific current transports with expendable bathythermograph data. J. Geophys. Res., 96, 3263–3277
  - Incorporating temperature–salinity constraints in the background error covariance of variational ocean data assimilation. *Mon. Weather Rev.*, **133**, 317–338
- 2000 Salt conservation, free surface, and varying levels: A new formulation for ocean general circulation models. *J. Geophys. Res.*, **105**, 23927–23942
- 2002 Global ocean circulation during 1992–1997, estimate from ocean observations and a general circulation model. *J. Geophys. Res.*, **107**, C93118, doi: 10.1029/2001JC000888
- 2003 Assimilation of sea-surface temperature and altimetric observations during 1992–1993 into an eddy permitting primitive equation model of the North Atlantic Ocean. J. Mar. Sys., 40-41, 291–316
- 1999 Use of temperature–salinity relation in a data assimilation context. *J. Atmos. Ocean. Technol.*, **16**, 2011–2025
- 2002 Salinity adjustments in the presence of temperature data assimilation. *Mon. Weather Rev.*, **130**, 89–102

Uppala, S. M., Kållberg, P. W., Simmons, A. J., Andrae, U., Bechtold, da Costa V., Fiorino, M., Gibson, J. K., Haseler, J., Hernandez, A., Kelly, G. A., Li, X., Onogi, K., Saarinen, S., Sokka, N., Allan, R. P., Andersson, E., Arpe, K., Balmaseda, M. A., Beljaars, A. C. M., van de Berg, L., Bidlot, J., Bormann, N., Caires, S., Chevallier, F., Dethof, A., Dragosavac, M., Fisher, M., Fuentes, M., Hagemann, S., Hólm, E., Hoskins, B. J., Isaksen, L., Janssen, P. A. E. M., Jenne, R., McNally, A. P., Mahfouf, J.-E., Morcrette J.-J., Rayner, N. A., Saunders, R. W., Simon, P., Sterl, A., Trenberth, K. E., Untch, A., Vasiljevic, D., Viterbo, P. and Woollen, J. Weaver, A. T. and Courtier, P.

Weaver, A. T., Vialard, J. and Anderson, D. L. T. 2003

- 2001 Correlation modelling on the sphere using a generalized diffusion equation. *Q. J. R. Meteorol. Soc.*, **127**, 1815–1846
  - Three- and four-dimensional variational assimilation with an ocean general circulation model of the tropical Pacific Ocean. Part 1: formulation, internal diagnostics and consistency checks. *Mon. Weather Rev.*, **131**, 1360–1378

2005 The ERA-40 re-analysis. Q. J. R. Meteorol. Soc., **131**, 2961–3012



# Ensemble estimation of background-error variances in a three-dimensional variational data assimilation system for the global ocean

N. Daget<sup>a</sup>, A. T. Weaver<sup>a</sup>\* and M. A. Balmaseda<sup>b</sup> <sup>a</sup> CERFACS / SUC URA 1875, Toulouse, France <sup>b</sup> ECMWF, Reading, UK

**ABSTRACT:** This paper studies the sensitivity of global ocean analyses to two flow-dependent formulations of the background-error standard deviations ( $\sigma^b$ ) for temperature and salinity in a three-dimensional variational data assimilation (3D-Var) system. The first formulation is based on an empirical parameterization of  $\sigma^b$  in terms of the vertical gradients of the background temperature and salinity fields, while the second formulation involves a more sophisticated approach that derives  $\sigma^b$  from the spread of an ensemble of background states. The ensembles are created by explicitly perturbing both the surface fluxes (wind stress, fresh water and heat) used to force the model and the observations (temperature and salinity profiles) used in the assimilation process. The two formulations are compared in two cycled 3D-Var experiments for the period 1993–2000. In both experiments, the observation-error standard deviations ( $\sigma^o$ ) are geographically dependent and estimated from a model-data comparison prior to assimilation. An additional 3D-Var experiment that employs the parametrized  $\sigma^b$  but a simpler  $\sigma^o$  formulation, and a control experiment involving no data assimilation, were also conducted and used for comparison.

All 3D-Var experiments produce a significant reduction in the mean and standard deviation of the temperature and salinity innovations compared to those of the control experiment. The largest differences between the two  $\sigma^{b}$  formulations occur in the upper 150 m, where the parametrized  $\sigma^{b}$  are notably larger than the ensemble  $\sigma^{b}$ . In this region, the innovation statistics are slightly better for the parametrized  $\sigma^{b}$ . Statistical consistency checks indicate that both schemes underestimate  $\sigma^{b}$ , the underestimation being stronger with the ensemble formulation. The error growth between cycles, however, is much reduced with the ensemble  $\sigma^{b}$ , suggesting that the analyses produced with the ensemble  $\sigma^{b}$  are in better balance than those produced with the parametrized  $\sigma^{b}$ . This claim is supported by independent data comparisons involving model fields not directly constrained by the assimilated temperature and salinity profiles. In particular, sea-surface height (SSH) anomalies in the northwest Atlantic and zonal velocities in the equatorial Pacific are clearly better with the ensemble  $\sigma^{b}$  than with the parametrized  $\sigma^{b}$ . Results also show that while some aspects of those variables are improved with data assimilation (SSH anomalies and currents in the central and eastern Pacific), other aspects are degraded (SSH anomalies in the northwest Atlantic, currents in the western Pacific). Areas for improving the ensemble method and for making better use of the ensemble information are discussed. Copyright © 2009 Royal Meteorological Society

KEY WORDS oceanography; ocean reanalysis; 3D-Var; 4D-Var; covariance estimation

Received 6 June 2008; Revised 3 February 2009; Accepted 25 February 2009

#### 1. Introduction

An important feature of an ensemble data assimilation system is its capacity to provide flow-dependent information on analysis and background error. This information can be exploited in a cycled assimilation system to improve the estimate of the background-error covariance matrix on each cycle. The simplest way to use the ensemble information is to build a lowrank approximation to the background-error covariance matrix on a given cycle from the sample covariance of the ensemble of model forecast states initiated from the previous cycle. The matrix is rank deficient since

Copyright © 2009 Royal Meteorological Society

the number of ensemble members is typically several orders of magnitude smaller than the number of background state variables. In the Ensemble Kalman Filter (EnKF), this rank deficiency can be exploited to produce computationally efficient implementations of the standard Kalman filter analysis equation. (Houtekamer and Mitchell (2005) and Evensen (2007) provide reviews of the different variants of the EnKF.) However, using a small ensemble to estimate the covariance matrix directly in a high-dimensional system can lead to noisy variances and spurious long-range correlations due to sampling error. Various filtering and localization procedures have been proposed to alleviate this problem in practical implementations of the EnKF (Houtekammer and Mitchell, 2001; Keppenne and Reinecker, 2002; Ott et al., 2004; Buehner and Charron, 2007; Oke et al., 2007).

<sup>\*</sup>Correspondence to: A. T. Weaver, CERFACS, 42 avenue Gaspard Coriolis, 31057 Toulouse Cedex 01, France. E-mail: Anthony.Weaver@cerfacs.fr

Lorenc (2003b) and Buehner (2005) illustrate how an ensemble-estimated background-error covariance matrix, with or without localization, can be used in a variational assimilation scheme. The procedure involves using the square root of the (localized) ensemble covariance matrix to transform the control vector into a vector of background-state increments. The basic transformation is designed to precondition the minimization problem and is standard in variational assimilation systems that employ more conventional background-error covariance formulations based on covariance models (Derber and Bouttier, 1999; Lorenc, 2003a; Weaver et al., 2005). Methods to define the background-error covariance matrix as a linear combination of an ensemble-estimated matrix and a covariance model matrix have also been proposed (Hamill and Synder, 2000; Lorenc, 2003b; Buehner, 2005).

Rather than using the ensemble directly to construct an estimate of the covariance matrix, it may be used indirectly to calibrate specific parameters of a covariance model (Fisher, 2003; Žagar et al., 2005; Belo Pereira and Berre, 2006; Berre et al., 2006; Küçükkaraca and Fisher, 2006). The use of a covariance model has the advantage of providing a full-rank (implicit) representation of the covariance matrix and thus allows the assimilation method to produce corrections to the background state in a much larger space than that spanned by a limited number of ensemble members. There is also no need for a separate localization procedure since covariance models are constructed to permit only spatially limited covariance functions. The use of ensembles in combination with a variational assimilation scheme is relatively unexplored in ocean data assimilation. The main purpose of this study is to investigate the potential of an ensemble of ocean analyses to provide useful flow-dependent estimates of the background-error variances in a three-dimensional variational assimilation (3D-Var) system. This study can be viewed as a first step towards making more comprehensive use of an ensemble for calibrating additional parameters of the covariance model.

The paper is organized as follows. Section 2 gives a description of the ocean data assimilation system. The sensitivity experiments presented in this paper involve different formulations of both the observationerror variances and background-error variances. These formulations, including the background-error variance formulation based on the ensemble method, are described in section 3. Results from cycled 3D-Var experiments that compare the relative impact of the different variance formulations are presented in section 4. A summary and conclusions are given in section 5. Appendix A provides a derivation of the formula used to estimate geographically dependent observation-error variances. Appendix B presents the mathematical basis of the ensemble method used for estimating background-error covariances.

#### 2. The assimilation system

A variational data assimilation system has been developed at CERFACS for climate research applications. The system, known as OPAVAR, is based on an incremental variational algorithm (Courtier et al., 1994) and version 8.2 of the Océan Parallélisé (OPA) ocean general circulation model (Madec et al., 1998). Three- and fourdimensional variational assimilation (3D-Var and 4D-Var) versions of the system were initially developed for tropical Pacific basin applications (Weaver et al., 2003; Vialard et al., 2003; Vossepoel et al., 2004; Ricci et al., 2005). The system was later extended to a global configuration in the European project ENACT (Enhanced Ocean Data Assimilation and Climate Prediction; http://www. ecmwf.int/research/EU\_projects/ENACT), where it was applied to produce multi-decadal ocean analyses for seasonal hindcast initialization and studies of ocean climate variability (Davey et al., 2006; Carton and Santorelli, 2008). Important advances were made to the system during ENACT, one of the most noteworthy being the development of a fully multivariate background-error covariance model based on balance operators (Weaver et al., 2005). More recently the system has been extended in the European project ENSEMBLES (Ensemble-based Predictions of Climate Changes and their Impacts; http:// www.ecmwf.int/research/EU\_projects/ENSEMBLES) to generate a nine-member ensemble of multi-decadal ocean analyses. The ensemble was produced using multiple atmospheric forcing fields whose differences were constructed to be consistent with estimates of the actual uncertainty in these fields. In ENSEMBLES, the ocean analysis ensemble has been used to contribute to the production of probabilistic forecasts on seasonal to decadal time-scales (Weisheimer et al., 2007). The ensemble 3D-Var system used in this study is based on the system developed for ENSEMBLES. The basic components of the system are described in the remainder of this section.

#### 2.1. Ocean model and forcing fields

The ocean model is a global, free-surface configuration of the ocean general circulation model OPA8.2 (Madec et al., 1998). The model solves the primitive equations for horizontal currents,  $\mathbf{u}_{\rm h} = (u, v)$ , potential temperature, T, salinity, S, and sea-surface height (SSH),  $\eta$ . The freesurface formulation is described in Roullet and Madec (2000). The equations are formulated in orthogonal curvilinear z-coordinates and discretized using finite differences on an Arakawa C-grid. The horizontal grid is stretched in the Northern Hemisphere and contains two poles located on the North American and Asian continents. Outside the equatorial region, the grid mesh is approximately isotropic (Mercator-like) with zonal  $\times$ meridional resolution approximately  $2^{\circ} \times 2^{\circ} \cos \phi$ , where  $\phi$  is latitude. Within the equatorial region, the meridional resolution is increased, with the grid size reaching a value of  $0.5^{\circ}$  at the Equator. Increased resolution is also used in the Mediterranean Sea  $(1^{\circ} \times 1^{\circ})$  and Red Sea  $(\approx 1^{\circ} \times 2^{\circ})$ . The number of horizontal grid points is  $182 \times 149$ . The model has 31 levels of which 21 are in the upper 1000 m. The thickness of the levels varies from 10 m within the upper 100 m to 500 m below the 3000 m level. The maximum depth is 5500 m.

Lateral and vertical subgrid-scale mixing is parametrized using Laplacian diffusion. Vertical diffusion coefficients for momentum, heat, and salt are computed using a Turbulent Kinetic Energy mixing scheme. Lateral mixing coefficients of momentum, heat and salt are geographically dependent. For heat and salt, the lateral diffusion acts along neutral surfaces and includes an additional tracer advection term following Gent and McWilliams (1990). The model is forced using windstress,  $\tau = (\tau^x, \tau^y)$ , heat flux, Q, and fresh-water (precipitation minus evaporation) flux, PmE, from ERA-40 (Uppala *et al.*, 2005). The fresh-water flux from ERA-40 is known to be inaccurate. Here the model is forced using bias-corrected ERA-40 precipitation from Troccoli and Kållberg (2004).

The ensemble experiments are performed over the 9-year period 1 January 1993 to 31 December 2001. The experiments are designed to test the impact of using the ensemble to update the background-error variances on each assimilation cycle. A separate set of ensemble experiments covering the 46-year period 1960-2005 has also been conducted as part of the ENSEMBLES project. The assimilation system used in those experiments is a close variant of the system used here, the main difference being that there was no attempt to use the ensemble to update the background-error covariance matrix as done in this study. The ENSEMBLES experiments also used a more recent version (EN3) of the quality-controlled in situ dataset described in section 2.2 and these data were not perturbed as in this study (section 3.2). In ENSEMBLES, the ocean analysis ensemble was used to provide initial conditions for seasonal and decadal ensemble forecasts. Results from the assimilation experiments conducted in ENSEMBLES will not be discussed in this paper, although results from a separate experiment that employs a system similar to the one used in ENSEMBLES will be used as a reference for evaluating the impact of the ensemble-generated background-error variances.

The experimental design follows closely the common procedures used in ENSEMBLES and in the earlier project ENACT (Davey et al. 2006). The initial conditions on 1 January 1993 were obtained by spinning up the model from rest and temperature and salinity states defined from the Levitus climatology (Levitus et al., 1998). Climatological ERA-40 forcing was used from 1 January 1978 to 31 December 1982, and daily ERA-40 forcing was used from 1 January 1983 to 31 December 1992. The model sea-surface temperature (SST) field is relaxed to model-gridded SST analysis products. During the spin-up from 1 January 1978 to 31 December 1982, the SST climatology from ERA-40 was used, while dailyinterpolated SST analyses from Reynolds OIv2 (Reynolds et al., 2002) were used from 1 January 1983 onwards. As in ENACT and ENSEMBLES, a globally uniform relaxation coefficient of  $-200 \text{ Wm}^{-2}\text{K}^{-1}$  is used, which corresponds to a relaxation time-scale of 12 days for a mixed-layer depth of 50 m. With this choice, the model SST is always close to the 'observed' SST. This is an important requirement for seasonal and decadal forecast

initialization for which the system has been applied in ENSEMBLES.

Subsurface relaxation to climatology has been applied to control model drift but has been chosen to be rather weak so as not to suppress interannual and decadal variability. A weak global subsurface relaxation to modelgridded temperature and salinity monthly climatology, smoothed with a 3-month running mean, is applied with a 3-year time-scale at all vertical levels (Davey et al., 2006). Within 1000 km of coastlines, the relaxation coefficient is reduced smoothly to zero directly at the coast since the smooth density gradients from the Levitus climatology are not dynamically consistent with the steep topographic gradients in the model. Poleward of 60°N/S, where the model is less reliable because of the absence of an active sea-ice model, the relaxation time-scale is reduced smoothly from 3 years to 50 days at 70°N/S and beyond.

The subsurface relaxation provides a weak relaxation to temperature and salinity climatology in the top ocean model level. For temperature, the relaxation is dominated by the much stronger relaxation to SST described above. For sea-surface salinity, no relaxation is applied other than the weak contribution at the surface from the relaxation to climatology. Imbalances in the fresh-water fluxes cause the globally averaged model SSH field to drift ( $\approx 0.7$  m in 15 years). Here the drift has been suppressed by applying a daily correction to the freshwater fluxes based on the sea-level drift that occurs on the previous day. As a result, the global mean SSH field is very close to zero on any given time step.

#### 2.2. Observations

The assimilation dataset consists of in situ temperature and salinity profiles from version EN2\_v1 of the ENACT/ENSEMBLES quality-controlled dataset (Ingleby and Huddleston, 2007). The data are obtained primarily from the World Ocean Database 2001 (WOD01; Conkright et al., 2002). After 1990, they are supplemented with data from the World Marine Environmental Laboratory (Johnson et al., 2002) and the Global Temperature-Salinity Profile Program. The dataset is essentially composed of bathythermographs (MBTs and XBTs), hydrographic profiles (conductivity-temperaturedepth (CTD) and predecessors), moored buoys from the Tropical Atmosphere–Ocean/Triangle Trans-Ocean Buoy Network (TAO/TRITON) and Prediction and Research Moored Array in the Atlantic (PIRATA) arrays, profiling floats and Argo data. Observations determined by the quality control as 'definitely wrong' or 'probably wrong' were not assimilated. Additional screening has been done directly in the assimilation system. Observations have been rejected in closed seas, in some semi-enclosed seas (Mediterranean, Red, Baltic and Japan Seas), below 1000 m and poleward of 65°N/S. The reason for rejecting the data in those regions was based on the inadequacy of the model or assimilation system to use the observational information effectively, rather than on the actual quality of the observations. Vertical thinning of profiles

was performed to restrict the number of individual measurements between two model levels to a maximum of five. The model background state was not used in any of the quality control decisions, so that all observations were assimilated regardless of their difference from their background counterpart. This was done to ensure that exactly the same observations were assimilated in each of the sensitivity experiments.

#### 2.3. Data assimilation method

The data assimilation method is a variant of the multivariate incremental 3D-Var FGAT (First-Guess at Appropriate Time) method described in Weaver *et al.* (2003, 2005) and Ricci *et al.* (2005). A short description is given below to highlight those features of the method that are important in this study.

Let  $\mathbf{w} = (T, S)^{T}$  denote the model vector of temperature and salinity, both *T* and *S* being understood to be row-vectors defined on the three-dimensional (3D) model grid. The superscript *T* denotes transpose. Let  $\mathbf{w}^{b} = (T^{b}, S^{b})^{T}$  be a background estimate of  $\mathbf{w}$ , and let  $\delta \mathbf{w} = (\delta T, \delta S)^{T}$  be an increment defined such that  $\mathbf{w} = \mathbf{w}^{b} + \delta \mathbf{w}$ . Given profile observations of temperature  $(T_{i}^{o})$  and salinity  $(S_{i}^{o})$  distributed over a time window  $t_{0} \leq t_{i} \leq t_{N}$ , 3D-Var FGAT produces an increment  $\delta \mathbf{w}^{a}$ by approximately minimizing the quadratic cost function

$$J[\delta \mathbf{w}] = \frac{1}{2} \delta \mathbf{w}^{\mathrm{T}} \mathbf{B}_{(\mathbf{w})}^{-1} \delta \mathbf{w} + \frac{1}{2} (\mathbf{H} \delta \mathbf{w} - \mathbf{d})^{\mathrm{T}} \mathbf{R}^{-1} (\mathbf{H} \delta \mathbf{w} - \mathbf{d}), \qquad (1)$$

where

$$\mathbf{d} = \begin{pmatrix} \mathbf{d}_0 \\ \vdots \\ \mathbf{d}_i \\ \vdots \\ \mathbf{d}_N \end{pmatrix} = \begin{pmatrix} \mathbf{y}_0^{\mathrm{o}} - \mathbf{H}_0 \mathbf{w}^{\mathrm{b}}(t_0) \\ \vdots \\ \mathbf{y}_i^{\mathrm{o}} - \mathbf{H}_i \mathbf{w}^{\mathrm{b}}(t_i) \\ \vdots \\ \mathbf{y}_N^{\mathrm{o}} - \mathbf{H}_N \mathbf{w}^{\mathrm{b}}(t_N) \end{pmatrix}, \quad (2)$$

 $\mathbf{y}_i^{o} = (T_i^{o}, S_i^{o})^{\mathrm{T}}$  being the observation vector at measurement time  $t_i$ , and  $\mathbf{H}_i \mathbf{w}^{b}(t_i)$  the background counterpart of the observation vector at  $t_i$ . The background state at  $t_i$ ,  $\mathbf{w}^{b}(t_i) = (T_i^{b}, S_i^{b})^{\mathrm{T}}$ , is a subset of the complete model background state vector,  $\mathbf{x}^{b}(t_i) = (T_i^{b}, S_i^{b}, \eta_i^{b}, u_i^{b}, v_i^{b})^{\mathrm{T}}$ , that is obtained by integrating the model from  $t_0$  to  $t_i$  from the background initial condition  $\mathbf{x}^{b}(t_0)$  available at the start of the window. The model integration can be represented as

$$\mathbf{x}^{\mathbf{b}}(t_i) = M(t_i, t_{i-1})[\mathbf{x}^{\mathbf{b}}(t_{i-1}), \mathbf{f}_i], \qquad (3)$$

where  $M(t_i, t_{i-1})$  denotes the nonlinear model operator between  $t_{i-1}$  and  $t_i$ , and  $\mathbf{f}_i = (\tau_i^x, \tau_i^y, Q_i, PmE_i)^T$ denotes the vector of external atmospheric surface fluxes used to force the ocean model on the interval  $t_{i-1}$  to  $t_i$ . These surface fluxes have been made explicit in Equation (3) in order to clarify the description of the ensemble method given in section 3 and Appendix B. The observation operators  $\mathbf{H}_i$  in Equation (2) are 3D interpolation operators at each measurement time  $t_i$  and are formulated as the product of a horizontal  $(\mathbf{H}_i^{\rm h})$  and vertical  $(\mathbf{H}_i^{\rm z})$  interpolation operator. Here,  $\mathbf{H}_i^{\rm z}$  is a cubic spline and  $\mathbf{H}_i^{\rm h}$  is a bilinear interpolation operator, specially adapted to irregular grids (such as the global OPA grid) following the remapping technique of Jones (1998). For TAO/TRITON temperature data, which are daily averaged in the EN2v1 database, the observation operator also includes a daily averaging of the model temperature field. The observation matrix in Equation (1) is given by

$$\mathbf{H} = \begin{pmatrix} \mathbf{H}_0 \\ \vdots \\ \mathbf{H}_i \\ \vdots \\ \mathbf{H}_N \end{pmatrix}. \tag{4}$$

Note that the  $\mathbf{H}_i$  operators in Equation (4) act on the *same* increment  $\delta \mathbf{w}$  in Equation (1), whereas in Equation (2) they act on the different background states  $\mathbf{w}^{b}(t_i)$  in the computation of the innovation vectors  $\mathbf{d}_i = \mathbf{y}_i^{o} - \mathbf{H}_i \mathbf{w}^{b}(t_i)$  in Equation (2).

The matrices  $\mathbf{B}_{(\mathbf{w})}$  and  $\mathbf{R}$  contain estimates of the background- and observation-error covariances, respectively. Observation errors are assumed to be mutually uncorrelated so that  $\mathbf{R} = \mathbf{D}_{(\mathbf{y})} = \mathbf{D}_{(\mathbf{y})}^{1/2} \mathbf{D}_{(\mathbf{y})}^{1/2}$  where  $\mathbf{D}_{(\mathbf{y})}^{1/2} = \text{diag}\{\sigma_T^{\text{o}}, \sigma_S^{\text{o}}\}, \sigma_T^{\text{o}}$  and  $\sigma_S^{\text{o}}$  denoting row-vectors that contain estimates of the standard deviations of temperature and salinity observation error. The specification of the observation-error standard deviations is described in section 3. Background errors are assumed to be correlated. The covariance matrix is described by the symmetric product of operators

$$\mathbf{B}_{(\mathbf{w})} = \mathbf{K}_{(\mathbf{w})} \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} \mathbf{F}_{(\widehat{\mathbf{w}})} \mathbf{F}_{(\widehat{\mathbf{w}})}^{T} \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} \mathbf{K}_{(\mathbf{w})}^{T}$$
(5)

$$= \mathbf{U}_{(\mathbf{w})} \mathbf{U}_{(\mathbf{w})}^{\mathrm{T}}, \qquad (6)$$

where

$$\mathbf{F}_{(\widehat{\mathbf{w}})} = \begin{pmatrix} \mathbf{F}_{TT} & \mathbf{0} \\ \mathbf{0} & \mathbf{F}_{S_U S_U} \end{pmatrix}, \tag{7}$$

$$\mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} = \begin{pmatrix} \mathbf{D}_T^{1/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}_{S_U}^{1/2} \end{pmatrix}, \tag{8}$$

$$\mathbf{K}_{(\mathbf{w})} = \begin{pmatrix} \mathbf{I} & 0\\ \mathbf{K}_{ST} & \mathbf{I} \end{pmatrix},\tag{9}$$

and

$$\mathbf{U}_{(\mathbf{w})} = \mathbf{K}_{(\mathbf{w})} \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} \mathbf{F}_{(\widehat{\mathbf{w}})}.$$
(10)

The matrix product  $\mathbf{B}_{(\widehat{\mathbf{w}})} = \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} \mathbf{F}_{(\widehat{\mathbf{w}})} \mathbf{F}_{(\widehat{\mathbf{w}})}^{T} \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2}$  in Equation (5) is block diagonal (univariate) and can be interpreted as a background-error covariance matrix for the vector  $\widehat{\mathbf{w}}^{b} = (T^{b}, S_{U}^{b})^{T}$  where  $S_{U}^{b}$  is an 'unbalanced' background salinity variable that is constructed to be

approximately uncorrelated with  $T^{b}$  (Weaver *et al.*, 2005). The transformation of background errors from  $\widehat{\mathbf{w}}$ -space to w-space is achieved using the linear balance operator  $\mathbf{K}_{(\mathbf{w})}$ . Here,  $\mathbf{K}_{(\mathbf{w})}$  is formulated so that it leaves  $T^{b}$ errors unchanged but estimates  $S^{b}$  errors as the sum of balanced  $(S^{b}_{B})$  and unbalanced  $(S^{b}_{U})$  errors where the balanced component is computed directly from  $T^{b}$  errors using the operator  $\mathbf{K}_{ST}$ . Following Ricci et al. (2005),  $\mathbf{K}_{ST}$  has been parametrized in terms of the vertical gradients of  $T^{b}$  and  $S^{b}$  so that local salinity changes can be produced in response to local temperature changes to allow approximate preservation of the background watermass properties. The degree to which the water-mass properties are preserved is controlled by the backgrounderror standard deviation matrices  $\mathbf{D}_T^{1/2} = \text{diag}\{\sigma_T^{b}\}$  and  $\mathbf{D}_{S_{\mathrm{U}}}^{1/2} = \mathrm{diag}\{\sigma_{S_{\mathrm{U}}}^{\mathrm{b}}\}$  where  $\sigma_{T}^{\mathrm{b}}$  and  $\sigma_{S_{\mathrm{U}}}^{\mathrm{b}}$  are row-vectors containing estimates of the standard deviations of temperature and unbalanced salinity background errors. The main purpose of this study is to explore the potential of an ensemble 3D-Var to provide flow-dependent estimates of these standard deviations.

The block matrices  $\mathbf{F}_{TT}$  and  $\mathbf{F}_{S_US_U}$  are 3D univariate smoothing operators, each constructed as the product of a 1D and 2D anisotropic diffusion operator (Weaver and Courtier, 2001). The product of  $\mathbf{F}_{(\widehat{\mathbf{w}})}$  with its adjoint  $\mathbf{F}_{(\widehat{\mathbf{w}})}^{\mathrm{T}}$ is, with appropriate normalization, a 3D correlation operator. The correlation functions implied by the diffusion model are approximately Gaussian. The parameters of the 3D diffusion model are the same as those used for the univariate T correlations in Weaver et al. (2003), except for the vertical correlation scales which have been slightly reduced here (they are proportional to the local vertical grid depths). Identical correlation parameters are used for T and  $S_{\rm U}$ . The ensemble 3D-Var could also be used to estimate parameters of the diffusion model although this interesting possibility goes beyond the scope of the current study.

The cost function *J* is minimized iteratively using a conjugate gradient algorithm (Fisher, 1998; Tshimanga *et al.*, 2008). To improve the convergence properties of the minimization, a preconditioning transformation  $\delta \mathbf{v} = \mathbf{U}_{(\mathbf{w})}^{-1} \delta \mathbf{w}$ , where  $\mathbf{U}_{(\mathbf{w})}^{-1} = \mathbf{F}_{(\widehat{\mathbf{w}})}^{-1/2} \mathbf{K}_{(\mathbf{w})}^{-1}$ , is employed in Equation (1) resulting in the modified cost function

$$J[\delta \mathbf{v}] = \frac{1}{2} \delta \mathbf{v}^{\mathrm{T}} \delta \mathbf{v} + \frac{1}{2} (\mathbf{H} \mathbf{U}_{(\mathbf{w})} \delta \mathbf{v} - \mathbf{d})^{\mathrm{T}} \mathbf{R}^{-1} (\mathbf{H} \mathbf{U}_{(\mathbf{w})} \delta \mathbf{v} - \mathbf{d}). \quad (11)$$

Forty iterations are performed on each assimilation cycle, which typically results in a reduction of nine orders of magnitude in the Euclidean norm of the gradient relative to its initial value. If  $\delta \mathbf{v}^a$  denotes the minimizing solution of Equation (11) then the minimizing solution of Equation (1) is determined from  $\delta \mathbf{w}^a = \mathbf{U}_{(\mathbf{w})} \, \delta \mathbf{v}^a$ . To produce balanced increments for the other model state variables  $\eta$ , u and v, a more general variable transform is applied to the solution  $\delta \mathbf{v}^a$ :

$$\delta \mathbf{x}^{a} = \mathbf{K}_{(\mathbf{x})} \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} \mathbf{F}_{(\widehat{\mathbf{w}})} \, \delta \mathbf{v}^{a} \,, \tag{12}$$

where  $\delta \mathbf{x}^{a} = (\delta T^{a}, \delta S^{a}, \delta \eta^{a}, \delta u^{a}, \delta v^{a})^{T}$  is the analysis increment for the complete model state vector, and

$$\mathbf{K}_{(\mathbf{x})} = \begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{K}_{ST} & \mathbf{I} \\ \mathbf{K}_{\eta T} & \mathbf{K}_{\eta S} \\ \mathbf{K}_{uT} & \mathbf{K}_{uS} \\ \mathbf{K}_{vT} & \mathbf{K}_{vS} \end{pmatrix}$$
(13)

is the full balance operator. The matrix

$$\mathbf{B} = \mathbf{U} \mathbf{U}^{\mathrm{T}},\tag{14}$$

where

$$\mathbf{U} = \mathbf{K}_{(\mathbf{x})} \, \mathbf{D}_{(\widehat{\mathbf{w}})}^{1/2} \, \mathbf{F}_{(\widehat{\mathbf{w}})}, \tag{15}$$

can be interpreted as a reduced-rank error covariance matrix for the complete background state  $\mathbf{x}^{b}$ . The operators  $\mathbf{K}_{\eta T}$  and  $\mathbf{K}_{\eta S}$  in Equation (13) compute a balanced SSH increment,  $\delta \eta^{a}$ , by integrating a density increment from a reference depth (1500 m) to the surface, where the density increment is computed from  $\delta T^{a}$  and  $\delta S^{a}$  using a linearized equation of state. The operators  $\mathbf{K}_{uT}$ ,  $\mathbf{K}_{uS}$ ,  $\mathbf{K}_{vT}$  and  $\mathbf{K}_{vS}$  compute balanced horizontal velocity increments,  $\delta u^{a}$  and  $\delta v^{a}$ , from the geostrophic relation. Near the Equator,  $\delta v^{a}$  is reduced to zero while  $\delta u^{a}$  is balanced geostrophically using a  $\beta$ -plane approximation (Lagerloef *et al.*, 1999). A detailed description of the multivariate balance operator can be found in Weaver *et al.* (2005).

The increment  $\delta \mathbf{w}$  and background-error covariance matrix  $\mathbf{B}_{(w)}$  are formally defined with respect to  $w^b$ . In 3D-Var,  $\mathbf{w}^{b}$  can be chosen from any background state  $\mathbf{w}^{b}(t_{i})$  within the window. It is usually taken in the middle of the window to minimize the effects of the approximations in 3D-Var. This is particularly important in non-FGAT formulations where this static background state is used to compare directly with observations occurring at different times (i.e. using  $\mathbf{w}^{b}$  instead of  $\mathbf{w}^{b}(t_{i})$  in Equation (2)). Here, following Weaver et al. (2003) and Ricci et al. (2005), we take, as in 4D-Var,  $\mathbf{w}^{b}$  to be the background state at the start of the window ( $\mathbf{w}^{b} \equiv \mathbf{w}^{b}(t_{0})$ ). This choice was made mainly for simplifying the technical implementation of 3D-Var in our system, which also supports 4D-Var. It also provides a useful interpretation of 3D-Var as a limiting case of incremental 4D-Var in which the tangent-linear operator that propagates the increment in 4D-Var is replaced by the identity operator in 3D-Var. The background state  $\mathbf{w}^{b}$  is used here in two places: first, to define the linearization state in the T-S balance  $(\mathbf{K}_{ST})$ ; and second, in the parametrized formulation of the background-error variances with respect to which the ensemble-generated variances will be compared.

The technique of Incremental Analysis Updates (IAU; Bloom *et al.* 1996) is used to introduce the analysis increment gradually into the ocean model in order to minimize spurious adjustment processes. In this study, IAU is applied over the entire window; i.e. given  $\delta \mathbf{x}^a$ , the model integration from  $t_0$  to  $t_N$  is repeated using a prognostic equation of the form

$$\mathbf{x}^{a}(t_{i}) = M(t_{i}, t_{i-1})[\mathbf{x}^{a}(t_{i-1}), \mathbf{f}_{i}] + F_{i}\delta\mathbf{x}^{a}, \qquad (16)$$

where  $\mathbf{x}^{a}(t_{0}) = \mathbf{x}^{b}(t_{0})$ , and  $F_{i}$  is a weighting function defined such that  $\sum_{i=1}^{N} F_{i} = 1$ . The weighting function has been formulated to give maximum weight in the centre of the window, with the weighting reduced linearly to a small value at the window end-points. Such weighting provides a smooth transition of the analysis trajectory from one assimilation cycle to the next. An assimilation window of  $t_{N} = 10$  days has been used for the experiments in this study. Note that, in the ensemble system, the entire 3D-Var cycle (the integration of Equations (3) and (16) and the computation of the analysis increment via Equation (12)) must be performed separately for each ensemble member l, as described in the next section.

## **3.** Specification of the observation- and backgrounderror variances

## 3.1. Observation-error variance matrix: $\mathbf{D}_{(\mathbf{y})}$

Two formulations of the observation-error variance matrix have been tested in this study. The first formulation, denoted  $\mathbf{D}_{(\mathbf{y})}^{(1)}$ , is based on a simple analytical function that depends, except near coastlines, on depth only. The function has been constructed to provide an approximate fit to the vertical profiles of globally averaged temperature and salinity observation-error standard deviations  $(\sigma^{o})$  estimated by Ingleby and Huddleston (2007) (their Table 3). For temperature,  $\sigma^{o}$  is a maximum at 75 m depth where it has a value of 1 °C compared to 0.75 °C at the surface and its minimum value of 0.07 °C in the deep ocean. For salinity,  $\sigma^{o}$  decreases exponentially with depth from 0.18 psu at the surface to a minimum value of 0.02 psu in the deep ocean. Near coastlines, where our coarse resolution model is a poor representation of the real ocean, the  $\sigma^{o}$  profiles have been inflated. The inflation factor has been set to a value of two directly at the coastline and decreases smoothly to a value of one beyond 300 km of the coastline.

The second formulation, denoted  $\mathbf{D}_{(\mathbf{y})}^{(2)}$ , employs geographically dependent temperature and salinity  $\sigma^{\circ}$  that have been estimated using a statistical method originally proposed by Fu *et al.* (1993). The method has been widely used in ocean data assimilation (Fukumori, 2000; Menemenlis and Chechelnitsky, 2000; Leeuwenburgh, 2007). Given a vector  $\mathbf{w}^{c} = (T^{c}, S^{c})^{T}$  of temperature and salinity fields computed from a model integration without data assimilation (the control run in this study), the Fu *et al.* method estimates the observation-error variances from the covariance between co-located observation and observation-minus-control anomalies:

$$\mathbf{D}_{(\mathbf{y})}^{(2)} = \operatorname{diag}\left\{\overline{\mathbf{y}_{i}^{o'}\left(\mathbf{y}_{i}^{o'}-\mathbf{H}_{i}\mathbf{w}_{i}^{c'}\right)^{\mathrm{T}}}\right\},\qquad(17)$$

where the overbar indicates an appropriate time and spatial average, and the prime indicates anomaly with respect to this average. Appendix A provides a derivation of Equation (17) and a discussion of the various assumptions involved.

The variance computation has been performed using all in situ data between January 1962 and December 2002 contained in the ENSEMBLES data-set (section 2.2). Estimates have been made at each model grid point by averaging covariances within that model grid cell. In some regions, such as the deep ocean and Southern Hemisphere, the  $\sigma^{0}$  are grossly underestimated due to the sparseness of the data. To avoid this problem, the Ingleby and Huddleston variances were imposed as minimum values. The  $\sigma^{o}$  were then smoothed in each level by applying a local two grid-point Shapiro filter. Finally, the model-gridded  $\sigma^{o}$  were interpolated to the observation locations using the observation operator, and inflated near coastlines as in  $\mathbf{D}_{(\mathbf{y})}^{(1)}$ . Both  $\mathbf{D}_{(\mathbf{y})}^{(1)}$  and  $\mathbf{D}_{(\mathbf{y})}^{(2)}$  provide estimates of only the stationary component of  $\sigma^{\circ}$ . No attempt was made to estimate a time-varying component of  $\sigma^{o}$  due to the sparseness of the data.

The globally averaged profiles of  $\sigma^{o}$  computed from Equation (17) (not shown) have similar characteristics to those of Ingleby and Huddleston (2007) although are noticeably larger above 1500 m. For temperature, the largest difference between the two estimates is 0.3 °C and occurs near the maximum value of  $\sigma^{\circ}$  at 75 m. For salinity, the largest difference is 0.05 psu and occurs at the surface. The geographical distribution of  $\sigma^{o}$  is illustrated in Figure 2 in Daget et al. (2008). The largest  $\sigma^{\circ}$  (up to 3 °C) occur in regions characterized by strong internal variability. In particular, large values are obtained in the thermocline in the tropical Pacific and Atlantic Oceans, as well as in western boundary current regions (in particular, the Gulf Stream, Kuroshio, Agulhas and Malvinas Current regions) where there is significant mesoscale activity that our coarse-resolution model cannot resolve and thus where representativeness error is large. This important feature is absent in  $\mathbf{D}_{(\mathbf{v})}^{(1)}$ .

## 3.2. Background-error variance matrix: $\mathbf{D}_{(\widehat{\mathbf{w}})}$

Two flow-dependent formulations of the backgrounderror variance matrix have been tested in this study. The first formulation, denoted  $\mathbf{D}_{(\widehat{\mathbf{w}})}^{(1)}$ , is based on an empirical parametrization, while the second formulation, denoted  $\mathbf{D}_{(\widehat{\mathbf{w}})}^{(2)}$ , is derived from an ensemble method. The formulations are described in detail in the remainder of this section.

# 3.2.1. Parametrized error variance matrix: $\mathbf{D}_{(\widehat{\mathbf{w}})}^{(1)}$

For temperature, the background-error standard-error deviations ( $\sigma^{b}$ ) are parametrized in terms of the vertical gradient of the background temperature field so that large  $\sigma^{b}$  are concentrated at the level of the thermocline where thermal variability is greatest. Weaver *et al.* (2003) illustrate how this simple variance parametrization can capture some of the dynamical effects implicit in 4D-Var. A similar parametrization is used in the operational ocean data assimilation systems at the National Centers for Environmental Prediction (NCEP; Behringer *et al.*, 1998) and

ECMWF (Balmaseda *et al.*, 2008). The parametrization is described by the equation

$$\sigma_T^{\rm b} = \begin{cases} \max\left(\widetilde{\sigma}_T^{\rm b}, \ \sigma_T^{\rm ml}\right) & \text{in the mixed layer,} \\ \max\left(\widetilde{\sigma}_T^{\rm b}, \ \sigma_T^{\rm do}\right) & \text{below the mixed layer,} \end{cases}$$
(18)

where

$$\widetilde{\sigma}_{T}^{b} = \min\left\{ \left| \left( \frac{\partial T}{\partial z} \right|_{T=T^{b}} \right) \delta z \right|, \ \sigma_{T}^{\max} \right\}, \qquad (19)$$

 $\sigma_T^{\text{max}}$  being the maximum-allowed value of  $\sigma_T^{\text{b}}$ ,  $\delta z$  a vertical scale, and  $\sigma_T^{\text{ml}}$  and  $\sigma_T^{\text{do}}$  lower bounds in the mixed layer and deep ocean, respectively. In this study, as in Weaver *et al.* (2005),  $\sigma_T^{\text{max}} = 1.5 \,^{\circ}\text{C}$ ,  $\delta z = 10 \,\text{m}$ ,  $\sigma_T^{\text{ml}} = 0.5 \,^{\circ}\text{C}$ , and  $\sigma_T^{\text{do}} = 0.07 \,^{\circ}\text{C}$ . Finally, the  $\sigma_T^{\text{b}}$  were smoothed in each model level using a diffusion (Gaussian) filter with geographically dependent length-scales identical to those specified in the horizontal correlation operator.

For unbalanced salinity,  $\sigma^{b}$  is defined in a somewhat *ad hoc* fashion according to the equation

$$\sigma_{S_{\rm U}}^{\rm b} = \begin{cases} \sigma_{S_{\rm U}}^{\rm max} & \text{if } z \ge z_{\rm max} ,\\ \sigma_{S_{\rm U}}^{\rm max} \alpha(z) & \text{if } z < z_{\rm max} , \end{cases}$$
(20)

where  $\sigma_{S_U}^{\text{max}} = 0.25$  psu,  $z_{\text{max}}$  is the depth of the maximum of

$$\left| \left( \frac{\partial S}{\partial T} \right|_{T=T^{b}} \right) \right| \equiv \left| \left( \frac{\partial S}{\partial z} \right|_{S=S^{b}} \right) \left( \frac{\partial z}{\partial T} \right|_{T=T^{b}} \right) \right|,$$

and

$$\alpha(z) = 0.1 + 0.45 \times \left[1 - \tanh\left\{2 \ln\left(z/z_{\max}\right)\right\}\right]. \quad (21)$$

The above parametrization thus defines the largest  $\sigma_{S_U}^b$ between the surface and the level of maximum S(T)gradients, and decreases  $\sigma_{S_U}^b$  monotonically below this level. The large values in the mixed layer are especially important since there salinity is described primarily by its unbalanced component (Ricci *et al.* 2005). The empirical formulation of  $\sigma^b$  will serve as a reference for evaluating the ensemble-generated  $\sigma^b$  described below.

# 3.2.2. Ensemble-estimated error variance matrix: $\mathbf{D}_{(\widehat{\mathbf{w}})}^{(2)}$

The ensemble method employed in this study is similar to that used in the meteorological variational data assimilation studies of Fisher (2003), Žagar *et al.* (2005) and Berre *et al.* (2006). Appendix B provides the mathematical basis of the method. In particular, it is shown how perturbing the input parameters of a cycled analysis/forecast system leads to linearized evolution equations for the analysis and forecast state perturbations which are identical to those for the true errors. Furthermore, assuming that the perturbations to the input parameters are random samples drawn from the probability distribution of the true errors, then the evolved analysis and forecast perturbations from the cycled ensemble will also be random samples from the distribution of the true errors. The covariance matrices estimated from a sample of perturbed-minus-unperturbed analysis and forecast differences then provide accurate estimates of the true analysisand forecast-error covariance matrices. In practice, these covariance matrices will only be approximate due to the finite sample of the ensemble and due to inaccuracies in the specification of the error covariance matrix of the input parameters.

The method for cycling the ensemble analysis/forecast system is summarized schematically in Figure 1. Assuming that the errors in the different ensemble members are uncorrelated then, as discussed in Appendix B (Equation (B.25)),  $\mathbf{D}_{(\widehat{\mathbf{w}})}$  can be estimated from the difference between background states  $\mathbf{w}_l^{\rm b}(t_0)$  of successive ensemble members,  $l = 0, \ldots, L - 1$ :

$$\mathbf{D}_{(\widehat{\mathbf{w}})}^{(2)} = \text{diag} \left\{ \frac{1}{2(L-1)} \sum_{l=0}^{L-1} \left[ \mathbf{K}_{(\mathbf{w})}^{-1} \left\{ \mathbf{w}_{l}^{b}(t_{0}) - \mathbf{w}_{l+1}^{b}(t_{0}) \right\} \right] \\ \times \left[ \mathbf{K}_{(\mathbf{w})}^{-1} \left\{ \mathbf{w}_{l}^{b}(t_{0}) - \mathbf{w}_{l+1}^{b}(t_{0}) \right\} \right]^{\mathrm{T}} \right\}, \quad (22)$$

where

$$\mathbf{K}_{(\mathbf{w})}^{-1} = \begin{pmatrix} \mathbf{I} & 0\\ -\mathbf{K}_{ST} & \mathbf{I} \end{pmatrix}$$
(23)

and  $\mathbf{w}_{L}^{b}(t_{0}) = \mathbf{w}_{0}^{b}(t_{0})$ . Equation (22) can be related to Equation (B.25) by noting that  $\mathbf{w}_{l}^{b}(t_{0}) = \mathbf{w}_{l,c}^{b}(t_{0}) =$  $\mathbf{w}_{l,c-1}^{a}(t_{N})$  where *c* is the cycle number. Equation (23) is the inverse of the balance operator (Equation (9)) and is needed in order to estimate  $\sigma^{b}$  for  $\hat{\mathbf{w}}$  as required by the covariance model (Equation (5)).

Key to the design of the ensemble system is the construction of the perturbations for the system input parameters. In Appendix B, the ensemble method is developed while considering a general set of input parameters consisting of the external surface forcing fields, initial state, observations, and model-error source terms. Ideally, the perturbations should be chosen to sample the true statistical uncertainty in these parameters. The true error statistics of the input parameters are unknown and must be approximated in practice. In this study, the perturbations  $\tilde{\epsilon}_{l,i}^{f}$ , l = 1, ..., L - 1, to the surface fields (wind-stress, heat flux, PmE) are defined from differences between different analysis products (see below). The perturbations  $\widetilde{\epsilon}^{\mathrm{o}}_{l,i}$  to the observations are drawn from a Gaussian distribution with covariance matrix equal to the diagonal R-matrix used in the assimilation system. The background initial state perturbations  $\widetilde{\epsilon}_{l}^{b}(t_{0}) = \widetilde{\epsilon}_{l}^{b}(t_{0})$ are set to zero on the first cycle (c = 1). On subsequent cycles, these perturbations are defined implicitly as the difference between the perturbed and unperturbed back-ground states ( $\tilde{\epsilon}_{l}^{b}(t_{0}) = \mathbf{x}_{l,c}^{b}(t_{0}) - \mathbf{x}_{0,c}^{b}(t_{0})$ ). Perturbations associated with model error  $\tilde{\epsilon}_{l,i}^{q}$  are neglected altogether in this study.

The perturbations to the surface forcing fields have been derived by ECMWF where they are used to produce ensembles of initial conditions for operational seasonal forecasting (Balmaseda *et al.*, 2008). They have also been used by various groups for ocean analysis


Figure 1. Schematic illustration of the ensemble 3D-Var system. The ensemble of analysis states  $\mathbf{x}_{l,c-1}^{a}(t_N)$ ,  $l = 0, \ldots L - 1$ , at the end of cycle c-1 are used to initialize the background trajectories of each ensemble member on the next cycle c. The background trajectory of each member l is produced by integrating the model with a perturbed set of forcing fields (wind stress, heat flux, PmE),  $\mathbf{f}_{l,c,i} = \mathbf{f}_{c,i} + \tilde{\mathbf{e}}_{l,c,i}^{o}$ , from the initial condition  $\mathbf{x}_{l,c}^{b}(t_0) = \mathbf{x}_{l,c-1}^{a}(t_N)$ . Each background trajectory is compared with a set of perturbed observations  $\mathbf{y}_{l,c,i}^{o} = \mathbf{y}_{c,i}^{o} + \tilde{\mathbf{e}}_{l,c,i}^{o}$  to produce an innovation vector for each member l. A 3D-Var (FGAT) analysis is then performed for each ensemble member using the appropriate innovation vector and a background-error variance matrix  $\mathbf{D}_{(\widehat{\mathbf{w}}),c}$  that has been estimated from the ensemble of background initial states  $\mathbf{x}_{l,c}^{b}(t_0)$ . Ensemble member l = 0 is unperturbed:  $\tilde{\mathbf{e}}_{0,c,i}^{f} = \mathbf{0}$  and  $\tilde{\mathbf{e}}_{0,c,i}^{o} = \mathbf{0}$ . The resulting analysis increment is then used to produce an analysis state trajectory as described at the end of section 2.3.

production in the ENSEMBLES project. For wind stress, the perturbations are computed from differences between monthly mean anomalies from the ERA-40 and NCEP/National Center for Atmospheric Research (NCAR) reanalysis products. Perturbations to the fresh-water flux have been introduced in the precipitation field only, and are computed from differences between monthly mean anomalies of bias-corrected ERA-40 and NCEP/NCAR precipitation fields. To define the forcing perturbations for a given date and a given ensemble member, the perturbations are chosen randomly among the various difference fields that have the same calendar month (a sample of 44). Finally, daily perturbations of wind-stress and fresh-water flux are computed from the monthly fields using linear interpolation.

Perturbations of SST are used as a proxy for perturbations in heat flux, and are derived from differences between daily anomalies from different Reynolds products (2D-VAR and OIv2). The SST perturbations for a given date and ensemble member are constructed following the same random selection procedure used for the wind-stress and fresh-water flux perturbations. The procedure leads to a set of daily SST perturbations that, for a given member, are uncorrelated from one day to the next. To remove the temporal discontinuity, the daily SST perturbations have been smoothed in time using a two-pass recursive filter which is equivalent to correlating the perturbations with a second-order auto-regressive function (Purser et al., 2003). A filtering time-scale of 7 days was used. The perturbations were then rescaled to ensure that the globally averaged standard deviation was the same before and after filtering.

Four sets of surface forcing field perturbations were generated using the procedure above. Eight perturbed forcing fields were then produced by adding and subtracting the four forcing perturbations from the unperturbed fields. A different set of randomly perturbed observations was defined for each of the eight branches involving different forcing fields. The eight perturbed branches plus the unperturbed branch give a 9-member ensemble. Variances computed from this small number of ensemble members were too noisy to be used directly in the assimilation system. In order to increase the sample size, a sliding window was used to include the ensemble of background states from the previous 9 cycles (90 days) in the computation of the variances for the current cycle. This effectively increased the ensemble size to 81. Assuming Gaussian statistics, the standard error in the estimated standard deviation for an ensemble size L is  $1/\sqrt{2L}$  (e.g. Barlow (1989), p. 89). Thus, with L = 81, the error is 8% compared to 24% with L = 9. A 17-member ensemble with four perturbed observation branches for each perturbed forcing branch was also tested (with and without a 9-cycle sliding window) but did not lead to noticeable improvements over the 9-member ensemble (with 9-cycle sliding window) to justify the extra computational cost. The use of a sliding window represents a compromise between the desire to have, on the one hand, truly flow-dependent background-error variances and, on the other, to reduce sampling error. In particular, with the 90-day window used here, background-error variations on intraseasonal time-scales are filtered out and those on seasonal time-scales are strongly damped.

Minimum values were set for the ensemble  $\sigma^{b}$  to avoid excessively small values in the deep ocean where the surface forcing perturbations and limited number of perturbed temperature and salinity profiles were not very effective in maintaining an adequate spread. The minimum-allowed values were taken to be 0.07 °C and



Figure 2. Vertical profiles of (a, b)  $\sigma^{b}$  and (c, d)  $\sigma^{o}$  for (a, c) temperature and (b, d) salinity in B1R1 (grey shaded areas), B1R2 (solid curves) and B2R2 (dashed curves). The solid and dashed curves coincide in (c, d). (e) and (f) show the corresponding ratios  $(\sigma^{b})^{2}[(\sigma^{b})^{2} + (\sigma^{o})^{2}]^{-1}$ . Both  $(\sigma^{b})^{2}$  and  $(\sigma^{o})^{2}$  have been computed at observation points, temporally averaged over the 1994–2000 period, and spatially averaged over the global region and within vertical model grid cells.

0.01 psu and correspond to globally averaged climatological estimates of  $\sigma_T^b$  and  $\sigma_S^b$  at 5000 m given by Ingleby and Huddleston (2007). These are also the deep-ocean values used in the parametrized  $\sigma^b$ , and only affect the ensemble spread below the thermocline (Figure 2).

## 4. Results

Four experiments were performed over the period 1993–2000 to test the sensitivity of the analyses to the different

background- and observation-error variance formulations presented in the previous section. Experiment B1R1 uses the parametrized  $\sigma^{b}$  and simplified  $\sigma^{o}$ . Experiment B1R2 uses the parametrized  $\sigma^{b}$  and the  $\sigma^{o}$  estimated using the Fu *et al.* method. The reanalysis experiments conducted by CERFACS in ENACT (Davey *et al.*, 2006) and ENSEMBLES used the variance specifications in B1R1 and B1R2, respectively. Experiments B1R1, B1R2 and the control (CTL) are our reference experiments. Experiment B2R2 uses the ensemble  $\sigma^{b}$ , and the  $\sigma^{o}$ from the Fu *et al.* method. The parametrized  $\sigma^{b}$  were

Table I. Summary of the background- and observation-error variance matrix formulations used in the different experiments. The matrices  $\mathbf{D}_{(\widehat{\mathbf{w}})}^{(1)}$  and  $\mathbf{D}_{(\widehat{\mathbf{w}})}^{(2)}$  contain the parametrized and ensemble-estimated background-error variances, respectively. The matrices  $\mathbf{D}_{(y)}^{(1)}$  and  $\mathbf{D}_{(y)}^{(2)}$  contain the simplified and Fu

et al. estimated observation-error variances, respectively.

Experiment name	$\mathbf{D}_{(\widehat{\mathbf{w}})}^{(1)}$	$\mathbf{D}^{(2)}_{(\widehat{\mathbf{w}})}$	$\mathbf{D}_{(\mathbf{y})}^{(1)}$	$\mathbf{D}_{(\mathbf{y})}^{(2)}$
B1R1	Х		Х	
B1R2	Х			Х
B2R2		Х		Х

used to initialize B2R2 on 1 January 1993 but were then replaced with the ensemble  $\sigma^{b}$  180 days after cycling. All time-averaged statistics presented in this section exclude the first year of the experiments. The different ensemble members of B2R2 produced statistically similar results. Unless stated otherwise, results from B2R2 will be presented from the unperturbed member. The assimilation experiments are summarized in Table I.

Our objective in this paper is to provide an overall assessment of the relative performance of the different experiments, so we focus mainly on globally averaged diagnostics in this section. An exception is in section 4.5 where results involving comparisons with independent data are presented for the northwest Atlantic and tropical Pacific regions.

## 4.1. Vertical profiles of $\sigma^{b}$ and $\sigma^{o}$

The vertical profiles of the prescribed  $\sigma^{b}$  and  $\sigma^{o}$  are illustrated in this section for the different experiments. For consistency with the observation-space diagnostics presented later in this paper, both  $\sigma^{b}$  and  $\sigma^{o}$  have been evaluated by first computing the variances  $(\sigma^{b})^{2}$  and  $(\sigma^{o})^{2}$  at observation points, averaging the variances in space and time, and then taking the square root to obtain the standard deviations. Here, the spatial averaging is performed over the global region and within the vertical model grid cells, and the time averaging is performed over the 1994–2000 period.

The specified background-error variances  $(\sigma^b)^2$  at observation points correspond to the diagonal elements of  $\mathbf{HB}_{(\mathbf{w})}\mathbf{H}^T$ . To compute the diagonal of  $\mathbf{HB}_{(\mathbf{w})}\mathbf{H}^T$  requires a specific algorithm since this matrix is only available in operator form in our system. The diagonal elements can be estimated at a reasonable cost using a randomization algorithm (Andersson *et al.*, 2000). Specifically, given an ensemble of Gaussian random vectors  $\mathbf{v}_m$ ,  $m = 1, \ldots, M$ , drawn from a population with zero mean and unit variance  $(E[\mathbf{v}_m] = 0 \text{ and } E[\mathbf{v}_m \mathbf{v}_m^T] = \mathbf{I}$  where  $E[\cdot]$  is the expectation operator) then

$$\mathbf{HB}_{(\mathbf{w})}\mathbf{H}^{\mathrm{T}} \approx \frac{1}{M-1} \sum_{m=1}^{M} (\mathbf{HU}_{(\mathbf{w})}\mathbf{v}_{m}) (\mathbf{HU}_{(\mathbf{w})}\mathbf{v}_{m})^{\mathrm{T}}, \quad (24)$$

where  $\mathbf{U}_{(\mathbf{w})}$  is given by Equation (10). On each cycle, Equation (24) was used with an ensemble of M = 100

Copyright © 2009 Royal Meteorological Society

random vectors to produce an estimate of  $\sigma^{b}$  at observation points, with an estimated error of approximately 7%.

Figure 2 shows vertical profiles of the specified  $\sigma^{b}$ and  $\sigma^{o}$  for temperature and salinity. At all depths, but especially in the upper 200 m, the ensemble-estimated  $\sigma^{b}$ of B2R2 are smaller than the parametrized  $\sigma^{b}$  of B1R1 and B1R2, while the Fu et al. estimated  $\sigma^{o}$  of B1R2 and B2R2 are larger than the simplified  $\sigma^{\circ}$  of B1R1. The ratio  $(\sigma^{b})^{2}[(\sigma^{b})^{2} + (\sigma^{o})^{2}]^{-1}$ , displayed in the lower panels, roughly indicates the average weight given to an innovation at a particular depth in determining the analysis increment (Equations (B.10) and (B.11)). For B2R2, the weights are noticeably smaller and more uniform with depth compared to those from B1R1 and B1R2. As a result, the analysis on each cycle of B2R2 will tend to remain closer to the background state than it will in either B1R1 or B1R2 which will tend to pull it more to the observations, especially in the upper 200 m. It is not possible *a priori* to say which of these  $\sigma^{b}$  and  $\sigma^{o}$ profiles are most appropriate. The diagnostics presented in the remainder of the paper will examine their relative impact on different aspects of the ocean analyses.

#### 4.2. Assimilation statistics

The innovation vector, **d** (Equation (2)), and analysis increment,  $\delta \mathbf{w}^{a}$ , provide valuable information for assessing the statistical performance and internal consistency of the assimilation system (Desroziers et al., 2005). In this section, we examine mean statistics of **d** and the analysis residual,  $\mathbf{r} = \mathbf{d} - \mathbf{H} \delta \mathbf{w}^{a}$ , where these vectors, with the time index omitted, are understood to contain the innovation vectors and analysis residuals from all cycles in the 1994–2000 period. The analysis residual r (simply called the *residual* in what follows) corresponds to the value of the difference field in the observation term of the 3D-Var FGAT cost function (Equation (1)) at the end of minimization. Whereas  $\mathbf{r}$  quantifies the fit to the data achieved by the assimilation method, it does not represent the actual fit to the data achieved after correcting the model integration using IAU, which is given by  $\tilde{\mathbf{r}} = \mathbf{y}^{o} - \mathbf{H}\mathbf{w}^{a}$ . By construction, the IAU procedure does not produce a close fit to the data near the beginning of each cycle so that, in general,  $\|\widetilde{\mathbf{r}}\| > \|\mathbf{r}\|$ .

Figure 3 shows the vertical profiles of the time mean of the globally averaged residual and innovation vector for temperature and salinity. A non-zero mean in the innovations and residuals is an indication of bias (systematic error) in the system (Dee and Todling, 2000; Balmaseda et al., 2007). In CTL there is a large negative bias above 200 m in both the temperature and salinity innovations (Figures 3(c, d)), where the model without data assimilation is, on average, too warm (up to  $0.7 \,^{\circ}$ C) and too salty (up to 0.6 psu) compared to observations. The mean temperature innovations change sign near 250 m, suggesting that the model is biased cold below this level. The mean salinity innovations are very small below 200 m, possibly as a result of the subsurface relaxation to climatology. The mean innovations are reduced substantially, especially for salinity, in all assimilation experiments. The



Figure 3. Vertical profiles of the 1994–2000 time mean of the globally averaged analysis (a, b) residuals ( $\mathbf{r} = \mathbf{d} - \mathbf{H}\delta\mathbf{w}^{a}$ ) and (c, d) innovations ( $\mathbf{d} = \mathbf{y}^{o} - \mathbf{H}\mathbf{w}^{b}$ ) for (a, c) temperature and (b, d) salinity for CTL (thin dotted curves), B1R1 (grey shaded areas), B1R2 (solid curves) and B2R2 (dashed curves). Values have been averaged onto model levels. For CTL the innovation and residual are identical ( $\delta\mathbf{w}^{a} = 0$ ).

mean residuals are slightly smaller than the mean innovations. They are smallest for B1R1 (grey shade) which is understandable since the  $\sigma^{o}$  in B1R1 are smaller than those used in B1R2 and B2R2, so that the assimilation method will tend to give more weight to the observations in B1R1 than in B1R2 and B2R2. In all experiments, the remaining biases, while much smaller than in CTL, are still significant, the largest being at the surface in B2R2 where the maximum innovation biases are approximately  $-0.3 \,^{\circ}$ C and -0.11 psu.

Figure 4 shows vertical profiles of the standard deviation (sd) of the residual and innovation vectors:

$$\operatorname{sd}(\mathbf{z}) = \sqrt{(\mathbf{z} - \overline{\mathbf{z}})^2}$$
 (25)

where  $\mathbf{z} = \mathbf{d}$ ,  $\mathbf{r}$  or  $\tilde{\mathbf{r}}$ , and the overbar indicates spatial average over the globe and within vertical model grid cells, and temporal average over the 1994–2000 period. The standard deviation indicates how well the model fits the observed temporal and spatial variability. CTL exhibits large errors in both temperature and salinity, particularly in the upper 150 m where signals associated with seasonal and interannual variability are largest. Maximum differences are 2.25 °C for temperature and 1.65 psu for salinity. Relative to CTL, all assimilation experiments improve the fit to the observed temperature and salinity variability at all depths. This is true on the global average (Figure 4) although in the equatorial Pacific (not shown) the salinity variability below 50 m was found to be slightly degraded in B1R1 and B1R2, but not in B2R2, which points to a deficiency in the parametrized estimates of  $\sigma_s^b$ . Differences between B1R1 and B1R2 are small (shaded and solid curves). B1R1 displays slightly smaller  $sd(\mathbf{r})$  in salinity around 100 m and in temperature at all depths, whereas B1R2 displays slightly smaller sd(d) in both temperature and salinity in the upper 100 m. This illustrates that a better fit to the data (achieved in B1R1 by reducing  $\sigma^{o}$ ) does not necessarily translate into a better model forecast. The differences arising from using the ensemble  $\sigma^{b}$  (B2R2; dashed curves) are larger, with both  $sd(\mathbf{r})$  and  $sd(\mathbf{d})$  being increased relative to those in B1R1 and B1R2, especially near the surface.

At first sight it appears that the use of the ensemble  $\sigma^{b}$  has slightly degraded the performance of the assimilation system. Closer inspection of Figure 4, however, reveals that while the innovations are larger in B2R2, the difference between the residuals and innovations is smaller than in B1R1 and B1R2, particularly in the upper 100 m where the difference is 0.1 °C and 0.05 psu smaller. This result indicates that the error growth in a 10-day forecast cycle is smaller in B2R2 than in B1R1 and B1R2, which in turn suggests that the analyses in B2R2 are better



Figure 4. As Figure 3, but showing the standard deviation of the analysis residuals  $(sd(\mathbf{r}))$  and innovations  $(sd(\mathbf{d}))$  as defined by Equation (25).

balanced. Since the error growth is compensated by the assimilation increment, a smaller error growth in B2R2 should be indicative of smaller increments in B2R2. This is confirmed by Figures 5(a, b) which display the vertical profiles of the root-mean-square (rms) of the analysis increments at observation points (rms( $H\delta w^a$ )) from the different experiments. For both temperature and salinity, the increments are smallest in B2R2 and largest in B1R1.

The smaller analysis increments in B2R2 could also be an indication that the assimilation system is underaffected by the observations. It is instructive therefore to compare the rms of the analysis increments with the actual 10-day forecast improvement as measured by the innovations. To do so, we define an 'efficiency' (E) index,

$$\mathbf{E} = \frac{\mathrm{rms}(\mathbf{d}_{\mathrm{c}}) - \mathrm{rms}(\mathbf{d})}{\mathrm{rms}(\mathbf{H}\delta\mathbf{w}^{\mathrm{a}})},$$
(26)

which measures the ratio of the difference between the rms of the 10-day forecast error from the control and from the assimilation experiment, to the 'work done' by the assimilation method (at observation points) to reduce the forecast error. Small (large) innovations and increments will act to increase (decrease) the E index. For example, one system will be more efficient than another (have a larger E value) if it can achieve, on average, a similar reduction in the innovations but with smaller increments. Positive (negative) values of the E index imply that the assimilation is beneficial (detrimental) to the model. Note that the index depends on the forecast lead-time (which influences the numerator in Equation (26)) as well as the width of the assimilation window (which influences the denominator in Equation (26)). Therefore, the E index cannot be used to compare experiments with different assimilation windows. Here, the forecast lead-time and assimilation window width are both equal to ten days. Note also that the E index is defined for any assimilation experiment that is affected by observations at least once, so the denominator is always non-zero.

Vertical profiles of the E index for the three assimilation experiments are shown in Figure 5. The E index is positive at all depths for all experiments, with highest values obtained in B2R2 and lowest values in B1R1. For temperature, the E index is largest at the mean level of the thermocline (100 m), whereas for salinity it is largest nearer the surface. In all experiments, there is a decrease in the temperature E index near the surface. This is related to the strong SST relaxation term used in both CTL and the assimilation experiments, which acts to reduce the value of the numerator in Equation (26).

## 4.3. Specified versus diagnosed $\sigma^{b}$ and $\sigma^{o}$

The difficulty in defining background- and observationerror statistics means that they are likely to be



Figure 5. Vertical profiles of (a, b) the rms of the assimilation increments at observation points  $(rms(H\delta w^a))$  and of (c, d) the efficiency index (Equation (26)) for (a, c) temperature and (b, d) salinity in B1R1 (grey shaded areas), B1R2 (solid curves) and B2R2 (dashed curves). Values have been averaged onto model levels.

incorrectly specified in a practical data assimilation system. Desroziers *et al.* (2005) discuss how the innovations and analysis increments generated by a data assimilation system can be used to diagnose *a posteriori* the covariances of observation error and background error in observation space. Assuming that the background and observation errors are mutually uncorrelated, and that their covariance matrices are good approximations to the true error covariance matrices, then the covariance matrix of the innovation vector satisfies

$$E\left[\mathbf{d}\mathbf{d}^{\mathrm{T}}\right] \approx \mathbf{H}\mathbf{B}_{(\mathbf{w})}\mathbf{H}^{\mathrm{T}} + \mathbf{R}.$$
 (27)

This classical result is easily derived using the expression for the innovation vector in terms of the background and observation errors, given by Equation (B.12) in Appendix B. Furthermore, using the analysis equation (Equation (B.10)), it is straightforward to show that the components of Equation (27) satisfy

$$E\left[\mathbf{d}\left(\mathbf{H}\delta\mathbf{w}^{\mathrm{a}}\right)^{\mathrm{T}}\right] \approx \mathbf{H}\mathbf{B}_{(\mathbf{w})}\mathbf{H}^{\mathrm{T}}$$
 (28)

and

$$E\left[\mathbf{d}\left(\mathbf{d}-\mathbf{H}\delta\mathbf{w}^{\mathrm{a}}\right)^{\mathrm{T}}\right] \approx \mathbf{R}.$$
 (29)

Copyright © 2009 Royal Meteorological Society

The left-hand sides of Equations (28) and (29) can be estimated using statistics from the assimilation system, while the right-hand sides of these equations are the specified covariance matrices presented earlier. In this section, these expressions are used to check the consistency of the specified *standard deviations* ( $\sigma^{b}$  and  $\sigma^{o}$ ) with those diagnosed using assimilation statistics. The analysis focuses on the time- and horizontally averaged component of the standard deviations. As in Equation (25), the mean bias has been removed from **d** and  $H\delta w^{a}$  in estimating the standard deviations from Equations (28) and (29).

Figure 6 shows vertical profiles from B2R2 of the specified  $\sigma^{b}$  and  $\sigma^{o}$  (solid curves) and the diagnosed  $\sigma^{b}$  and  $\sigma^{o}$  (dashed curves) estimated from Equations (28) and (29) using the innovation and analysis increments from all cycles between 1994 and 2000. The specified  $\sigma^{b}$  are identical to those displayed earlier in Figure 2 (dashed curves). In B2R2, the specified  $\sigma_{T}^{b}$  and  $\sigma_{S}^{b}$  are everywhere *underestimated* compared to the diagnosed values (Figures 6(a, c)), whereas the specified  $\sigma_{T}^{o}$  and  $\sigma_{S}^{o}$  are *overestimated* compared to the diagnosed values, apart from the upper 30 m where the  $\sigma_{S}^{o}$  are slightly underestimated (Figures 6(b, d)). The maximum specified-minusdiagnosed differences are  $-0.45 \,^{\circ}$ C and  $-0.4 \,^{\circ}$ psu for  $\sigma_{T}^{b}$ and  $\sigma_{S}^{b}$ , and  $0.4 \,^{\circ}$ C and 0.15 psu for  $\sigma_{T}^{o}$  and  $\sigma_{S}^{o}$ . It is interesting to note that the structure and amplitude of the



Figure 6. Vertical profiles of (a, c)  $\sigma^{b}$  and (b, d)  $\sigma^{o}$  for (a, b) temperature and (c, d) salinity in B2R2. Solid curves correspond to the  $\sigma^{b}$  and  $\sigma^{o}$  that were *specified* in the assimilation experiment; dashed curves correspond to the  $\sigma^{b}$  and  $\sigma^{o}$  that were *diagnosed a posteriori* using Equations (28) and (29). The specified  $\sigma^{b}$  and  $\sigma^{o}$  are identical to those displayed by dashed curves in Figure 2.

diagnosed  $\sigma_T^b$ , and to a lesser extent the diagnosed  $\sigma_S^b$ , are closer to those of the parametrized  $\sigma_T^b$  and  $\sigma_S^b$  than the ensemble  $\sigma_T^b$  and  $\sigma_S^b$  (cf. Figure 2). The ensemble and diagnosed  $\sigma_S^b$  in particular exhibit large differences in the upper 200 m. Compared to B2R2, there is better consistency between the diagnosed and specified  $\sigma^b$  in B1R2 (Figures 7(a, c)), although this seems to be achieved at the expense of degrading the consistency between the diagnosed and specified  $\sigma^o$  (Figures 7(b, d)).

The results in Figure 6 suggest that the ensemble 3D-Var system produces background (and analysis) perturbations with inadequate spread on a global average. This apparent deficiency is not unique to our system but is a common problem in other ensemble data assimilation systems as well (e.g. Houtekamer and Mitchell (2005) give a discussion within the context of the EnKF). This issue is discussed further in section 5. The apparent overestimation of  $\sigma^{o}$ , on the other hand, points to limitations in our simple model of the observation-error covariances, which ignores spatial and temporal correlations and employs flow-independent variance estimates derived from a method that is itself subject to assumptions of questionable validity. Although Equation (29) is used purely for diagnostic purposes in this study, it provides the basis of an iterative algorithm for calibrating  $\sigma^{o}$  using the innovations and analysis increments generated by the

assimilation system (Desroziers *et al.*, 2005). In a similar way, Equation (28) can be used to calibrate observationspace values of  $\sigma^{b}$ . Unlike  $\sigma^{o}$ , however, these are not direct inputs to the ensemble 3D-Var system. How best to use Equation (28) to improve the estimates of  $\sigma^{b}$  in the space of the analysis variables and how to combine this information effectively with ensemble estimates of  $\sigma^{b}$  are open questions.

4.4. Temporal variability of the ensemble and assimilation statistics

The results presented so far have highlighted timeaveraged aspects of the assimilation performance. In this section, time-varying aspects will now be evaluated, focusing on results from the ensemble experiment B2R2. Figures 8(a, b) show time series of the 1993–2000 ensemble *spread* (the square root of the ensemble variance) of the observation-space analysis  $\mathbf{H}_i \mathbf{w}_i^{\text{a}}(t_i)$  and background  $\mathbf{H}_i \mathbf{w}_i^{\text{b}}(t_i)$ , computed with respect to all (L = 9) ensemble members:

spread{
$$\mathbf{H}_{i}\mathbf{w}^{a,b}$$
} =  
 $\sqrt{\frac{1}{L-1}\sum_{l=0}^{L-1} \left(\mathbf{H}_{i}\mathbf{w}_{l}^{a,b}(t_{i}) - \frac{1}{L}\sum_{l=0}^{L-1}\mathbf{H}_{i}\mathbf{w}_{l}^{a,b}(t_{i})\right)^{2}},$  (30)



Figure 7. As Figure 6, but for B1R2.

where the overbar indicates spatial average over the globe and within vertical model grid cells, and temporal average over 30-day intervals. A well-defined ensemble should have a spread characteristic of the actual uncertainty in the model state. Figures 8(a, b) show that the spread in both temperature and salinity is systematically smaller in the analysis than in the background, as one would expect. The spread appears to stabilize around mean values of 0.1 °C and 0.035 psu, after an initial increase during the first 6 months of the experiment. In other words, there is no evidence of ensemble collapse. The decrease in the spread from mid-1993 onwards corresponds to the time when the parametrized  $\sigma^{b}$  are replaced with the ensemble  $\sigma^{b}$ . The variability of the spread is larger for salinity than for temperature, which is mainly associated with increased sampling error due to the smaller number of salinity observations. It is interesting to note that the values of the mean spread are similar to those computed in the stochastic EnKF system of Leeuwenburgh (2007; his Figure 3 for the tropical Pacific region). His system was based on a different ocean model as well as a different assimilation method, but employed a similar perturbation strategy to ours, involving random perturbations to the atmospheric forcing fields and observations.

Figures 8(c, d) show corresponding time series of the  $sd(\tilde{\mathbf{r}}_i)$  and innovation  $sd(\mathbf{d}_i)$  of the unperturbed ensemble member l = 0, as given by Equation (25) but with the temporal averaging operator defined as in Equation (30).

Both sd( $\tilde{\mathbf{r}}_i$ ) and sd( $\mathbf{d}_i$ ) are about one order of magnitude larger than the analysis spread of the (observation space) analysis and background (Figures 8(a, b)). The spread of the background state at observation points roughly corresponds to the prescribed values of  $\sigma^{b}$  at observation points, as can be seen by comparing the magnitudes of the temperature and salinity spread in Figures 8(a, b) with those of the prescribed mean  $\sigma_T^{b}$  and  $\sigma_S^{b}$  profiles in Figures 2(a, b) (dashed curves). For both temperature and salinity, the magnitude of sd( $\mathbf{d}_i$ ) is at all times comparable to that of the mean  $\sigma^{o}$  in Figures 2(c, d), which is consistent with Equation (27) in view of the relatively small ensemble spread that defines  $\sigma^{b}$ . Despite the small spread, sd( $\mathbf{d}_i$ ) (and sd( $\tilde{\mathbf{r}}_i$ )) of B2R2 is consistently much smaller than sd( $\mathbf{d}_i$ ) of CTL.

#### 4.5. Comparison with independent data

The diagnostics presented in the previous sections have focused on the model variables (temperature and salinity) that are directly constrained by the observations. In this section, model variables (SSH and velocity) that are not directly constrained by the observations are examined and validated against independent data. First, SSH anomalies from the various experiments are compared with SSH anomalies from TOPEX/Poseidon (T/P), where the anomalies are computed with respect to the 1993– 2000 mean SSH of each product. Correlation coefficients



Figure 8. 1993–2000 time series of (a, b) the ensemble spread at observation points for the background,  $\mathbf{H}_i \mathbf{w}^{b}(t_i)$  (black shade), and analysis,  $\mathbf{H}_i \mathbf{w}^{a}(t_i)$  (light grey shade), in B2R2, and of (c, d) the standard deviation of the innovation vector,  $sd(\mathbf{d}_i)$  (black shade), and of the residuals,  $sd(\tilde{\mathbf{r}}_i)$  (light grey shade), in B2R2. The standard deviation of the innovation in CTL (dark grey shade) is also shown in (c, d). Temperature and salinity are displayed in (a, c) and (b, d), respectively. Values have been been computed for the global region and averaged into 30-day intervals.

Table II. Correlation coefficient and rms error in the northwest extratropical Atlantic (75–40°W, 30–60°N) and the NINO3.4 region of the tropical Pacific (170–120°W, 5°S–5°N) between SSH anomalies from TOPEX/Poseidon data and those from the model in the various experiments.

	NW.EXTROP.ATL		Ν	NINO3.4		
Experiment name	Correlation	Rms error (m)	Correlation	Rms error (m)		
CTL	0.97	0.012	0.98	0.022		
B1R1	0.62	0.040	0.99	0.012		
B1R2	0.73	0.033	0.99	0.012		
B2R2	0.87	0.023	0.99	0.013		

and rms errors are displayed in Table II for the northwest extratropical Atlantic and NINO3.4 in the tropical Pacific. In the northwest Atlantic, the CTL has the highest correlation and lowest rms error of all experiments, which suggests that data assimilation is degrading the SSH field to some extent in this region. Of the assimilation experiments, B2R2 compares best with T/P, while B1R1 compares worst. Since the closest fit to the *in situ* data was achieved in B1R1, followed by B1R2 and then B2R2 (Figures 3 and 4), this further suggests that the SSH field degrades in this region as the model fit to the *in situ* data improves. One possible explanation for this behaviour is the presence of model bias, which is not explicitly taken into account in the assimilation algorithm yet known to be important in this region. For example, Dee (2005) shows that the interaction of bias with a non-stationary observing



Figure 9. 1993–2000 time series of SSH anomalies in (a) the northwest extratropical Atlantic and (b) the NINO3.4 region of the tropical Pacific from CTL (dotted curve), B1R2 (thick solid curve), B2R2 (dashed curve), and TOPEX/Poseidon (thin solid curve).

system can lead to spurious time variability. Other factors such as inadequacies in the balance operator, incorrect background- or observation-error covariances, or simply 'bad' data could explain the degradation but pinpointing the specific reasons is difficult without conducting more targeted studies. In contrast, in NINO3.4, the assimilation experiments give similar statistical performance. Relative to CTL, they exhibit a slight improvement in correlation (the correlation of CTL is already very high) and a larger reduction in the rms error.

The 1993-2000 time series of the SSH anomalies in these regions, displayed in Figure 9, show clearly that the dominant variability is seasonal in the northwest Atlantic (Figure 9(a)) and interannual in the tropical Pacific (Figure 9(b)). Compared to T/P, the seasonal variations in the northwest Atlantic are reproduced in the assimilation experiments but with smaller amplitude, especially in B1R2 during 1996–1998. The observed seasonal variability is better reproduced in B2R2. Experiment B1R2 also displays a pronounced decreasing trend after 1999, which is weaker in B2R2 and not present in T/P. In NINO3.4, the interannual variations of CTL are slightly damped relative to those in T/P, especially during the 1997 El Niño event where the assimilation experiments, especially B1R2, reproduce the large amplitude of the observed SSH anomalies much better.

At the Equator in the Pacific, the quality of the velocity field can be assessed by comparing it to current meter data from the TAO array. Figure 10 shows vertical profiles of the correlation coefficients and rms errors between zonal current data from TAO at three locations (165°E, 140°W and 110°W) and the corresponding zonal velocity field from CTL, B1R2 and B2R2. The assimilation of temperature and salinity profiles improves the intensity of the equatorial surface currents and equatorial undercurrent in the central Pacific, as indicated by the reduced rms errors

in B1R2 and B2R2 compared to those of CTL in the upper 100 m at 140°W (Figure 10(d)). The upper ocean currents in B2R2 are also improved relative to CTL in the eastern Pacific (Figures 10(e, f)) but slightly degraded in the western Pacific (Figures 10(a, b)). The zonal currents in B1R2 are degraded in the upper 150 m of the western Pacific (Figures 10(a, b)), and show no clear improvement over CTL in the eastern Pacific (Figures 10(e, f)). Experiment B2R2 outperforms B1R2 at nearly all depths at all three locations.

#### 5. Summary and conclusions

An ensemble 3D-Var system for global ocean analysis has been described in this paper. The global 3D-Var system is based on an earlier 3D-Var system for the tropical Pacific (Weaver et al., 2003; Vialard et al., 2003) but includes many new features such as a fully multivariate background-error covariance model (Weaver et al., 2005), the use of a state-of-the-art qualitycontrolled in situ dataset (Ingleby and Huddleston, 2007), revised background- and observation-error variance formulations, and the capacity to generate ensembles of ocean analyses. On a given assimilation cycle, the ensemble of analyses are created by adding perturbations to the surface forcing fields (wind stress, fresh-water flux, and SST - a proxy for heat flux) and observations (temperature and salinity profiles) used in the assimilation process. These perturbations are based on estimates of the actual uncertainty in these input fields. The ocean initial conditions on each cycle are also perturbed, but this is done implicitly as a result of the parallel cycling of the 3D-Var system with different perturbed forcing and observations for each ensemble member. The purpose of the analysis ensemble is to sample uncertainty in the



Figure 10. (a, c, e) Correlation coefficient and (b, d, f) rms error (m s<sup>-1</sup>) between equatorial zonal currents from TAO data and those from CTL (dotted curves), B1R2 (solid curves) and B2R2 (dashed curves) over the 1993–2000 period at (a, b)  $165^{\circ}$ E, (c, d)  $140^{\circ}$ W and (e, f)  $110^{\circ}$ W.

ocean model state. Applications of the analysis ensemble include initialization of coupled ocean–atmosphere models for probabilistic climate forecasting, uncertainty estimation for historical ocean reanalysis, and the estimation of flow-dependent background-error covariances.

The main purpose of this paper was to explore the use of the ensemble 3D-Var for providing flow-dependent estimates of the background-error *standard deviations* ( $\sigma^{b}$ ). A 9-member ensemble was constructed and tested in a cycled 3D-Var framework over the period 1993–2000.

On each 10-day cycle, the  $\sigma^{b}$  of all members were updated based on the ensemble spread of background states. To reduce sampling error, a 9-cycle (90-day) sliding window was used to include additional ensemble members from the recent past in the  $\sigma^{b}$  computation. The larger sample size (81 in total) was achieved at the expense of filtering out intraseasonal variations in background error. This constraint could be relaxed in the future by increasing the number of ensemble members and/or by employing alternative filtering techniques for reduced sampling noise, such as those described in recent articles by Buehner and Charron (2007) and Berre *et al.* (2007).

A control experiment, in which no data were assimilated, produced large differences in the mean state and variability of the temperature and salinity fields when compared to the profile observations that were assimilated in the 3D-Var experiment. These differences were substantially reduced in the ensemble 3D-Var experiment. Evaluation of fields not directly constrained by the assimilated observations gave mixed results. Results showed that, in general, the ensemble 3D-Var experiment improved equatorial currents in the central and eastern Pacific and the representation of interannual variability of SSH. However, there were other regions where the assimilation degraded the results (equatorial currents in the western Pacific, SSH anomalies in the northwest Atlantic), possibly because of problems related to large systematic model error in these regions. Comparisons with a separate 3D-Var experiment that employed a simpler, empirically based flow-dependent  $\sigma^{b}$  parametrization showed that, on the global average, both led to similar reductions in the profile innovations (the mean and standard deviation) below 150 m but the parametrized  $\sigma^{\rm b}$  gave slightly smaller innovations above 150 m. Fields not directly constrained by the assimilated observations, however, were clearly better (closer to independent observations) with the ensemble  $\sigma^{b}$  than with the parametrized  $\sigma^{b}$ . Moreover, the error growth between assimilation cycles was much reduced using the ensemble  $\sigma^{b}$  suggesting that the ensemble  $\sigma^{b}$  produced analyses that were in better balance than those generated using the parametrized  $\sigma^{b}$ . This result could have important implications on the degree to which the assimilated information is retained by the model during the forecast step, but further investigation of this issue is needed; e.g. by computing statistics of the observation-minus-background differences on time periods that go beyond the 10-day forecast cycling period, or by testing the impact on seasonal forecasts using coupled models.

Diagnostics designed to check the consistency of the prescribed covariances with those estimated a posteriori from assimilation statistics indicated that the  $\sigma^{b}$  were underestimated by the ensemble, especially above 150 m. The parametrized  $\sigma^{b}$  were also underestimated but to a lesser extent. Simple procedures to inflate the ensemblegenerated  $\sigma^{b}$  in the upper ocean were tested (results not presented in this study) but did not give satisfactory results. The apparent underestimation of the ensemble spread is likely due to several factors including the small size of the ensemble and deficiencies in the perturbation strategy. The SST relaxation term, for example, has the tendency to produce excessive damping of temperature perturbations near the surface. The direct assimilation of SST data (via the cost function), which will be implemented in future versions of our assimilation system, should alleviate this problem. The absence of model-error perturbations is also a weakness, particularly for the relatively low-resolution model used in this study which can be expected to have a significant model-error

component associated with the unresolved mesoscale. Techniques to include model-error perturbations, such as those described by Hamill and Whitaker (2005), could be explored in future work. Despite these apparent shortcomings, results from this study are encouraging and suggest that useful information about background error can be extracted from a suboptimal ensemble.

This study has focused on using the ensemble to estimate the background-error standard deviations, but other parameters of the background-error covariance model could be estimated as well. Pannekoucke *et al.* (2008) present a practical method for estimating geographically dependent correlation length-scales from ensemble differences. In our quasi-Gaussian correlation model based on a diffusion operator, these length-scales are related to the elements of the diffusion tensor (Pannekoucke and Massart, 2008). Preliminary results from applying the Pannekoucke *et al.* method to estimate the tensor elements from time-averaged ensembles generated by our system are encouraging, although further work is needed to evaluate the impact of the new length-scale estimates in a cycled assimilation experiment.

The ensemble procedure has been tested in a 3D-Var framework in this study but is applicable to 4D-Var as well. In incremental 4D-Var, the background-error covariances are propagated implicitly within assimilation cycles using a linearized version of the model and its adjoint. This feature is absent in 3D-Var where the background-error covariances are stationary within assimilation cycles. An ensemble 4D-Var method could be used, as in ensemble 3D-Var, to propagate backgrounderror information between assimilation cycles, and is therefore complementary to the deterministic propagation of covariances achieved by 4D-Var within cycles. Practical ensemble 4D-Var applications, however, would likely require approximations in the ensemble-generation strategy due to the substantial extra cost of 4D-Var. In general, the extra computational expense of producing ensembles of analyses may be justified if these analyses can be used simultaneously for probabilistic forecasting as well as background-error covariance estimation.

## Acknowledgement

The development of the global ocean assimilation system was a contribution to the European projects ENACT (contract No. EVK2-CT-2001-00117) and ENSEMBLES (contract No. GOCE-CT-2003-505539). Additional support was obtained from the Groupe Mission Mercator/Coriolis (GMMC) and LEFE-ASSIM. We wish to thank Frédéric Vitart for his help in constructing the atmospheric forcing perturbations, and Kristian Mogensen for his help in designing the prepIFS/SMS suite that was used for running the ensemble experiments. We are grateful to two anonymous reviewers for their valuable comments that helped us improve the manuscript.

## Appendix A

Observation-error covariance estimation using the Fu et al. (1993) method

Following Janjić and Cohn (2006), we define the true state vector  $\mathbf{x}^{t}(t_i)$  at time  $t_i$  to be the component  $\Pi_{(\mathbf{x})} \mathbf{x}^{t}_{C}(t_i)$  of the true continuum state  $\mathbf{x}^{t}_{C}(t_i)$  where  $\Pi_{(\mathbf{x})}$  is a projection operator from the continuum onto the finite-dimensional subspace resolved by the numerical model. The resolved component is the quantity that we wish to estimate through data assimilation. The observation vector  $\mathbf{y}^{o}_{i}$  can be related to  $\mathbf{x}^{t}(t_i)$  through an equation of the form (Janjić and Cohn, 2006)

$$\mathbf{y}_{i}^{\mathrm{o}} = \mathbf{H}_{i}\mathbf{x}^{\mathrm{t}}(t_{i}) + \boldsymbol{\epsilon}_{i}^{\mathrm{m}} + \boldsymbol{\epsilon}_{i}^{\mathrm{r}} + \boldsymbol{\epsilon}_{i}^{\mathrm{i}}, \qquad (A.1)$$

where  $\mathbf{H}_i$  is the discrete observation operator which is taken to be linear as in the assimilation system described in section 2. The discrepancy between  $\mathbf{y}_i^0$  and  $\mathbf{H}_i \mathbf{x}^t(t_i)$ can be attributed to errors in the measurement process,  $\epsilon_i^{\mathrm{m}}$ , representativeness errors associated with the unresolved scales,  $\epsilon_i^{\mathrm{r}} = \mathcal{H}_i [\mathbf{x}_{\mathrm{C}}^{\mathrm{t}}(t_i) - \mathbf{x}^{\mathrm{t}}(t_i)]$ , where  $\mathcal{H}_i$  is the continuum observation operator, and interpolation errors associated with approximating the continuum observation operator by  $\mathbf{H}_i$ ,  $\epsilon_i^{\mathrm{i}} = (\mathcal{H}_i - \mathbf{H}_i) \mathbf{x}^{\mathrm{t}}(t_i)$ . The sum of these errors is the total observation error,

$$\epsilon_i^{\rm o} = \epsilon_i^{\rm m} + \epsilon_i^{\rm r} + \epsilon_i^{\rm i} \,. \tag{A.2}$$

The method of Fu *et al.* (1993) is designed to estimate the static component of the observation-error covariance matrix by comparing time-averaged statistics of the observation vector with those of its model equivalent  $\mathbf{H}_i \mathbf{x}^c(t_i)$ , where  $\mathbf{x}^c(t_i)$  is the state vector computed from a model integration without data assimilation (the control run in this study). At any time  $t_i$ , the unconstrained model state can be related to the true state through

$$\mathbf{x}^{c}(t_{i}) = \mathbf{x}^{t}(t_{i}) + \boldsymbol{\epsilon}^{c}(t_{i}), \qquad (A.3)$$

where  $\epsilon^{c}(t_i)$  represents the unconstrained model-state error. For notational convenience, the time parameter will be dropped in the rest of this appendix. Using Equations (A.1)–(A.3), the auto- and cross-covariances of  $\mathbf{y}^{o}$  and  $\mathbf{Hx}^{c}$  can be computed as follows:

$$E[\widetilde{\mathbf{y}}^{\circ}(\widetilde{\mathbf{y}}^{\circ})^{\mathrm{T}}] = \mathbf{H}E[\widetilde{\mathbf{x}}^{\mathrm{t}}(\widetilde{\mathbf{x}}^{\mathrm{t}})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}} + E[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}] + \mathbf{H}E[\widetilde{\mathbf{x}}^{\mathrm{t}}(\epsilon^{\circ})^{\mathrm{T}}] + E[\epsilon^{\circ}(\widetilde{\mathbf{x}}^{\mathrm{t}})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}},$$

$$(A.4)$$

$$E[\mathbf{H}\widetilde{\mathbf{x}}^{\circ}(\mathbf{H}\widetilde{\mathbf{x}}^{\circ})^{\mathrm{T}}] = \mathbf{H}E[\widetilde{\mathbf{x}}^{\mathrm{t}}(\widetilde{\mathbf{x}}^{\mathrm{t}})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}} + \mathbf{H}E[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}} + \mathbf{H}E[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}} + \mathbf{H}E[\widetilde{\mathbf{x}}^{\mathrm{t}}(\epsilon^{\circ})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}},$$

$$(A.5)$$

$$E[\widetilde{\mathbf{y}}^{\circ}(\mathbf{H}\widetilde{\mathbf{x}}^{\circ})^{\mathrm{T}}] = \mathbf{H}E[\widetilde{\mathbf{x}}^{\mathrm{t}}(\widetilde{\mathbf{x}}^{\mathrm{t}})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}} + \mathbf{H}E[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}]\mathbf{H}^{\mathrm{T}},$$

$$(A.6)$$

where E[.] denotes the expectation operator and  $\tilde{\mathbf{z}} = \mathbf{z} - E[\mathbf{z}]$ . The errors are assumed to be unbiased:  $E[\epsilon^{o}] =$ 

 $E[\epsilon^{c}] = 0$ . The *unknown* auto-covariance of the true state,  $E[\mathbf{\tilde{x}}^{t}(\mathbf{\tilde{x}}^{t})^{T}]$ , can be eliminated by subtracting Equation (A.6) from Equations (A.4) and (A.5) to yield

$$E\left[\widetilde{\mathbf{y}}^{\circ}(\widetilde{\mathbf{y}}^{\circ} - \mathbf{H}\widetilde{\mathbf{x}}^{c})^{\mathrm{T}}\right] = E\left[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}\right] + \mathbf{Z}_{1}, \qquad (A.7)$$
$$E\left[\mathbf{H}\widetilde{\mathbf{x}}^{\circ}(\widetilde{\mathbf{v}}^{\circ} - \mathbf{H}\widetilde{\mathbf{x}}^{c})^{\mathrm{T}}\right] = -\mathbf{H}E\left[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}\right]\mathbf{H}^{\mathrm{T}} + \mathbf{Z}_{2}, \quad (A.8)$$

where

$$\mathbf{Z}_{1} = \mathbf{H}E[\mathbf{\tilde{x}}^{t}(\boldsymbol{\epsilon}^{o})^{T}] - E[\boldsymbol{\epsilon}^{o}(\boldsymbol{\epsilon}^{c})^{T}]\mathbf{H}^{T} - \mathbf{H}E[\mathbf{\tilde{x}}^{t}(\boldsymbol{\epsilon}^{c})^{T}]\mathbf{H}^{T}, \qquad (A.9)$$
$$\mathbf{Z}_{2} = E[\boldsymbol{\epsilon}^{o}(\mathbf{\tilde{x}}^{t})^{T}]\mathbf{H}^{T} + E[\boldsymbol{\epsilon}^{o}(\boldsymbol{\epsilon}^{c})^{T}]\mathbf{H}^{T} - \mathbf{H}E[\boldsymbol{\epsilon}^{c}(\mathbf{\tilde{x}}^{t})^{T}]\mathbf{H}^{T}. \qquad (A.10)$$

The left-hand sides of Equations (A.7) and (A.8) are quantities that can be estimated, under the ergodic assumption, from time- and spatially averaged observations and their unconstrained model counterpart. The unknown quantity of interest here is the observation-error covariance matrix,  $\hat{\mathbf{R}} \equiv E[\epsilon^{\circ}(\epsilon^{\circ})^{\mathrm{T}}]$ , in Equation (A.7).

Following Fu et al. (1993), an approximate equation for  $\widehat{\mathbf{R}}$  is obtained by assuming that  $\mathbf{Z}_1 \approx 0$  or at least that this term is small (in a matrix-norm sense) compared to  $\mathbf{R}$ . The validity of this assumption can be appreciated by examining each term in Equation (A.9). The first term can be neglected by assuming that the true state (the resolved scales) and observation error are approximately uncorrelated,  $E[\mathbf{\tilde{x}}^{t}(\boldsymbol{\epsilon}^{o})^{T}] \approx 0$ . This is a safe assumption for the measurement component of the observation error which has no reason to be correlated with the true state. It also implies that the resolved and unresolved scales are entirely decoupled, which is more restrictive. The second term can be neglected by assuming that the observation error and unconstrained model-state error are approximately uncorrelated,  $E[\epsilon^{o}(\epsilon^{c})^{T}] \approx 0$ . From Equation (A.3), this is ensured if  $E[\epsilon^{o}(\mathbf{\tilde{x}}^{t})^{T}] \approx 0$ , as already discussed above, and if  $E[\epsilon^{o}(\tilde{\mathbf{x}}^{c})^{T}] \approx 0$ , which is a reasonable assumption since the model (control) integration is not constrained by the observations. It is more difficult, however, to justify ignoring the third term  $(E[\mathbf{\tilde{x}}^{t}(\boldsymbol{\epsilon}^{c})^{T}] \approx 0)$ , as already pointed out by Menemenlis and Chechelnitsky (2000) who provide evidence in their analysis of TOPEX/Poseidon altimeter data that suggests that this term is not negligible. This third assumption is made purely for practical convenience and should be treated with caution. Equation (A.7) (with  $\mathbf{Z}_1 = 0$ ) has been used in this study to estimate the variances of observation error (Equation (17)) although in principle it could be used to estimate the correlations as well.

The assumptions described above also imply that  $\mathbf{Z}_2 \approx 0$  in Equation (A.8), thereby yielding an approximate expression for  $\mathbf{HB}_{(\mathbf{x}^c)}\mathbf{H}^T$  where  $\mathbf{B}_{(\mathbf{x}^c)} \equiv E[\epsilon^c(\epsilon^c)^T]$ . Equation (A.8) may provide useful information for initializing the background-error covariance matrix on the first assimilation cycle, where the background state is obtained from an unconstrained model (spin-up) integration, but is of questionable relevance for defining background-error covariances in the presence of data assimilation. This

expression has not been exploited in this study where instead the ensemble method has been used to enrich a quasi-static covariance model with flow-dependent estimates of the variances.

#### Appendix B

Background-error covariance estimation using an ensemble method

The purpose of this appendix is to illustrate how differences between members of a suitably generated 3D-Var ensemble can be used to estimate the covariances of background error. First, expressions for the first-order evolution of the true background- and analysis-state errors will be derived. These expressions will then be related to the first-order evolution of background- and analysis-state perturbations in the ensemble system. The presentation is similar to that of Berre *et al.* (2006) but extended here to a nonlinear framework and tailored to account for particular features of our ensemble 3D-Var system.

# B.1. First-order evolution of the true background- and analysis-state errors

As discussed in section 2.3, the background state on a given 3D-Var cycle corresponds to the IAU-corrected state at the end of the previous cycle  $(\mathbf{x}_c^{\rm b}(t_0) = \mathbf{x}_{c-1}^{\rm a}(t_N))$ . For notational convenience, the index *c* will be ignored except when clarification is necessary.

The background state evolves from  $t_{i-1}$  and  $t_i$  according to Equation (3) where  $\mathbf{x}^{b}(t_0) = \mathbf{x}^{b}_{c}(t_0) = \mathbf{x}^{a}_{c-1}(t_N)$ . Using the notation established in Appendix A, the evolution of the true continuum state  $\mathbf{x}^{t}_{C}(t_i)$  can be described by the equation

$$\mathbf{x}_{\mathbf{C}}^{\mathsf{t}}(t_i) = \mathcal{M}(t_i, t_{i-1}) \left[ \mathbf{x}_{\mathbf{C}}^{\mathsf{t}}(t_{i-1}), \mathbf{f}_{\mathbf{C}, i}^{\mathsf{t}} \right], \qquad (B.1)$$

where  $\mathcal{M}(t_i, t_{i-1})$  is the true continuum model operator from  $t_{i-1}$  to  $t_i$ , and  $\mathbf{f}_{C,i}^t$  is the true continuum surface forcing vector acting from  $t_{i-1}$  to  $t_i$ . The evolution equation of the true resolved state  $\mathbf{x}^t(t_i) \equiv \Pi_{(\mathbf{x})} \mathbf{x}_C^t(t_i)$  can be represented in terms of the discrete model operator  $\mathcal{M}(t_i, t_{i-1})$ and the true resolved forcing vector  $\mathbf{f}_i^t \equiv \Pi_{(\mathbf{f})} \mathbf{f}_{C,i}^t$ , where  $\Pi_{(\mathbf{f})}$  is a projection operator from the atmospheric continuum onto the finite-dimensional subspace of the model forcing field, as

$$\mathbf{x}^{\mathsf{t}}(t_i) = M(t_i, t_{i-1}) \big[ \mathbf{x}^{\mathsf{t}}(t_{i-1}), \mathbf{f}_i^{\mathsf{t}} \big] - \boldsymbol{\epsilon}_i^{\mathsf{q}}, \qquad (B.2)$$

where  $\epsilon_i^{q}$  is the model error. Following Cohn (1997) and Janjić and Cohn (2006),  $\epsilon_i^{q}$  can be neatly expressed as the sum  $\epsilon_i^{q} = \epsilon_i^{qd} + \epsilon_i^{qu}$  where

$$\boldsymbol{\epsilon}_{i}^{\mathrm{qd}} = \left\{ \boldsymbol{M}(t_{i}, t_{i-1}) - \boldsymbol{\Pi}_{(\mathbf{X})} \mathcal{M}(t_{i}, t_{i-1}) \right\} \begin{bmatrix} \mathbf{X}^{\mathrm{t}}(t_{i-1}), \mathbf{f}_{i}^{\mathrm{t}} \end{bmatrix} \quad (\mathrm{B.3})$$

is the model error due to discretizaton, and

$$\boldsymbol{\epsilon}_{i}^{\text{qu}} = -\Pi_{(\mathbf{x})}\mathcal{M}(t_{i}, t_{i-1}) \Big[ \mathbf{x}_{\mathrm{C}}^{\text{t}}(t_{i-1}) - \mathbf{x}^{\text{t}}(t_{i-1}), \mathbf{f}_{\mathrm{C},i}^{\text{t}} - \mathbf{f}_{i}^{\text{t}} \Big]$$
(B.4)

Copyright © 2009 Royal Meteorological Society

is the model error due to the unresolved scales. Notice that our definition of model error through Equations (B.3) and (B.4) does not include the contribution from the surface forcing field error

$$\boldsymbol{\epsilon}_i^{\mathrm{f}} = \mathbf{f}_i - \mathbf{f}_i^{\mathrm{t}}, \qquad (\mathrm{B.5})$$

which is treated separately in what follows. The forcing errors include errors inherent in the (re)analysis procedure used to produce the atmospheric fluxes as well as errors associated with the interpolation procedure used to map the fluxes onto the model grid and time step.

An equation for the time evolution of the background error,

$$\boldsymbol{\epsilon}^{\mathrm{b}}(t_i) = \mathbf{x}^{\mathrm{b}}(t_i) - \mathbf{x}^{\mathrm{t}}(t_i), \qquad (B.6)$$

can be derived by subtracting Equation (B.2) from Equation (3) to yield

$$\mathbf{x}^{\mathbf{b}}(t_i) - \mathbf{x}^{\mathbf{t}}(t_i) = M(t_i, t_{i-1}) [\mathbf{x}^{\mathbf{b}}(t_{i-1}), \mathbf{f}_i] - M(t_i, t_{i-1}) [\mathbf{x}^{\mathbf{t}}(t_{i-1}), \mathbf{f}_i^{\mathbf{t}}] + \epsilon_i^{\mathbf{q}}.$$
 (B.7)

Expanding the second term on the right-hand side of Equation (B.7) about  $\mathbf{x}^{b}(t_{i-1})$  and  $\mathbf{f}_{i}$ , and using Equations (B.5) and (B.6), yields, to first order,

$$\boldsymbol{\epsilon}^{\mathsf{b}}(t_{i}) \approx \mathbf{M}_{\mathbf{x}^{\mathsf{b}}}(t_{i}, t_{i-1}) \boldsymbol{\epsilon}^{\mathsf{b}}(t_{i-1}) + \boldsymbol{\epsilon}_{i}^{\mathsf{p}}$$
$$\approx \mathbf{M}_{\mathbf{x}^{\mathsf{b}}}(t_{i}, t_{0}) \boldsymbol{\epsilon}^{\mathsf{b}}(t_{0}) + \sum_{j=1}^{i} \mathbf{M}_{\mathbf{x}^{\mathsf{b}}}(t_{i}, t_{j}) \boldsymbol{\epsilon}_{j}^{\mathsf{p}}, \quad (B.8)$$

where

$$\boldsymbol{\epsilon}_{i}^{\mathrm{p}} = \mathbf{M}_{\mathbf{f}}(t_{i}, t_{i-1}) \, \boldsymbol{\epsilon}_{i}^{\mathrm{f}} + \boldsymbol{\epsilon}_{i}^{\mathrm{q}}, \qquad (\mathrm{B.9})$$

is the total model error at time  $t_i$ . Here,

$$\begin{split} \mathbf{M}_{\mathbf{x}^{\mathbf{b}}}(t_{i}, t_{i-1}) &\equiv \left. \frac{\partial M}{\partial \mathbf{x}} \right|_{\mathbf{x}=\mathbf{x}^{\mathbf{b}}(t_{i-1})}, \\ \mathbf{M}_{\mathbf{x}^{\mathbf{b}}}(t_{i}, t_{j}) &\equiv \left. \mathbf{M}_{\mathbf{x}^{\mathbf{b}}}(t_{i}, t_{i-1}) \cdots \mathbf{M}_{\mathbf{x}^{\mathbf{b}}}(t_{j+1}, t_{j}) \right., \\ \mathbf{M}_{\mathbf{x}^{\mathbf{b}}}(t_{i}, t_{i}) &\equiv \mathbf{I}, \\ \text{and} \left. \mathbf{M}_{\mathbf{f}}(t_{i}, t_{i-1}) \right. &\equiv \left. \frac{\partial M}{\partial \mathbf{f}} \right|_{\mathbf{f}=\mathbf{f}_{i}}. \end{split}$$

The assimilation method transforms the innovation vector,  $\mathbf{d} = (..., \mathbf{d}_i^T, ...)^T$ , into an analysis increment. By minimizing the 3D-Var FGAT cost function exactly, the analysis increment can be expressed as

$$\delta \mathbf{x}^{\mathrm{a}} = \mathbf{K} \, \mathbf{d} \,, \tag{B.10}$$

where

$$\mathbf{K} = \mathbf{B} \, \mathbf{H}^{\mathrm{T}} \left( \mathbf{H} \mathbf{B} \, \mathbf{H}^{\mathrm{T}} + \mathbf{R} \right)^{-1}, \qquad (B.11)$$

is the Kalman gain matrix, **B** and **R** being the *pre-scribed* background- and observation-error covariance matrices detailed in section 2. The innovation vector can

be expressed in terms of the background error (Equation (B.6)) and observation error (Equations (A.1) and (A.2)) by noting that

$$\mathbf{d}_{i} = \mathbf{y}_{i}^{\mathrm{o}} - \mathbf{H}_{i} \mathbf{x}^{\mathrm{b}}(t_{i})$$
  
=  $\mathbf{y}_{i}^{\mathrm{o}} - \mathbf{H}_{i} \mathbf{x}^{\mathrm{t}}(t_{i}) + \mathbf{H}_{i} \mathbf{x}^{\mathrm{t}}(t_{i}) - \mathbf{H}_{i} \mathbf{x}^{\mathrm{b}}(t_{i})$   
=  $\boldsymbol{\epsilon}_{i}^{\mathrm{o}} - \mathbf{H}_{i} \boldsymbol{\epsilon}^{\mathrm{b}}(t_{i}).$  (B.12)

The analysis increment is applied to the model using IAU as described by Equation (16). The first-order evolution of the analysis error,

$$\boldsymbol{\epsilon}^{\mathrm{a}}(t_i) = \mathbf{x}^{\mathrm{a}}(t_i) - \mathbf{x}^{\mathrm{t}}(t_i), \qquad (B.13)$$

is obtained by subtracting Equation (B.2) from Equation (16) to yield

$$\mathbf{x}^{\mathbf{a}}(t_i) - \mathbf{x}^{\mathbf{t}}(t_i) = M(t_i, t_{i-1}) \left[ \mathbf{x}^{\mathbf{a}}(t_{i-1}), \mathbf{f}_i \right] + F_i \delta \mathbf{x}^{\mathbf{a}} - M(t_i, t_{i-1}) \left[ \mathbf{x}^{\mathbf{t}}(t_{i-1}), \mathbf{f}_i^{\mathbf{t}} \right] + \epsilon_i^{\mathbf{q}}.$$
 (B.14)

Expanding the third term on the right-hand side of Equation (B.14) about  $\mathbf{x}^{a}(t_{i-1})$  and  $\mathbf{f}_{i}$  and using Equation (B.9) gives, to first order,

$$\boldsymbol{\epsilon}^{\mathbf{a}}(t_{i}) \approx \mathbf{M}_{\mathbf{x}^{\mathbf{a}}}(t_{i}, t_{i-1}) \boldsymbol{\epsilon}^{\mathbf{a}}(t_{i-1}) + F_{i} \,\delta \mathbf{x}^{\mathbf{a}} + \boldsymbol{\epsilon}_{i}^{\mathbf{p}}$$
$$\approx \mathbf{M}_{\mathbf{x}^{\mathbf{a}}}(t_{i}, t_{0}) \boldsymbol{\epsilon}^{\mathbf{b}}(t_{0})$$
$$+ \sum_{j=1}^{i} \mathbf{M}_{\mathbf{x}^{\mathbf{a}}}(t_{i}, t_{j}) \left[ F_{j} \,\delta \mathbf{x}^{\mathbf{a}} + \boldsymbol{\epsilon}_{j}^{\mathbf{p}} \right], \quad (B.15)$$

where

and

$$\epsilon^{\mathbf{b}}(t_0) = \epsilon^{\mathbf{b}}_c(t_0) = \epsilon^{\mathbf{a}}_{c-1}(t_N)$$

Equation (B.15) is similar to Equation (B.8) for the background error but employs a different linearization state ( $\mathbf{x}^{a}(t_{i})$  instead of  $\mathbf{x}^{b}(t_{i})$ ) and includes the analysis increment as an extra component of 'model error'.

 $\mathbf{M}_{\mathbf{x}^{\mathrm{a}}}(t_{i}, t_{i-1}) \equiv \frac{\partial M}{\partial \mathbf{x}}$ 

## *B.2.* Ensemble representation of background- and analysis-state errors

Let the index *l* denote a particular ensemble member on a given cycle, and let  $\tilde{\epsilon}_{l}^{b}(t_{0})$ ,  $\tilde{\epsilon}_{l,i}^{f}$ ,  $\tilde{\epsilon}_{l,i}^{q}$  and  $\tilde{\epsilon}_{l,i}^{o}$  define a set of perturbations to the system input parameters such that

$$\mathbf{x}_{l}^{\mathbf{b}}(t_{0}) = \mathbf{x}^{\mathbf{b}}(t_{0}) + \widetilde{\boldsymbol{\epsilon}}_{l}^{\mathbf{b}}(t_{0}), \ \widetilde{\boldsymbol{\epsilon}}_{l}^{\mathbf{b}}(t_{0}) \sim N(0, \widetilde{\mathbf{P}}^{\mathbf{b}}(t_{0})), \ (\mathbf{B}.16)$$
  
$$\mathbf{f}_{l} = \mathbf{f}_{l} + \widetilde{\boldsymbol{\epsilon}}_{l}^{\mathbf{f}} \qquad \widetilde{\boldsymbol{\epsilon}}_{l}^{\mathbf{f}} \sim N(0, \widetilde{\mathbf{F}}_{l}) \qquad (\mathbf{B}.17)$$

$$\mathbf{q}_{l,l} = \mathbf{\tilde{c}}_{l}^{\mathbf{q}} + \mathbf{\tilde{c}}_{l,l}^{\mathbf{q}}, \qquad \mathbf{\tilde{c}}_{l,l}^{\mathbf{q}} + \mathbf{V}(\mathbf{0}, \mathbf{\tilde{r}}_{l}), \qquad (\mathbf{B}, \mathbf{R})$$

$$\mathbf{q}_{l,i} = \mathbf{e}_{l,i}, \qquad \mathbf{e}_{l,i} = \mathbf{v}(\mathbf{0}, \mathbf{Q}_l), \qquad (\mathbf{B}_{l,10})$$

$$\mathbf{y}_{l,i}^{0} = \mathbf{y}_{i}^{0} + \boldsymbol{\epsilon}_{l,i}^{0}, \qquad \boldsymbol{\epsilon}_{l,i}^{0} \sim N(0, \mathbf{R}_{i}), \qquad (B.19)$$

where we assume that the perturbations are normally distributed with  $E[\tilde{\epsilon}] = 0$  and  $E[\tilde{\epsilon}\tilde{\epsilon}^T] = \mathbf{A}$ . From Equations (3) and (16), the equations describing the time evolution of the perturbed background state  $\mathbf{x}_i^{\rm b}(t_i)$  and perturbed analysis state  $\mathbf{x}_i^{\rm a}(t_i)$  can be written as

$$\mathbf{x}_{l}^{b}(t_{i}) = M(t_{i}, t_{i-1}) \left[ \mathbf{x}_{l}^{b}(t_{i-1}), \mathbf{f}_{l,i} \right] + \mathbf{q}_{l,i} , \qquad (B.20)$$

$$\mathbf{x}_{l}^{a}(t_{i}) = M(t_{i}, t_{i-1}) \left[ \mathbf{x}_{l}^{a}(t_{i-1}), \mathbf{f}_{l,i} \right] + F_{i} \,\delta \mathbf{x}_{l}^{a} + \mathbf{q}_{l,i},$$
 (B.21)

Copyright © 2009 Royal Meteorological Society

where  $\mathbf{x}_{l}^{a}(t_{0}) = \mathbf{x}_{l}^{b}(t_{0})$ , and  $\delta \mathbf{x}_{l}^{a} = \mathbf{K} \mathbf{d}_{l}$ , with  $\mathbf{d}_{l} = (..., \mathbf{d}_{l,i}^{T}, ...)^{T}$  and  $\mathbf{d}_{l,i} = \mathbf{y}_{l,i}^{o} - \mathbf{H}_{i} \mathbf{x}_{l}^{b}(t_{i})$ , is the analysis increment produced using the perturbed observations and perturbed background trajectory of ensemble member *l*.

Subtracting Equation (3) from Equation (B.20) and Equation (16) from Equation (B.21), and linearizing terms, gives

$$\widetilde{\boldsymbol{\epsilon}}_{l}^{b}(t_{i}) \approx \mathbf{M}_{\mathbf{x}^{b}}(t_{i}, t_{0}) \widetilde{\boldsymbol{\epsilon}}_{l}^{b}(t_{0}) + \sum_{j=1}^{i} \mathbf{M}_{\mathbf{x}^{b}}(t_{i}, t_{j}) \widetilde{\boldsymbol{\epsilon}}_{l,j}^{p}, \qquad (B.22)$$
$$\widetilde{\boldsymbol{\epsilon}}_{l}^{a}(t_{i}) \approx \mathbf{M}_{\mathbf{x}^{a}}(t_{i}, t_{0}) \widetilde{\boldsymbol{\epsilon}}_{l}^{b}(t_{0}) + \sum_{j=1}^{i} \mathbf{M}_{\mathbf{x}^{a}}(t_{i}, t_{j}) \Big[ F_{j} \, \delta \widetilde{\mathbf{x}}_{l}^{a} + \widetilde{\boldsymbol{\epsilon}}_{l,j}^{p} \Big], (B.23)$$

where 
$$\widetilde{\boldsymbol{\epsilon}}_{l,i}^{p} = \mathbf{M}_{\mathbf{f}}(t_{i}, t_{i-1}) \widetilde{\boldsymbol{\epsilon}}_{l,i}^{f} + \widetilde{\boldsymbol{\epsilon}}_{l,i}^{q},$$
  
 $\delta \widetilde{\mathbf{x}}_{l}^{a} = \mathbf{K} \widetilde{\mathbf{d}}_{l},$   
 $\widetilde{\mathbf{d}}_{l} = (..., \widetilde{\mathbf{d}}_{l,i}^{T}, ...)^{T},$   
and  $\widetilde{\mathbf{d}}_{l,i} = \widetilde{\boldsymbol{\epsilon}}_{l,i}^{o} - \mathbf{H} \widetilde{\boldsymbol{\epsilon}}_{l}^{b}(t_{i}).$ 

Comparing Equations (B.22) and (B.23) with Equations (B.8) and (B.15) shows that the ensemble perturbations,  $\tilde{\epsilon}_l^{\rm b}(t_i)$  and  $\tilde{\epsilon}_l^{\rm a}(t_i)$ , and the true errors  $\epsilon^{\rm b}(t_i)$  and  $\epsilon^{\rm a}(t_i)$ , obey identical first-order evolution equations. Furthermore, if the covariance matrices of the input perturbations in Equations (B.16)–(B.19) are equal to the covariance matrices of the true errors,

$$\begin{aligned} \widehat{\mathbf{P}}^{\mathbf{b}}(t_0) &\equiv E\left[\epsilon^{\mathbf{b}}(t_0)(\epsilon^{\mathbf{b}}(t_0))^{\mathrm{T}}\right], \\ \widehat{\mathbf{F}}_i &\equiv E\left[\epsilon^{\mathrm{f}}_i(\epsilon^{\mathrm{f}}_i)^{\mathrm{T}}\right], \\ \widehat{\mathbf{Q}}_i &\equiv E\left[\epsilon^{\mathrm{f}}_i(\epsilon^{\mathrm{q}}_i)^{\mathrm{T}}\right], \\ \text{and} \ \widehat{\mathbf{R}}_i &\equiv E\left[\epsilon^{\mathrm{o}}_i(\epsilon^{\mathrm{o}}_i)^{\mathrm{T}}\right], \end{aligned}$$

then it follows from Equations (B.22) and (B.23) that the evolving covariance matrices

$$\widetilde{\mathbf{P}}^{b}(t_{i}) = E\left[\widetilde{\epsilon}_{l}^{b}(t_{i})(\widetilde{\epsilon}_{l}^{b}(t_{i}))^{T}\right]$$
  
and 
$$\widetilde{\mathbf{P}}^{a}(t_{i}) = E\left[\widetilde{\epsilon}_{l}^{a}(t_{i})(\widetilde{\epsilon}_{l}^{a}(t_{i}))^{T}\right]$$

will be identical to those of the true errors

a

$$\widehat{\mathbf{P}}^{b}(t_{i}) \equiv E\left[\epsilon^{b}(t_{i})(\epsilon^{b}(t_{i}))^{T}\right]$$
  
and 
$$\widehat{\mathbf{P}}^{a}(t_{i}) \equiv E\left[\epsilon^{a}(t_{i})(\epsilon^{a}(t_{i}))^{T}\right].$$

Of particular interest here is the covariance matrix

$$\widetilde{\mathbf{P}}^{a}(t_{N}) = E\left[\widetilde{\boldsymbol{\epsilon}}_{l}^{a}(t_{N})(\widetilde{\boldsymbol{\epsilon}}_{l}^{a}(t_{N}))^{\mathrm{T}}\right]$$

of the analysis-state error  $\tilde{\epsilon}^{a}(t_{N})$  at the end of the cycle, since this matrix should be used to define the background-error covariance matrix for the next cycle (Figure 1):  $\tilde{\mathbf{P}}^{a}(t_{N}) = \tilde{\mathbf{P}}^{a}_{c}(t_{N}) = \tilde{\mathbf{P}}^{b}_{c+1}(t_{0})$ . This matrix can

be estimated from a sample of L - 1 perturbed analysis states as

$$\widetilde{\mathbf{P}}^{\mathbf{a}}(t_N) \equiv \frac{1}{L-1} \sum_{l=1}^{L-1} \left\{ \mathbf{x}_l^{\mathbf{a}}(t_N) - \mathbf{x}^{\mathbf{a}}(t_N) \right\} \times \left\{ \mathbf{x}_l^{\mathbf{a}}(t_N) - \mathbf{x}^{\mathbf{a}}(t_N) \right\}^{\mathrm{T}}, \qquad (B.24)$$

where each  $\mathbf{x}_l^a(t_N)$  is generated by perturbing the system input parameters as in Equations (B.16)–(B.19). Rather than using Equation (B.24), Fisher (2003), Žagar *et al.* (2005) and Berre *et al.* (2006) suggest estimating  $\widetilde{\mathbf{P}}^a(t_N)$  from differences between ensemble members. Assuming that the errors of the different members are mutually uncorrelated, then

$$\widetilde{\mathbf{P}}^{\mathbf{a}}(t_N) = \frac{1}{2} E \left[ \left\{ \widetilde{\boldsymbol{\epsilon}}_l^{\mathbf{a}}(t_N) - \widetilde{\boldsymbol{\epsilon}}_{l+1}^{\mathbf{a}}(t_N) \right\} \times \left\{ \widetilde{\boldsymbol{\epsilon}}_l^{\mathbf{a}}(t_N) - \widetilde{\boldsymbol{\epsilon}}_{l+1}^{\mathbf{a}}(t_N) \right\}^{\mathrm{T}} \right] \\ \approx \frac{1}{2(L-1)} \sum_{l=0}^{L-1} \left\{ \mathbf{x}_l^{\mathbf{a}}(t_N) - \mathbf{x}_{l+1}^{\mathbf{a}}(t_N) \right\} \times \left\{ \mathbf{x}_l^{\mathbf{a}}(t_N) - \mathbf{x}_{l+1}^{\mathbf{a}}(t_N) \right\}^{\mathrm{T}}, \quad (B.25)$$

where  $\mathbf{x}_L^a(t_N) = \mathbf{x}_0^a(t_N) = \mathbf{x}^a(t_N)$ . The multiplicative factor 1/2 arises since ensemble members are effectively used twice in Equation (B.25). For historical reasons, Equation (B.25) rather than Equation (B.24) has been used in this study.

#### References

- Andersson E, Fisher M, Munro R, McNally A. 2000. Diagnosis of background errors for radiances and other observable quantities in a variational data assimilation system, and the explanation of a case of poor convergence. Q. J. R. Meteorol. Soc. 126: 1455–1472.
- Balmaseda MA, Dee DP, Vidard A, Anderson DLT. 2007. A multivariate treatment of bias for sequential data assimilation: application to the tropical oceans. *Q. J. R. Meteorol. Soc.* **133**: 167–179.
- Balmaseda MA, Vidard A, Anderson DLT. 2008. The ECMWF ocean analysis system: ORA-S3. Mon. Weather Rev. 136: 3018–3034.
- Barlow RJ. 1989. Statistics: A Guide to the Use of Statistical Methods in the Physical Sciences. Wiley: Chichester, UK.
- Behringer DW, Ji M, Leetmaa A. 1998. An improved coupled model for ENSO prediction and implications for ocean initialization. Part I: The ocean data assimilation system. *Mon. Weather Rev.* 126: 1013–1021.
- Belo Pereira M, Berre L. 2006. The use of an ensemble approach to study the background error covariances in a global NWP model. *Mon. Weather Rev.* **134**: 2466–2489.
- Berre L, Ştefănescu SE, Belo Pereira M. 2006. The representation of the analysis effect in three error simulation techniques. *Tellus* 58A: 196–209.
- Berre L, Pannekoucke O, Desroziers G, Ştefănescu SE, Chapnik B, Raynaud L. 2007. 'A variational assimilation ensemble and the spatial filtering of its error covariances: increase of sample size by local spatial averaging'. *Proceedings of workshop on flow-dependent* aspects of data assimilation. ECMWF: Reading, UK. pp 151–168.
- Bloom SC, Takacs LL, da Silva AM, Ledvina D. 1996. Data assimilation using Incremental Analysis Updates. *Mon. Weather Rev.* 124: 1256–1271.
- Buehner M. 2005. Ensemble-derived stationary and flow-dependent background-error covariances: Evaluation in a quasi-operational NWP setting. Q. J. R. Meteorol. Soc. 131: 1013–1044.

- Buehner M, Charron M. 2007. Spectral and spatial localization of background-error correlations for data assimilation. Q. J. R. Meteorol. Soc. 133: 615–630.
- Carton JA, Santorelli A. 2008. Global decadal upper-ocean heat content as viewed in nine analyses. J. Climate **21**: 6015–6035.
- Cohn SE. 1997. An introduction to estimation theory. J. Meteorol. Soc. Japan 75: 257–288.
- Conkright ME, Antonov JI, Baranova O, Boyer TP, Garcia HE, Gelfeld R, Johnson D, Locarnini RA, Murphy PP, O'Brien TD, Smolyar I, Stephens C. 2002. World Ocean Database 2001, Volume 1: Introduction. NOAA Atlas NESDIS 42. US Government Printing Office: Washington, DC.
- Courtier P, Thépaut JN, Hollingsworth A. 1994. A strategy for operational implementation of 4D-Var, using an incremental approach. Q. J. R. Meteorol. Soc. 120: 1367–1388.
- Daget N, Weaver AT, Balmaseda MA. 2008. 'An ensemble threedimensional variational data assimilation system for the global ocean: Sensitivity to the observation- and background-error variance formulation'. Tech. Memo. No. 562. ECMWF: Reading, UK. Available at http://www.ecmwf.int/publications/library/do/references/ list/14.
- Davey M, Huddleston M, Ingleby B, Haines K, Le Traon P-Y, Weaver AT, Vialard J, Anderson DLT, Troccoli A, Vidard A, Burgers G, Leeuwenburgh O, Bellucci A, Masina S, Bertino L, Korn P. 2006. Multi-model multi-method multi-decadal ocean analyses from the ENACT project. *CLIVAR Exchanges* 11: 22–25.
- Dee DP. 2005. Bias and data assimilation. Q. J. R. Meteorol. Soc. 131: 3323–3343.
- Dee DP, Todling R. 2000. Data assimilation in the presence of forecast bias: the GEOS moisture analysis. *Mon. Weather Rev.* **128**: 3268–3282.
- Derber J, Bouttier F. 1999. A reformulation of the background error covariance in the ECMWF global data assimilation system. *Tellus* **51A**: 195–221.
- Desroziers G, Berre L, Chapnik B, Poli P. 2005. Diagnosis of observation, background and analysis-error statistics in observation space. Q. J. R. Meteorol. Soc. 131: 3385–3396.
- Evensen G. 2007. Data Assimilation: The Ensemble Kalman Filter. Springer: Berlin.
- Fisher M. 1998. 'Minimization algorithms for variational data assimilation'. In Proceedings of seminar on recent developments in numerical methods for atmospheric modelling. ECMWF: Reading, UK. pp 364–385.
- Fisher M. 2003. 'Background error covariance modelling'. In Proceedings of seminar on recent developments in data assimilation for atmosphere and ocean. ECMWF: Reading, UK. pp 35–63.
- Fu LL, Fukumori I, Miller RN. 1993. Fitting dynamic models to the Geosat sea level observations in the tropical Pacific Ocean. Part II: A linear, wind-driven model. J. Phys. Oceanogr. 23: 2162–2181.
- Fukumori I. 2000. Data Assimilation by Models. In *Satellite Altimetry* and Earth Sciences: A Handbook of Techniques and Applications. Academic Press. pp 237–265.
- Gent PR, McWilliams JC. 1990. Isopycnal mixing in ocean circulation models. J. Phys. Oceanogr. 20: 150–155.
- Hamill TM, Snyder C. 2000. A hybrid ensemble Kalman fitler-3D variational analysis scheme. *Mon. Weather Rev.* 128: 2905–2919.
- Hamill TM, Whitaker JS. 2005. Accounting for the error due to unresolved scales in ensemble data assimilation: A comparison of different approaches. *Mon. Weather Rev.* 133: 3132–3147.
- Houtekamer PL, Mitchell HL. 2001. A sequential ensemble Kalman filter for atmospheric data assimilation *Mon. Weather Rev.* **129**: 123–137.
- Houtekamer PL, Mitchell HL. 2005. Ensemble Kalman filtering. Q. J. R. Meteorol. Soc. 131: 3269–3289.
- Ingleby B, Huddleston M. 2007. Quality control of ocean temperature and salinity profiles – Historical and real-time data. J. Mar. Sys. 65: 158–175.
- Janjić T, Cohn SE. 2006. Treatment of observation error due to unresolved scales in atmospheric data assimilation. *Mon. Weather Rev.* 134: 2900–2915.
- Johnson GC, Sloyan BM, Kessler WS, McTaggart KE. 2002. Direct measurements of upper ocean currents and water properties across the tropical Pacific during the 1990s. *Prog. Oceanogr.* 52: 31–36.
- Jones PW. 1998. A user's guide for SCRIP: A spherical coordinate remapping and interpolation package. Technical Report. Los Alamos National Laboratory: New Mexico, USA. http://climate.lanl. gov/Software/SCRIP/SCRIPusers.pdf.

- Keppenne CL, Rienecker MM. 2002. Initial testing of a massively parallel Ensemble Kalman Filter with the Poseidon isopycnal ocean general circulation model. *Mon. Weather Rev.* 130: 2951–2965.
- Küçükkaraca E, Fisher M. 2006. 'Use of analysis ensembles in estimating flow-dependent background error variances'. Tech. Memo. No. 492. ECMWF: Reading, UK. Available at http://www.ecmwf.int/publications/library/do/references/list/14.
- Lagerloef GSE, Mitchum GT, Lukas RB, Niiler PP. 1999. Tropical Pacific near-surface currents estimated from altimeter, wind, and drifter data. J. Geophys. Res. 104: 23313–23326.
- Leeuwenburgh O. 2007. Validation of an EnKF system for OGCM initialization assimilating temperature, salinity, and surface height measurements. *Mon. Weather Rev.* **135**: 125–139.
- Levitus S, Conkright ME, Boyer TP, O'Brien T, Antonov JI, Stephens C, Stathoplos L, Johnson D, Gelfeld R. 1998. World Ocean Database 1998a, Volume 1: Introduction. NOAA Atlas NESDIS 18. US Government Printing Office: Washington, DC.
- Lorenc AC. 2003a. Modelling of error covariances by 4D-Var data assimilation. Q. J. R. Meteorol. Soc. 129: 3167–3182.
- Lorenc AC. 2003b. The potential of the Ensemble Kalman Filter for NWP – A comparison with 4D-Var. Q. J. R. Meteorol. Soc. 129: 3183–3203.
- Madec G, Delecluse P, Imbard M, Levy C. 1998. 'OPA8.1 Ocean General Circulation Model reference manual'. Technical note No. 11. LODYC/IPSL. Université P. and M. Curie: Paris, France.
- Menemenlis D, Chechelnitsky M. 2000. Error estimates for an ocean general circulation model from altimeter and acoustic tomography data. *Mon. Weather Rev.* **128**: 763–785.
- Oke PR, Sakov P, Corney SP. 2007. Impacts of localisation in the EnKF and EnOI: Experiments with a small model. *Ocean Dyn.* **57**: 32–45.
- Ott E, Hunt BR, Szunyogh I, Zimin AV, Kostelich EJ, Corazza M, Kalnay E, Patil DJ, Yorke JA. 2004. A local ensemble Kalman filter for atmospheric data assimilation. *Tellus* **56A**: 415–428.
- Pannekoucke O, Massart S. 2008. Estimation of the local diffusion tensor and normalization for heterogeneous correlation modelling using a diffusion equation. Q. J. R. Meteorol. Soc. 134: 1425–1438.
- Pannekoucke O, Berre L, Desroziers G. 2008. Background error correlation length-scale estimates and their sampling statistics. Q. J. R. Meteorol. Soc. 134: 497–508.
- Purser RJ, Wu WS, Parrish DF, Roberts NM. 2003. Numerical aspect of the application of recursive filters to variational statistical analysis. Part I: Spatially homogeneous and isotropic Gaussian covariances. *Mon. Weather Rev.* 131: 1524–1535.
- Reynolds R, Rayner N, Smith T, Stokes D, Wang W. 2002. An improved *in situ* and satellite SST analysis for climate. *J. Climate* 15: 1609–1625.
- Ricci S, Weaver AT, Vialard J, Rogel P. 2005. Incorporating statedependent temperature-salinity constraints in the background-error covariance of variational ocean data assimilation. *Mon. Weather Rev.* 133: 317–338.

- Roullet G, Madec G. 2000. Salt conservation, free surface and varying volume: A new formulation for ocean GCMs. J. Geophys. Res. 105: 23927–23942.
- Troccoli A, Kållberg P. 2004. 'Precipitation correction in the ERA-40 reanalysis'. ERA-40 Project Report Series 13. ECMWF: Reading, UK.
- Tshimanga J, Gratton S, Weaver AT, Sartenaer A. 2008. Limitedmemory preconditioners, with application to incremental fourdimensional variational data assimilation. *Q. J. R. Meteorol. Soc.* 134: 753–771.
- Uppala SM, Kållberg PW, Simmons AJ, da Costa Bechtold U, Fiorino M, Gibson JK, Haseler J, Hernandez A, Kelly GA, Li X, Onogi K, Saarinen S, Sokka N, Allan RP, Andersson E, Arpe K, Balmaseda MA, Beljaars ACM, van de Berg L, Bidlot J, Bormann N, Caires S, Chevallier F, Dethof A, Dragosavac M, Fisher M, Fuentes M, Hagemann S, Holm E, Hoskins BJ, Isaksen L, Janssen PAEM, Jenne R, McNally AP, Mahfouf J-F, Morcrette J-J, Rayner NA, Saunders RW, Simon P, Sterl A, Trenberth KE, Untch A, Vasiljevic D, Viterbo P, Woollen J. 2005. The ERA-40 re-analysis. Q. J. R. Meteorol. Soc. 131: 2961–3012.
- Vialard J, Weaver AT, Anderson DLT, Delecluse P. 2003. Three- and four-dimensional variational assimilation with a general circulation model of the tropical Pacific Ocean. Part II: Physical validation. *Mon. Weather Rev.* **131**: 1379–1395.
- Vossepoel F, Weaver AT, Vialard J, Delecluse P. 2004. Adjustment of near-equatorial wind stress with four-dimensional variational data assimilation in a model of the Pacific Ocean. *Mon. Weather Rev.* 132: 2070–2083.
- Weaver AT, Courtier P. 2001. Correlation modelling on the sphere using a generalized diffusion equation. Q. J. R. Meteorol. Soc. 127: 1815–1846.
- Weaver AT, Vialard J, Anderson DLT. 2003. Three- and fourdimensional variational assimilation with a general circulation model of the tropical Pacific Ocean. Part I: Formulation, internal diagnostics, and consistency checks. *Mon. Weather Rev.* 131: 1360–1378.
- Weaver AT, Deltel C, Machu E, Ricci S, Daget N. 2005. A multivariate balance operator for variational ocean data assimilation. Q. J. R. Meteorol. Soc. 131: 3605–3625.
- Weisheimer A, Doblas-Reyes F, Rogel P, Da Costa E, Keenlyside N, Balmaseda MA, Murphy J, Smith D, Collins M, Bhaskaran B, Palmer TN. 2007. 'Initialisation strategies for decadal hindcasts for the 1960–2005 period within the ENSEMBLES project'. Tech. Memo. No. 521. ECMWF: Reading, UK. Available at http://www.ecmwf.int/publications/library/do/references/list/14.
- Žagar N, Andersson E, Fisher M. 2005. Balanced tropical data assimilation based on a study of equatorial waves in ECMWF shortrange forecast errors. Q. J. R. Meteorol. Soc. 131: 987–1011.



## On the diffusion equation and its application to isotropic and anisotropic correlation modelling in variational assimilation

A. T. Weaver<sup>\*</sup> and I. Mirouze

CERFACS/SUC URA 1875, Toulouse, France

\*Correspondence to: A. T. Weaver, CERFACS, 42 avenue Gaspard Coriolis, 31057 Toulouse Cedex 01, France. E-mail: anthony.weaver@cerfacs.fr

Differential operators derived from the explicit or implicit solution of a diffusion equation are widely used for modelling background-error correlations in geophysical applications of variational data assimilation. Key theoretical results underpinning the diffusion method are reviewed. Solutions to the isotropic diffusion problem on both the spherical space  $\mathbb{S}^2$  and the *d*-dimensional Euclidean space  $\mathbb{R}^d$  are considered first. In  $\mathbb{R}^d$  the correlation functions implied by explicit diffusion are approximately Gaussian, whereas those implied by implicit diffusion belong to the larger class of Matérn functions which contains the Gaussian function as a limiting case. The Daley length-scale, defined as  $D = \sqrt{-d/\nabla^2 c(r)}\Big|_{r=0}$  where  $\nabla^2$  is the *d*-dimensional Laplacian operator and  $r = |\mathbf{r}|$  is Euclidean distance, is used as a standard parameter for comparing the different isotropic functions c(r). Diffusion on  $\mathbb{S}^2$  is shown to be well approximated by diffusion on  $\mathbb{R}^2$  for length-scales of interest. As a result, fundamental parameters that define the correlation model on  $\mathbb{S}^2$  can be specified using more convenient expressions available on  $\mathbb{R}^2$ .

Anisotropic Gaussian or Matérn correlation functions on  $\mathbb{R}^d$  can be represented by a diffusion operator furnished with a symmetric and positive-definite diffusion tensor. For anisotropic functions  $c(\mathbf{r})$ , the tensor  $D = -(\nabla \nabla^{\mathrm{T}} c(\mathbf{r})|_{\mathbf{r}=\mathbf{0}})^{-1}$  where  $\nabla$  is the *d*-dimensional gradient operator, is a natural generalization of the (square of) the Daley length-scale for characterizing the spatial scales of the function. Relationships between this tensor, which we call the Daley tensor, and the diffusion tensor of the explicit and implicit diffusion operators are established. Methods to estimate the elements of the local Daley tensor from a sample of simulated background errors are presented and compared in an idealized experiment with spatially varying covariance parameters. Since the number of independent parameters needed to specify the local diffusion tensor is of the order of the total number of grid points N, sampling errors are inherently much smaller than those involved in the order  $N^2$ estimation problem of the full correlation function. While the correlation models presented in this paper are general, the discussion is slanted to their application to background-error correlation modelling in ocean data assimilation. Copyright (c) 2012 Royal Meteorological Society

*Key Words:* correlation functions; covariance modelling; background error; ocean data assimilation; diffusion tensor; ensemble estimation

Received 30 June 2011; Revised 13 March 2012; Accepted 20 March 2012; Published online in Wiley Online Library 17 May 2012

*Citation:* Weaver AT, Mirouze I. 2013. On the diffusion equation and its application to isotropic and anisotropic correlation modelling in variational assimilation. *Q. J. R. Meteorol. Soc.* **139**: 242–260. DOI:10.1002/qj.1955

#### 1. Introduction

Various methods have been proposed for modelling background-error correlations in geophysical applications of variational data assimilation (VDA) (see Bannister, 2008, for example, for a thorough review of methods used in atmospheric VDA). In ocean VDA, backgrounderror correlation models based on the diffusion equation are popular. The method has its origins in the work of Derber and Rosati (1989), who proposed the use of an iterative Laplacian grid-point filter in order to approximate a Gaussian correlation operator. Egbert et al. (1994) described a close variant of the algorithm in which the Laplacian grid-point filter could be interpreted as a pseudo-timestep integration of a diffusion equation with an explicit scheme. Weaver and Courtier (2001) (hereafter WC01) described the algorithm in more detail and proposed various extensions to account for more general correlation functions than the quasi-Gaussian of the original Derber and Rosati (1989) algorithm. Correlation models based on explicit diffusion methods have been used in various VDA systems in oceanography (Weaver et al., 2003; Di Lorenzo et al., 2007; Muccino et al., 2008; Daget et al., 2009; Kurapov et al., 2009; Moore et al., 2011), meteorology (Bennett et al. 1996), and atmospheric chemistry (Geer et al., 2006; Elbern et al., 2010).

An explicit diffusion scheme is appealing because of its simplicity, but can be expensive if many iterations are required to keep the scheme numerically stable. This can occur when the local diffusion scale is 'large' relative to the local grid size. To keep the explicit scheme affordable, the correlation length-scales must be bounded even if statistics or physical considerations suggest that larger values would be more appropriate. This limitation can be overcome by reformulating the diffusion model using an implicit scheme which has the advantage of being unconditionally stable.

One-dimensional (1D) implicit diffusion operators have been used for representing temporal and vertical correlation functions (Bennett et al., 1997; Chua and Bennett, 2001; Ngodock, 2005) and products of 1D implicit diffusion operators have been used for constructing two-dimensional (2D) and three-dimensional (3D) correlation models (Chua and Bennett, 2001; Zaron et al., 2009). The correlation kernels associated with the 1D implicit diffusion operator belong to the family of Mth-order autoregressive (AR) functions where M is the number of implicit iterations (Mirouze and Weaver, 2010; hereafter MW10). As discussed by MW10, the 1D implicit diffusion operator is closely linked to the recursive filter (Lorenc, 1992; Hayden and Purser, 1995), which has been developed extensively in meteorology for constructing correlation models in multiple dimensions (Wu et al., 2002; Purser et al., 2003a, 2003b; Liu et al., 2007). The recursive filter has also been employed in ocean data assimilation systems (Martin et al., 2007; Dobricic and Pinardi, 2008; Liu et al., 2009).

The 1D implicit diffusion approach for constructing 2D and 3D correlation models can be convenient for computational reasons, but has limitations. For example, with few iterations, the product of 1D implicit diffusion operators produces a well-known spurious anisotropic response (Purser *et al.*, 2003a). Unphysical features can also appear near complex boundaries, such as coastlines or islands in an ocean model, where correlation functions cannot always be reasonably represented by a product of

separable functions of the model's coordinates. Correlation models based on 2D or 3D implicit diffusion operators can overcome these limitations but are more difficult to implement since they involve the solution of a large linear system (matrices of dimension  $O(10^6 \times 10^6)$  or larger in VDA). Some progress in the development of this approach has been made by Weaver and Ricci (2004) and Massart *et al.*(2012), who used sparse matrix methods to solve a 2D implicit diffusion problem directly, and by Carrier and Ngodock (2010) and S. Gratton (2011, personal communication), who used iterative methods based on conjugate gradient or multi-grid to approximate the solution of a 2D or 3D implicit diffusion problem.

Multidimensional implicit diffusion correlation operators can be interpreted in terms of smoothing norm splines, which were introduced to atmospheric data assimilation by Wahba and Wendelberger (1982) and Wahba (1982), and discussed within an oceanographic context by McIntosh (1990). In the norm spline approach, the background term of the cost function in VDA is formulated in terms of a linear combination of weighted derivative operators that penalize explicitly the amplitude and curvature of the solution. When the weighting coefficients are given by binomial coefficients, the inverse of the background-error correlation operator implied by the norm spline can be expressed as the inverse of an implicit diffusion operator. The direct penalty approach was popular in some of the early studies of four-dimensional VDA (Thacker, 1988; Sheinbaum and Anderson, 1990) but generally leads to a poorly conditioned minimization problem (Lorenc et al., 2000). Effective preconditioning techniques for VDA require access to the backgrounderror covariance operator itself. An interesting exception is the recent study of Yaremchuk et al. (2011), who propose a variational formulation in which the inverse of the background-error covariance is modelled directly using the inverse of a low-order (two-iteration) 3D implicit diffusion operator. No apparent conditioning problems were reported in their examples from an ocean VDA system.

The present paper has a dual purpose: first, to provide a review of the diffusion equation as a basis for constructing anisotropic and inhomogeneous correlation models for data assimilation; and second, to illustrate how fundamental parameters that control spatial smoothness properties of these models can be estimated using ensemble methods. Section 2 brings together key results from data assimilation and geostatistics on the isotropic diffusion problem. Diffusion is considered both on the sphere and in the d-dimensional Euclidean space. Analytical expressions for the isotropic correlation functions implied by appropriately normalized explicit and implicit diffusion in these spaces are presented and compared. The Daley length-scale is used as a standard parameter for comparing the different functions, and expressions relating it to the parameters of the diffusion-model are established.

The results from section 2 provide the foundation for building anisotropic correlation models with the diffusion equation. This is discussed in sections 3 and 4. The Daley tensor is introduced, which is defined as the negative inverse of the tensor of second derivatives of the correlation function evaluated at zero distance (the Hessian tensor). The Daley tensor is an anisotropic generalization of the Daley length-scale. Expressions relating the Daley tensor to the diffusion tensor of the diffusion models are given. Section 4 discusses techniques for estimating the Daley tensor from statistics of a sample of simulated errors such as those that would be available from an ensemble data assimilation system. Idealized experiments are then presented to compare the effectiveness of two of the estimation techniques. Conclusions are given in section 5. Appendix A provides a derivation of the relationship between the Daley and diffusion tensors for the correlation functions represented by the implicit diffusion equation in  $\mathbb{R}^2$ . Appendix B provides a derivation of the key formulae for estimating the Daley tensor.

#### 2. Isotropic diffusion

Coordinate systems of global atmospheric and ocean models refer to the spherical-shell geometry of the atmosphere and ocean. From a mathematical perspective, this leads naturally to consideration of 2D 'horizontal' correlation functions on the spherical space  $\mathbb{S}^2$ . The product of a 2D correlation function on  $\mathbb{S}^2$  and a 1D correlation function on the bounded subset of the Euclidean space  $\mathbb{R}^1$  is commonly used to construct 3D correlation functions on the sphericalshell subspace of  $\mathbb{R}^3$  that defines the model domain. This approach of separating the horizontal and vertical correlation functions is usually justified by the fact that the global atmospheric and ocean circulations are characterized by scales that are much larger in the horizontal direction (along geopotential surfaces) than in the vertical direction (perpendicular to geopotential surfaces). In the remainder of this section, the correlation functions that can be represented by isotropic diffusion on  $\mathbb{S}^2$  and the general Euclidean space  $\mathbb{R}^d$  are described. Table 1 provides a brief description of the main symbols used in this section.

## 2.1. Explicit diffusion on $\mathbb{S}^2$

Consider the 2D diffusion equation applied to the scalar field  $\eta(\lambda, \phi, s)$ :

$$\frac{\partial \eta}{\partial s} - \kappa \nabla^2 \eta = 0, \tag{1}$$

where  $\kappa > 0$  is a diffusion coefficient, and

$$\nabla^2 = \frac{1}{a^2 \cos \phi} \frac{\partial}{\partial \phi} \left( \cos \phi \frac{\partial}{\partial \phi} \right) + \frac{1}{a^2 \cos^2 \phi} \frac{\partial^2}{\partial \lambda^2}$$

is the Laplacian operator in geographical coordinates  $(\lambda, \phi)$ ,  $\lambda$  denoting longitude  $(0 \le \lambda \le 2\pi)$ ,  $\phi$  latitude  $(-\pi/2 \le \phi \le \pi/2)$ , and *a* the radius of the sphere (the Earth's radius in our case). In the context of this paper, *s* is to be interpreted as a dimensionless pseudo-time coordinate. The diffusion coefficient then has physical units of length squared. The solution of Eq. (1) on S<sup>2</sup> can be interpreted as a covariance operator (e.g. see WC01). Let

$$\eta(\lambda,\phi,0) = \gamma^{s} \,\widetilde{\eta}(\lambda,\phi) \tag{2}$$

denote the initial condition, where  $\gamma^s$  is a normalization constant. The solution at some s > 0 can be expressed as the integral operator  $C^s : \tilde{\eta} \mapsto \eta(s)$ ,

$$\eta(\lambda,\phi,s) = \int_{\mathbb{S}^2} c^{s}(\theta) \,\widetilde{\eta}(\lambda',\phi') \, a^2 \cos \phi' \, \mathrm{d}\lambda' \, \mathrm{d}\phi', \quad (3)$$

Copyright © 2012 Royal Meteorological Society

Table 1. A list of the main generic symbols used in section 2. The specification of the superscripts  $\alpha$  and  $\beta$  is summarized in the bottom table. A quantity in  $\mathbb{R}^d$  is supplemented with a subscript *d* if it depends explicitly on the dimension of the space; otherwise it is omitted.

Symbol		Description
$ \begin{array}{l} \overline{\mathcal{C}^{\alpha}, \mathcal{C}^{\beta}} \\ \overline{\mathcal{C}^{\alpha}(\mathbf{x}, \mathbf{x}')} \\ \mathbf{x} \\ \overline{\mathbf{x}} \\ \overline{\mathbf{c}^{\alpha}_{d}(\mathbf{r})} \\ r \\ \overline{\mathbf{c}^{\alpha}_{d}(\mathbf{x})} \\ \mathbf{\hat{x}} \\ \overline{\mathbf{c}^{\beta}(\theta)} \\ \theta \\ \overline{\mathbf{c}^{\beta}_{n}} \\ n \\ D^{\alpha}, D^{\beta} \\ \overline{\mathbf{v}^{\beta}_{d}}, \overline{\mathbf{v}^{\beta}} \end{array} $		Correlation operators on $\mathbb{R}^d$ and $\mathbb{S}^2$ General correlation function on $\mathbb{R}^d$ Vector of Cartesian coordinates Isotropic correlation function on $\mathbb{R}^d$ Euclidean distance Fourier transform of $c_d^{\alpha}(r)$ Vector of spectral wave numbers Isotropic correlation function on $\mathbb{S}^2$ Angular separation Legendre coefficients for $c^{\beta}(\theta)$ Total wave number Daley length-scale of $c_d^{\alpha}(r)$ and $c^{\beta}(\theta)$ Normalization constants on $\mathbb{R}^d$ and $\mathbb{S}^2$
Superscript		Description
α	g	Regular diffusion on $\mathbb{R}^d$
β	w s	Implicit diffusion on $\mathbb{R}^{4}$ Regular diffusion on $\mathbb{S}^{2}$
r	h	Implicit diffusion on $\mathbb{S}^2$

where  $c^{s}(\theta)$  is an *isotropic* function that depends on the angular separation  $\theta$ ,  $0 \le \theta \le \pi$ , between points  $(\lambda, \phi)$  and  $(\lambda', \phi')$  on the sphere:

$$\cos\theta = \cos\phi\,\cos\phi'\,\cos(\lambda - \lambda') + \sin\phi\,\sin\phi'.$$
 (4)

The normalization constant  $\gamma^s$  in Eq. (2) has been absorbed into the function  $c^s(\theta)$  which has the specific form

$$c^{s}(\theta) = \sum_{n=0}^{\infty} c_{n}^{s} P_{n}^{0}(\cos\theta), \qquad (5)$$

where

$$c_n^{\rm s} = \frac{\gamma^{\rm s}}{4\pi a^2} \sqrt{2n+1} \exp\left(-\frac{\kappa s}{a^2} n(n+1)\right),$$
 (6)

*n* being the total wave number, and  $P_n^0(\cos \theta)$  the Legendre polynomials, normalized such that  $P_n^0(1) = \sqrt{2n+1}$ , following the usual convention in meteorology (Courtier *et al.*, 1998). All isotropic covariance functions on  $\mathbb{S}^2$  can be expressed, as in Eq. (5), as an expansion in terms of the Legendre polynomials (Weber and Talkner, 1993; Theorem 2.11 of Gaspari and Cohn, 1999). They are positive-definite functions on  $\mathbb{S}^2$  if the spectral coefficients are positive, which is clearly the case for all of the coefficients  $c_n^s$ . Equation (3) is thus a valid covariance operator on  $\mathbb{S}^2$ .

The covariance function is readily transformed into a correlation function ( $c^{s}(0) = 1$ ) by defining the normalization constant as

$$\gamma^{s} = 4\pi a^{2} \left( \sum_{n=0}^{\infty} (2n+1) \exp\left(-\frac{\kappa s}{a^{2}} n(n+1)\right) \right)^{-1}.$$
 (7)

The fundamental parameter controlling the shape of the correlation function is the *product*  $\kappa s$  in Eq. (6). To define

the length-scale of  $c^{s}(\theta)$ , we use the standard definition from Daley (1991, p. 110), the geometrical interpretation of which is discussed by Pannekoucke *et al.* (2008). For  $c^{s}(\theta)$ , the *Daley length-scale* reads

$$D^{s} = \sqrt{-\frac{2}{\nabla^{2} c^{s}|_{\theta=0}}}$$
  
=  $a \sqrt{\frac{2}{\sum_{n=0}^{\infty} n(n+1)\sqrt{2n+1} c_{n}^{s}}}.$  (8)

Equations (6) and (8) provide a relationship between  $\kappa s$  and  $D^s$  that allows us to control the correlation shape (length-scale).

Now consider a discretized version of Eq. (1) in which the first-order derivative is approximated using a forward-Euler (*explicit*) scheme. This yields

$$\eta(\lambda,\phi,s_m) = \eta(\lambda,\phi,s_{m-1}) + \kappa \,\Delta s \nabla^2 \eta(\lambda,\phi,s_{m-1}), \quad (9)$$

where *m* is a positive integer,  $\Delta s = s_m - s_{m-1}$  is the step size, and  $\nabla^2$  is understood to be the Laplacian operator in discretized form. For convenience, we can assume that  $s_m = m$  so that the step size  $\Delta s = 1$ . This parameter can thus be ignored hereafter without loss of generality. Repeated applications of Eq. (9) on the interval  $0 < m \le M$  leads to the linear operator

$$\eta(\lambda,\phi,M) = \left(1 + \kappa \nabla^2\right)^M \eta(\lambda,\phi,0), \tag{10}$$

where  $\eta(\lambda, \phi, 0)$  is given by Eq. (2). For clarity we let

$$\kappa = L^2, \tag{11}$$

to emphasize that the coefficient is positive and can be interpreted as the square of a scale parameter.

The key idea is that, on a numerical grid, the effect of the integral correlation operator (3) on an arbitrary scalar field  $\tilde{\eta}(\lambda, \phi)$  can be approximated by applying a discretized differential operator (10). This is the essence of the original Derber and Rosati (1989) scheme. The parameter  $\kappa s$  of the correlation function  $c^{s}(\theta)$  can be related to the parameters *M* and  $\kappa$  of Eq. (10) by noticing that  $\kappa s = \kappa M = ML^2$ . In practice, it is customary to prescribe the Daley lengthscale  $(D^{s})$ . Given  $D^{s}$ , the product  $ML^{2}$  can be determined by a non-trivial inversion of Eq. (8). This has been done by trial and error for the illustrative examples presented in this paper. To determine M and  $L^2$  from the product  $ML^2$ , we have an additional requirement that M must be sufficiently large ( $L^2$  sufficiently small) in order to maintain the numerical stability of the explicit scheme. Provided M is not too 'large', applying the discretized operator (10) is an efficient way of evaluating the integral operator (3). What defines an acceptable value of M will depend on the application.

To represent a larger family of correlation functions than Eqs (5) and (6), WC01 proposed a *generalized* diffusion model in which the scaled Laplacian in Eq. (1) is replaced by a linear combination of powers of scaled Laplacians:

$$-\kappa \nabla^2 \mapsto \sum_{p=1}^p \kappa_p (-\nabla^2)^p,$$
 (12)

Copyright © 2012 Royal Meteorological Society

where the diffusion coefficients  $\kappa_p > 0$  can be related to a general set of scale parameters  $L_p$  via the equation

$$\kappa_p = L_p^{2p}, \qquad p = 1, \dots, P. \tag{13}$$

The resulting correlation functions have the same basic form as Eq. (5) but with the  $c_n^s$  given by

$$c_n^s = \frac{\gamma^s}{4\pi a^2} \sqrt{2n+1} \exp\left(-\sum_{p=1}^p \frac{\kappa_p s}{a^{2p}} \left(n(n+1)\right)^p\right), \quad (14)$$

and the appropriate modification to  $\gamma^s$  to produce a unitamplitude function. Equation (6) is a special case of Eq. (14) with P = 1. Unlike the standard diffusion model, the generalized diffusion model can be used to represent correlation functions that change sign, as illustrated in Figure 1 of WC01. This is an appealing feature if there is compelling evidence of negative correlations in the error fields, although representing them with powers of Laplacian operators would clearly increase the cost of the correlation model.

#### 2.2. Explicit diffusion on $\mathbb{R}^d$

Now consider the diffusion equation (1) on the *d*-dimensional Euclidean space  $\mathbb{R}^d$ , where  $\nabla^2$  now represents the Laplacian operator in Cartesian coordinates  $\mathbf{x} = (x_1, \ldots, x_d)$ . While our particular interest concerns the spaces  $\mathbb{R}^1$ ,  $\mathbb{R}^2$  and  $\mathbb{R}^3$ , it is easier to consider them as special cases of the general diffusion problem in  $\mathbb{R}^d$ . The initial condition of the diffusion problem can be written as

$$\eta(\mathbf{x}, 0) = \gamma_d^{\mathsf{g}} \, \widetilde{\eta}(\mathbf{x}), \tag{15}$$

where  $\gamma_d^g > 0$  is a normalization constant and  $\tilde{\eta}(\mathbf{x})$  is assumed to be bounded at infinity. Using the Fourier transform (FT), the solution at 'time' s > 0 can be written as a convolution operator  $C_d^g : \tilde{\eta} \mapsto \eta(s)$ :

$$\eta(\mathbf{x}, s) = \int_{\mathbb{R}^d} C(\mathbf{x}, \mathbf{x}') \,\widetilde{\eta}(\mathbf{x}') \, \mathrm{d}\mathbf{x}', \qquad (16)$$

where  $C(\mathbf{x}, \mathbf{x}') = c^{g}(r)$  is the Gaussian function

$$c^{\mathrm{g}}(r) = \frac{\gamma_d^{\mathrm{g}}}{(4\pi\kappa s)^{d/2}} e^{-r^2/4\kappa s},$$
(17)

 $r = |\mathbf{x} - \mathbf{x}'|$  being the Euclidean distance between points  $\mathbf{x}$  and  $\mathbf{x}'$  on  $\mathbb{R}^d$ . Setting the normalization factor to

$$\gamma_d^{\rm g} = (4\pi\kappa s)^{d/2} \tag{18}$$

ensures that  $c^{g}(0) = 1$ .

The Daley length-scale for any twice differentiable, isotropic correlation function c(r) in d dimensions is given by

$$D = \sqrt{\frac{d}{\operatorname{tr}(-\nabla\nabla^{\mathrm{T}}c|_{r=0})}} = \sqrt{-\frac{d}{\nabla^{2}c|_{r=0}}},\qquad(19)$$

where  $\nabla \nabla^{T}$  is the outer product of the *d*-dimensional gradient operator  $\nabla = (\partial/\partial x_1 \dots \partial/\partial x_d)^{T}$  and its transpose. The



Figure 1. The grid-point power spectra  $\sqrt{2n+1}c_n^{\rm h}$  (lower panel) of sample correlation functions generated with Eqs (5) and (30) using different values of M. The scale parameters have been set to L = 353 km, L = 250 km and L = 125 km for the functions corresponding to M = 3 (dashed-dotted curves), M = 4(dashed curves) and M = 10 (dotted curves), respectively, in order to achieve a common Daley length-scale of  $D^{h} = 500 \text{ km}$  (see Eq. (8)). The Gaussian correlation function  $c^{g}(r(\theta))$  (Eq. (24)) with  $D^{g} = 500 \text{ km}$  is shown for reference (thick solid curves). The correlation functions in the upper panel are plotted as a function of chordal distance (Eq. (23)). A spectral truncation at n = 500 has been used. The lower panel is plotted on a log-log scale.

quantity within the trace operator is the correlation Hessian tensor (Chorti and Hristopulos, 2008). The Hessian tensor plays a fundamental role in characterizing the anisotropic correlation functions described later in this paper (sections 3 and 4). For the *d*-dimensional Gaussian function, the Daley length-scale is

$$D^{g} \sqrt{2\kappa s}$$
. (20)

In terms of  $D^{g}$ , the normalization factor is

a

$$\gamma_d^{g} (2\pi)^{d/2} (D^{g})^d.$$
 (21)

As before, we can approximate Eq. (16) with a differential operator based on a discretization of the diffusion equation using an M-step explicit scheme. In terms of the parameters M and  $L^2$  of the explicit diffusion operator, Eqs (20) and (21) become

$$D^{g} = \sqrt{2ML}$$
(22)  
nd  $\gamma_{J}^{g} = (4M\pi)^{d/2}L^{d}.$ 

an  

$$10^1$$
  $10^2$   
Wave number (n)  
values  $c^h(r(\theta))$  (upper panel) and the variance-  
the (lower panel) of cample correlation functions

Let us consider now the interpretation of the Gaussian function on  $\mathbb{S}^2$ . First, since  $\mathbb{S}^2$  is embedded in  $\mathbb{R}^3$ , a valid isotropic correlation function on  $\mathbb{S}^2$  can always be constructed from a valid isotropic correlation function in  $\mathbb{R}^3$  by restricting  $\mathbf{x} = (x_1, x_2, x_3)$  and  $\mathbf{x}' = (x_1', x_2', x_3')$  to be points on the sphere. Expressing these points in geographical coordinates  $\mathbf{x} = (a \cos \phi \cos \lambda, a \cos \phi \cos \lambda, a \sin \phi)$  and  $\mathbf{x}' = (a \cos \phi' \cos \lambda', a \cos \phi' \sin \lambda', a \sin \phi')$  leads to the chordal distance measure

$$r = r(\theta) = a\sqrt{2(1 - \cos\theta)}, \quad 0 \le \theta \le \pi,$$
 (23)

where  $\cos \theta$  is given by Eq. (4). The Gaussian correlation function on  $\mathbb{R}^3$  confined to the subspace  $\mathbb{S}^2$  is thus

$$c^{g}(r(\theta)) = e^{-(r(\theta))^{2}/2(D^{g})^{2}} = e^{-a^{2}(1-\cos\theta)/(D^{g})^{2}}$$

From Eq. (23) we notice that *r* depends only on  $\cos \theta$ , or alternatively  $\theta$ , and that  $\cos \theta = 1 - r^2/2a^2$ , where  $0 \le r \le 2a$ . We also recall that all isotropic correlation functions on  $\mathbb{S}^2$  can be expressed as a Legendre expansion that depends only on  $\cos\theta$  (Eq. (5)). It is then possible to represent any isotropic correlation function on  $\mathbb{S}^2$  as a function of either *r* or  $\theta$ .\*

In particular, consider the representation of the Gaussian on  $\mathbb{S}^2$  in terms of the Legendre polynomials. As shown in WC01:

$$c^{g}(\theta) = \sum_{n=0}^{\infty} c_n^{g} P_n^0(\cos\theta)', \qquad (24)$$

ere

a

$$c_n^{g} = \widetilde{\gamma}^{g} \sqrt{2n+1} \frac{I_{n+1/2}(\omega)}{I_{1/2}(\omega)}$$
  
nd  $\widetilde{\gamma}^{g} = \frac{e^{-\omega}\sinh(\omega)}{\omega},$ 

 $I_{n+1/2}(\omega)$  denoting the modified Bessel function of fractional order n + 1/2, and  $\omega = (a/D^g)^2$ . In view of the results on  $\mathbb{R}^d$ , one might expect that the correlation kernel  $c^{s}(\theta)$  implied by diffusion on  $\mathbb{S}^2$  (Eq. (5)) is similar to the Gaussian correlation function (24) on  $\mathbb{S}^2$ . Indeed, for a given length-scale  $D^g$ , it is possible to find a corresponding parameter  $\kappa s$  in Eq. (6) such that the difference between  $c^{s}(\theta)$  and  $c^{g}(\theta)$  is 'small' (Roberts and Ursell, 1960; Hartman and Watson, 1974). In particular, consider the scales of interest in atmospheric and ocean data assimilation for which  $\omega \gg 1$ . Matching the n = 0 coefficients  $c_0^s$  and  $c_0^g$  of the Legendre polynomials and noting that

$$\widetilde{\gamma}^{\mathrm{g}} \approx \frac{(D^{\mathrm{g}})^2}{2a^2}$$

for large  $\omega$ , we obtain the approximation to the normalization factor

$$\gamma^{\rm s} \approx 2\pi (D^{\rm g})^2. \tag{25}$$

Q. J. R. Meteorol. Soc. 139: 242-260 (2013)

Copyright © 2012 Royal Meteorological Society

<sup>\*</sup>Isotropic correlation functions on  $\mathbb{S}^2$  will be written explicitly as a function of  $r(\theta)$  whenever the context requires an interpretation in terms of chordal distance.

Now matching the n = 1 coefficients  $c_1^s$  and  $c_1^g$  and using (25) leads to the approximation

$$\kappa s \approx \frac{(D^{g})^2}{2}.$$
 (26)

WC01 illustrate the excellent agreement between  $c^{s}(\theta)$  and  $c^{g}(\theta)$ , particularly for the large scales, for a given length-scale  $D^{g}$  and with  $\kappa s$  approximated according to Eq. (26) (see their Figure A1).

Equations (25) and (26) are none other than those derived earlier for the diffusion problem in  $\mathbb{R}^2$  (cf. Eqs (21) and (20) with d = 2). In other words, for length-scales small compared to the radius of the Earth, we obtain the somewhat intuitive result that diffusion on the sphere (the  $C^s$  operator) is well approximated by diffusion on the 2D Cartesian plane (the  $C_2^g$  operator). For calibrating the correlation model, it is then possible to employ the simple expressions (26) and (25) for the length-scale and normalization factor in place of the more complicated expressions (8) and (7).

There are two main drawbacks with the generalized explicit diffusion model of WC01. First, the correlation functions that can be represented by the model have limited flexibility in the spectral domain, especially at high wave numbers where their decay rates are at least as fast as that of the Gaussian function. In data assimilation, this can result in excessive smoothing of small-scale features in the analysis (Purser et al., 2003b). Second, the explicit scheme is subject to a stability criterion that depends on the ratio of the length-scale and grid size, raised to the power of 2P. As a result, many iterations may be required when the lengthscale is large compared with the grid resolution. With the variable coefficient and anisotropic versions of the model discussed later, the computational cost of the algorithm can be especially high. A diffusion model based on an implicit formulation can overcome these limitations, as described next.

### 2.3. Implicit diffusion on $\mathbb{S}^2$

Consider again the diffusion equation (1) but this time discretized using a backward-Euler (*implicit*) scheme:

$$\eta(\lambda,\phi,s_m) = \eta(\lambda,\phi,s_{m-1}) + \kappa \,\Delta s \nabla^2 \eta(\lambda,\phi,s_m), \quad (27)$$

where, as in Eq. (10), we can assume  $s_m = m$  and hence  $\Delta s = 1$ , and interpret  $\kappa$  as the square of a scale parameter (Eq. (11)). Rearranging Eq. (27) and applying it repeatedly on the interval  $0 < m \le M$  leads to the 'reverse-time' or inverse diffusion operator

$$\left(1 - \kappa \nabla^2\right)^M \eta(\lambda, \phi, M) = \eta(\lambda, \phi, 0).$$
(28)

Equation (28) can be interpreted as a *roughening* operator as opposed to the diffusion operator itself, which is a *smoothing* operator.

Following Eq. (2), we define the initial condition as

$$\eta(\lambda,\phi,0) = \gamma^{h} \widetilde{\eta}(\lambda,\phi) \tag{29}$$

where  $\gamma^{h}$  is a normalization constant. Weaver and Ricci (2004) show that the differential operator  $(\mathcal{C}^{h})^{-1}: \eta(M) \mapsto \widetilde{\eta}$  is the *inverse* of a correlation operator  $\mathcal{C}^{h}: \widetilde{\eta} \mapsto \eta(M)$ , where the latter is given by an integral equation of the form (3), with isotropic correlation function  $c^{h}(\theta)$  of the Legendre form (5) as its kernel. The spectral coefficients of  $c^{h}(\theta)$  are strictly positive and given by

$$c_n^{\rm h} = \frac{\gamma^{\rm h}}{4\pi a^2} \sqrt{2n+1} \left(1 + \frac{L^2}{a^2}n(n+1)\right)^{-M}.$$
 (30)

The normalization factor is

$$\gamma^{\rm h} = 4\pi a^2 \left( \sum_{n=0}^{\infty} (2n+1) \left( 1 + \frac{L^2}{a^2} n(n+1) \right)^{-M} \right)^{-1} (31)$$

and the Daley length-scale is given by Eq. (8) with  $c_n^s$  replaced by  $c_n^h$ .

In the explicit diffusion model, the only free parameter was the product  $\kappa s_M = ML^2$  which controls the spatial scale of the quasi-Gaussian correlation kernel (Eq. (6) with  $s = s_M$ ). The implicit diffusion model, on the other hand, allows for greater control of the shape characteristics of the associated correlation kernels since both  $L^2$  and M are free parameters. Numerically, this extra flexibility is reflected by the important property of *unconditional stability* of the implicit scheme. In the limiting case of  $M \rightarrow \infty$ , with  $ML^2$ held fixed, the spectral coefficients (30) reduce to those of the quasi-Gaussian solution which is the only correlation function that can be represented by solving the diffusion equation explicitly.

The upper panel in Figure 1 displays correlation functions  $c^{h}(r(\theta))$  for different values of M and a constant Daley length-scale (500 km). The values are plotted as a function of chordal distance  $r(\theta)$ . The Gaussian function  $c^{g}(r(\theta))$  is also shown for reference. Increasing the value of M decreases the 'fatness' of the tail of  $c^{h}(r(\theta))$ , with the Gaussian providing the upper limit as  $M \to \infty$ . The total variance of  $c^{h}(r(\theta))$  and  $c^{g}(r(\theta))$  is given by their value at the origin, which is equal to one. The coefficients  $\sqrt{2n+1} c_n^h$  and  $\sqrt{2n+1} c_n^g$  give the contribution of each wave number n to the total variance of  $c^{h}(r(\theta))$  and  $c^{g}(r(\theta))$ , respectively, and thus define the variance-power spectra. The lower panel in Figure 1 shows a log-log plot of this spectra as a function of n. Here we see that the increased fatness in correlation shape for low values of M is associated with higher variance and a reduced damping rate in the small scales, slightly less variance in the intermediate scales, and increased variance in the large scales.

As with the generalized diffusion equation, a linear combination of powers of scaled Laplacian operators (12) can be introduced in Eq. (28) to yield a larger family of correlation functions, but at extra cost. The spectral coefficients of this larger family are given by

$$c_n^{\rm h} = \frac{\gamma^{\rm h}}{4\pi a^2} \sqrt{2n+1} \left( 1 + \sum_{p=1}^{P} \left( \frac{L_p^2}{a^2} \right)^p (n(n+1))^p \right)^{-M},$$
(32)

with  $\gamma^{h}$  modified accordingly so that  $c^{h}(0) = 1$ . The smoothing spline functions introduced by Wahba (1982) correspond to the special case of Eqs (5) and (32) for which M = 2.

Increasing the degree P of the polynomial of the Laplacian leads to correlation functions that oscillate about the zero axis. This is illustrated in the upper panel of Figure 2,

Copyright © 2012 Royal Meteorological Society

**Figure 2.** As Figure 1 but for sample correlation functions generated with Eqs (5) and (32) with a fixed value of M = 4 and different values of P. The scale parameters have been adjusted to yield a common Daley length-scale of 500 km:  $L_1 = 250$  km with P = 1 (dashed-dotted curves),  $L_1 = 0$  and  $L_2 = 206$  km with P = 2 (dashed curves),  $L_1 = L_2 = 0$ ,  $L_3 = 209$  km with P = 3 (dotted curves).

where the generalized  $c^{h}(r(\theta))$  are displayed with different values of *P* but a fixed value of M = 4. The amplitude of the negative lobes increases with increasing value of *P*. In spectral space, the negative lobes are associated with a decrease in variance in the large scales and an increase in variance in the intermediate scales. Increasing the value *P* also leads to a steepening of the decay rate of the variance in the smaller scales.

A straightforward variant of Eq. (32) that can be used to enhance the oscillations while maintaining a gradual spectral decay rate at high wave numbers is

$$c_n^{\rm h} = \frac{\gamma^{\rm h}}{4\pi a^2} \sqrt{2n+1} \left( 1 + \sum_{p=1}^{P} \rho_p \left( \frac{L_p^2}{a^2} \right)^p (n(n+1))^p \right)^{-M},$$
(33)

where  $\rho_p$  is a dimensionless coefficient that can take on both negative and positive values. This is equivalent to redefining the diffusion coefficients (13) as  $\kappa_p = \rho_p L^{2p}$ . Equation (33) yields positive coefficients by restricting *M* to be even. Examples are shown in Figure 3 for the case P = 2 and M = 2, and a single scale parameter  $L_1 = L_2 = L$ . Here



**Figure 3.** As Figure 1 but for sample correlation functions generated with Eqs (5) and (33) with a fixed value of M = 2, P = 2 and  $\rho_2 = 1$ , and different values of  $\rho_1$ . A single scale parameter  $L_1 = L_2 = L$  has been used and adjusted to yield a common Daley length-scale of 500 km: L = 308 km with  $\rho_1 = -1$  (dashed-dotted curves), L = 326 km with  $\rho_1 = -1.5$  (dashed curves), and L = 340 km with  $\rho_1 = -1.8$  (dotted curves).

 $\rho_2$  has been set to one and negative values have been used for  $\rho_1$ . Increasing the magnitude of  $\rho_1$  results in a significant increase in the amplitude of the oscillations and a much sharper spectral peak at intermediate scales. Notice that by setting  $\rho_1 = 2$  we recover the non-oscillatory correlation function governed by Eq. (30) with M = 4, which is displayed in Figure 1 (dashed curves).

On a numerical grid,  $C^h$  can be approximated by a discrete operator that solves the linear system (28)–(29) for a given right-hand side  $\tilde{\eta}(\lambda, \phi)$ . We refer to  $C^h$  as an *implicit* diffusion correlation operator. Although the cost of each iteration of an implicit diffusion operator will generally increase relative to that of the explicit scheme, the total cost of the implicit algorithm can easily decrease through the possibility of performing significantly fewer iterations.

### 2.4. Implicit diffusion on $\mathbb{R}^d$

The starting point is the following general fractional differential operator  $(\mathcal{C}_d^w)^{-1}: \psi \mapsto \widetilde{\psi}$  (Whittle, 1954, 1963; Guttorp and Gneiting, 2006):

$$(\gamma_d^{\mathrm{w}})^{-1} \left(1 - L^2 \nabla^2\right)^{\nu + d/2} \psi(\mathbf{x}) = \widetilde{\psi}(\mathbf{x}), \qquad (34)$$





where  $\widetilde{\psi}(\mathbf{x}) \in \mathbb{R}^d$  is assumed to be bounded at infinity,  $\nu > 0$  is a smoothness parameter, and  $\gamma_d^{w} > 0$  is a normalization constant that depends on the dimension of the space. The FT of Eq. (34) gives the relation

$$\widehat{\psi}(\widehat{\mathbf{x}}) = \widehat{c}_d^{\mathsf{w}}(\widehat{r})\,\widehat{\widehat{\psi}}(\widehat{\mathbf{x}}),\tag{35}$$

where  $\widehat{\psi}(\hat{\mathbf{x}})$  and  $\widetilde{\psi}(\hat{\mathbf{x}})$  denote the FTs of  $\psi(\mathbf{x})$  and  $\widetilde{\psi}(\mathbf{x})$ , respectively, and

$$\widehat{c}_{d}^{w}(\widehat{r}) = \frac{\gamma_{d}^{w}}{\left(1 + L^{2} \,\widehat{r}^{2}\right)^{\nu + d/2}},\tag{36}$$

 $\hat{r} = |\hat{\mathbf{x}}|$  being the norm of the vector of spectral wave numbers associated with  $\mathbf{x}$  (see also Yaglom, 1987, p. 363, Eq. (4.130); Stein, 1999, p. 49, Eq. (32); or Gneiting *et al.*, 2009, p. 16, Eq. (20)). Setting

$$\gamma_d^{w} = 2^d \pi^{d/2} \, \frac{\Gamma(\nu + d/2)}{\Gamma(\nu)} \, L^d, \tag{37}$$

where  $\Gamma(\nu)$  denotes the Gamma function, and applying the inverse FT to Eq. (35) leads to an integral solution of the general form (16), where  $C(\mathbf{x}, \mathbf{x}') = c^{w}(r)$  is a unit-amplitude isotropic function given by

$$c^{\mathsf{w}}(r) = \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{r}{L}\right)^{\nu} K_{\nu}\left(\frac{r}{L}\right), \qquad (38)$$

 $K_{\nu}(r/L)$  denoting the modified Bessel function of the second kind of order  $\nu$ , and  $r = |\mathbf{x} - \mathbf{x}'|$ . Since  $\hat{c}_d^{w}(\hat{r})$  is strictly positive,  $c^{w}(r)$  is a valid correlation function in  $\mathbb{R}^d$  (Bochner's theorem; see Theorem 2.10 in Gaspari and Cohn, 1999). Notice that the power spectrum  $\hat{c}_d^{w}(\hat{r})$  depends on d but the correlation function itself  $c^{w}(r)$  is independent of d.

Equation (38) is a class of correlation function well known in the geostatistical literature as the Whittle–Matérn or Matérn family (Gneiting, 1999; Stein, 1999; Guttorp and Gneiting, 2006). The link between this correlation family and the fractional differential operator (34) is attributed to Whittle (1954, 1963). Of particular interest here is the subclass of Matérn functions that correspond to (positive) *integer* values of the parameter M = v + d/2. For this subclass, the inverse correlation operator  $(C_d^w)^{-1}$  has a greatly simplified representation for numerical applications and can be interpreted as an *M*-step implicitly formulated diffusion operator (MW10), where (cf. Eqs (15) and (16))

$$\psi(\mathbf{x}) \mapsto \eta(\mathbf{x}, 0) = \gamma_d^{\mathsf{w}} \, \widetilde{\eta}(\mathbf{x}),$$
  
$$\psi(\mathbf{x}) \mapsto \eta(\mathbf{x}, M).$$

The correlation kernels and their associated FT are given by

$$\epsilon_d^{\rm w}(r) = \frac{2^{1-M+d/2}}{\Gamma(M-d/2)} \left(\frac{r}{L}\right)^{M-d/2} K_{M-d/2}\left(\frac{r}{L}\right), \quad (39)$$

and

$$\widehat{c}_d^{\mathsf{w}}(\widehat{r}) = \frac{\gamma_d^{\mathsf{w}}}{\left(1 + L^2 \, \widehat{r}^2\right)^M}$$

Equation (39) yields valid correlation functions if M > d - 1 ( $\nu > 0$ ; Guttorp and Gneiting, 2006). Notice

also that in contrast to the full Matérn family, the implicitdiffusion kernels depend on d (which has been made explicit by adding the subscript d in  $c_d^{w}(r)$ ) but their normalized power spectrum  $\hat{c}_d^{w}(\hat{r})/\hat{c}_d^{w}(0)$  is independent of d. For odd values of d, Eq. (39) reduces to a polynomial of order M - (d + 1)/2 times an exponential function; this is the well-known class of AR functions.

Of relevance here are the spaces  $\mathbb{R}^1$ ,  $\mathbb{R}^2$  and  $\mathbb{R}^3$ . The implicit-diffusion kernels on these spaces can be written explicitly as

$$_{1}^{w}(r) = \sum_{j=0}^{M-1} \beta_{j,M} \left(\frac{r}{L}\right)^{j} e^{-r/L},$$
 (40)

$$c_2^{\rm w}(r) = \frac{2^{2-M}}{(M-2)!} \left(\frac{r}{L}\right)^{M-1} K_{M-1}\left(\frac{r}{L}\right), \quad (41)$$

and 
$$c_3^{W}(r) = \sum_{j=0}^{M-2} \beta_{j,M-1} \left(\frac{r}{L}\right)^j e^{-r/L},$$
 (42)

where

С

$$\beta_{j,M} = \frac{2^{j}(M-1)! (2M-j-2)!}{j! (M-j-1)! (2M-2)!}$$

From Eq. (37), the expressions for the normalization constants become

$$\gamma_{1}^{w} = \frac{2^{2M-1}[(M-1)!]^{2}}{(2M-2)!}L,$$

$$\gamma_{2}^{w} = 4\pi (M-1)L^{2} \qquad (43)$$
and
$$\gamma_{3}^{w} = \frac{2^{2M-1}\pi [(M-2)!]^{2}(M-1)}{(2M-4)!}L^{3}.$$

Using Eq. (19), the Daley length-scale of the implicit diffusion kernels in  $\mathbb{R}^d$  can be evaluated as

$$D_d^{w} = \sqrt{2M - d - 2} L.$$
 (44)

Equation (44) is derived in MW10 for d = 1 and in Appendix A for d = 2. The generalization to d > 2 follows by noting that the correlation functions associated with odd d all have the form (40) with  $M \mapsto M - (d - 1)/2$ , while those functions with even d all have the form (41) with  $M \mapsto M - (d - 2)/2$ . Equation (44) imposes further restrictions on the choice of M where now we require

$$M > \begin{cases} (d+1)/2 & \text{if } d \text{ odd,} \\ (d+2)/2 & \text{if } d \text{ even.} \end{cases}$$

In  $\mathbb{R}^2$ , for example, we require M > 2. Finally, even values of M are more convenient than odd values of M since they greatly simplify the derivation of a 'square-root' factor of the diffusion operator, which is important for estimating normalization factors and for preconditioning in variational assimilation (WC01).

The explicit diffusion kernels are the limiting case of the implicit diffusion kernels as  $M \to \infty$  with  $D_d^w$  fixed. This is easily deduced from Eqs (36), (37) and (44) where, for  $M = \nu + d/2$  large,  $D_d^w \mapsto D^g$  (Eq. (22)),  $\gamma_d^w \mapsto \gamma_d^g$  (Eq. (21)), and  $\hat{c}_d^w(\hat{r}) \mapsto \gamma_d^g \hat{c}^g(\hat{r})$ , where  $\hat{c}^g(\hat{r}) = \exp(-\hat{r}^2(D^g)^2/2)$  is the FT of the *d*-dimensional Gaussian function (17). Based on the similarity of the



**Figure 4.** A comparison of the correlation functions  $c^{h}(r(\theta))$  (Eqs (5) and (30)),  $c_{2}^{w}(r(\theta))$  (Eq. (41)) and  $c_{3}^{w}(r(\theta))$  (Eq. (42)) for M = 4. The four sets of curves correspond to Daley length-scales of 500 km, 1000 km, 2000 km and 4000 km (curves from left to right, respectively). Correlations are plotted as a function of chordal distance. A spectral truncation at n = 500 has been used for  $c^{h}(r(\theta))$ .

explicit diffusion kernels on  $\mathbb{S}^2$  and  $\mathbb{R}^2$  when  $(D^g/a)^2 \ll 1$ , one would expect similar agreement between the implicit diffusion kernels on  $\mathbb{S}^2$  and  $\mathbb{R}^2$  when  $(D_2^w/a)^2 \ll 1$ .

Figure 4 shows the correlation function  $c^{h}(r(\theta))$  for M = 4 plotted as a function of chordal distance  $r(\theta)$  for four different Daley length-scales D<sup>h</sup> (solid curves). For comparison, the correlation functions  $c_2^{W}(r(\theta))$  for M = 4are also shown (dashed curves). The Daley scales  $D_2^w$ have been set to  $D^{h}$  and then the corresponding L have been computed from Eq. (44). As expected, the curves are virtually indistinguishable for length-scales of primary interest (<1000 km). Only when the length-scale exceeds 2000 km do noticeable differences appear, and these mainly occur in the tail of the function. In other words, the differential operator  $(\mathcal{C}^h)^{-1}$  on  $\mathbb{S}^2$  can be well approximated locally by the differential operator  $(\mathcal{C}_2^{\mathrm{w}})^{-1}$  acting on the tangent plane  $\mathbb{R}^2$ . The simple relations (43) and (44) can then be used in place of the spectral expansions (31) and (8) to provide a good approximation of the normalization factor and length-scale. This is convenient especially with grid-point ocean models where spectral expansions cannot be readily computed due to the presence of complex land boundaries.

It is important to stress, however, that  $c_2^{w}(r(\theta))$  itself is not a valid correlation function on  $\mathbb{S}^2$ . A valid correlation function on  $\mathbb{S}^2$  from the Matérn family is the AR function  $c_3^{w}(r(\theta))$ . For example, Gaspari and Cohn (1999) discuss the second-order AR (SOAR) function on  $\mathbb{S}^2$  (see their Eq. (2.36)). Figure 4 shows the fourth-order AR function for different length-scales (dashed-dotted curves). The differences between  $c_3^{w}(r(\theta))$  and  $c^{h}(r(\theta))$  are larger than those between  $c_2^{w}(r(\theta))$  and  $c^{h}(r(\theta))$  but still quite small for length-scales less than 1000 km.

A more general set of correlation functions on  $\mathbb{R}^d$  can be modelled using a linear combination of implicit diffusion operators or a generalized implicit diffusion operator constructed from the inverse of a polynomial of Laplacian operators raised to the power of M (MW10; Yaremchuk and Smith, 2011; or see Purser *et al.*, 2003b, for related approaches involving the recursive filter). The correlation functions generated by the first approach are described by a linear combination of Matérn functions where the weighting coefficients for each function are specified such that the combined function is positive definite. Gregori *et al.* (2008) provide general conditions on the model parameters for achieving this. MW10 provide an example in  $\mathbb{R}^1$  in which two SOAR functions are combined to produce a correlation function with negative lobes.

The second approach is analogous to the one outlined in section 2.3 for the problem on  $\mathbb{S}^2$  (Eqs (32) and (33)). Hristopulos (2003), Hristopulos and Elogne (2007) and Yaremchuk and Smith (2011) have studied extensively the special case M = 1 and P = 2 on  $\mathbb{R}^d$  for which the parameter settings  $\kappa_1 = \rho L^2$  and  $\kappa_2 = L^4$  with  $\rho < 0$ and satisfying  $\rho^2 < 4$  yield a family of positive-definite, oscillatory functions such as those illustrated in Figure 3 on  $\mathbb{S}^2$ . With all of these approaches, however, the advantages of increasing the flexibility in the correlation model have to be carefully measured against the increase in computational cost that results from the need to solve additional or more complicated large linear systems, and the difficulty of having to estimate additional parameters.

### 3. Anisotropic diffusion

Isotropic correlation models are commonly used in data assimilation algorithms because of their simplicity and computational convenience. There is no reason, however, to expect actual background-error correlations to be isotropic in geophysical fluids such as the ocean. On the contrary, one would expect them to be strongly anisotropic, particularly near coastlines, bathymetry, or ocean fronts. General anisotropic correlation models allow for preferential stretching or shrinking of the correlation functions along arbitrary directions. With a diffusion-based correlation model this can be done using a diffusion tensor, as outlined in this section. To fix the concepts and definitions, we focus mainly on the homogeneous and anisotropic problem. Methods for estimating the parameters of a general inhomogeneous and anisotropic diffusion model are described in section 4.

### 3.1. Homogeneity and anisotropy

Consider the 2D diffusion equation on  $\mathbb{R}^2$ ,

$$\frac{\partial \eta}{\partial s} - \nabla \cdot \boldsymbol{\kappa} \nabla \eta = 0, \qquad (45)$$

where  $\kappa \in \mathbb{R}^2 \times \mathbb{R}^2$  is an anisotropic, but constant *diffusion tensor* 

$$\boldsymbol{\kappa} = \begin{pmatrix} \kappa_{xx} & \kappa_{xy} \\ \kappa_{yx} & \kappa_{yy} \end{pmatrix}$$
(46)

which is assumed to be symmetric ( $\kappa_{yx} = \kappa_{xy}$ ) and positive definite ( $\kappa_{xx}\kappa_{yy} > \kappa_{xy}^2$ ) so that  $\kappa$  is guaranteed to be invertible. The diagonal terms of the tensor determine the strength of the diffusion in the coordinate directions *x* and *y*, while the off-diagonal elements allow the principal axes of the diffusion to be rotated relative to *x* and *y*.

The solution of Eq. (45) is a straightforward extension of the solution to the isotropic problem (Pannekoucke and Massart, 2008; Pannekoucke, 2009). Given the initial

condition  $\eta(x, y, 0) = \gamma_2^g \tilde{\eta}(x, y)$ , the solution can be expressed as Eq. (16) (with d = 2) where the kernel is given by the Gaussian function

$$c^{\mathsf{g}}(\widetilde{r}) = \frac{\gamma_2^{\mathsf{g}}}{4\pi s |\boldsymbol{\kappa}|^{1/2}} e^{-\widetilde{r}^2/4s},\tag{47}$$

 $|\kappa|$  denoting the determinant of  $\kappa$ , and  $\tilde{r}$  the nondimensional distance measure

$$\widetilde{r} = \sqrt{(\mathbf{x} - \mathbf{x}')^{\mathrm{T}} \, \boldsymbol{\kappa}^{-1} \, (\mathbf{x} - \mathbf{x}')},\tag{48}$$

with  $\mathbf{x} = (x, y)^{\mathrm{T}}$ . From this definition,  $\boldsymbol{\kappa}$  can also be interpreted as the *aspect tensor* of the Gaussian function (47) (Purser *et al.*, 2003b). The elements of  $\boldsymbol{\kappa}$  have physical units of length squared. Setting

$$\gamma_2^{\rm g} = 4\pi s \, |\boldsymbol{\kappa}|^{1/2}$$

ensures that  $c^{g}(0) = 1$ .

For a homogeneous and at least twice-differentiable correlation function, we can define the Hessian tensor (Swerling, 1962; Hristopulos, 2002; Chorti and Hristopulos, 2008), which for the 2D Gaussian function is

$$H^{g} = -\nabla \nabla^{T} c^{g} \Big|_{\tilde{r}=0}, \qquad (49)$$

where  $\nabla \nabla^{T}$  is the outer product of the 2D gradient operator  $\nabla = (\partial/\partial x \ \partial/\partial y)^{T}$  and its transpose. The correlation Hessian tensor is of interest here since it is a quantity that can be estimated from sample statistics of background error (see section 4). Following the basic procedure described in Appendix A, it is straightforward to verify that

$$H^{\rm g} = \frac{1}{2s} \kappa^{-1}.$$
 (50)

In the isotropic case,  $\kappa = \kappa I$  and hence  $H^{g} = (D^{g})^{-2}I$ where  $(D^{g})^{2} = 2\kappa s$  is the square of the Daley length-scale. The inverse of the tensor (49)

$$D^{g} := (H^{g})^{-1} \tag{51}$$

can thus be considered as a generalization of the (square of the) Daley length-scale to the anisotropic case. We will thus refer to this quantity as a *Daley tensor*.

For the 3D diffusion equation, the diffusion tensor  $\kappa \in \mathbb{R}^3 \times \mathbb{R}^3$  contains six independent elements:

$$\boldsymbol{\kappa} = \begin{pmatrix} \kappa_{xx} & \kappa_{xy} & \kappa_{xz} \\ \kappa_{yx} & \kappa_{yy} & \kappa_{yz} \\ \kappa_{zx} & \kappa_{zy} & \kappa_{zz} \end{pmatrix}$$
(52)

where  $\kappa_{yx} = \kappa_{xy}$ ,  $\kappa_{zx} = \kappa_{xz}$  and  $\kappa_{zy} = \kappa_{yz}$ . In direct analogy with the 2D problem, the integral solution involves a 3D Gaussian kernel with aspect tensor given by (52), and normalization constant given by  $\gamma_3^{\rm g} = (4\pi s)^{3/2} |\kappa|^{1/2}$  (cf. Eq. (18)). The relationships (49)–(51) hold for the 3D problem with  $\nabla$  now interpreted as the 3D gradient operator.

To approximate a 2D or 3D anisotropic and homogeneous Gaussian correlation operator numerically, we can solve Eq. (45) with an explicit scheme,

$$\eta(\mathbf{x}, M) = \gamma_d^{\mathsf{g}} (1 + \nabla \cdot \boldsymbol{\kappa} \nabla)^M \widetilde{\eta}(\mathbf{x}),$$

Copyright © 2012 Royal Meteorological Society

where from Eqs (50) and (51)

$$\kappa = \frac{1}{2M} D^{\mathrm{g}},\tag{53}$$

and the operator  $1 + \nabla \cdot \kappa \nabla$  is understood to be in discrete form. If the non-diagonal tensor elements of  $\kappa$  are zero, which can always be achieved by rotating the model coordinates to be aligned with the principal axes of the ellipse or ellipsoid implied by Eq. (48) (see, for example, Xu, 2005), then the 2D or 3D Gaussian operator can be replaced by a product of 1D Gaussian operators acting independently along each direction *x*, *y* and *z*. Ignoring boundary conditions, each 1D Gaussian operator can in turn be approximated by a 1D diffusion operator discretized using an *M*-step explicit scheme.

Extending these results to the *d*-dimensional implicit case, we can define a set of anisotropic and homogeneous Matérn correlation operators, with v = M - d/2, as solutions to the following linear system (cf. Eq. (34)):

$$(\gamma_d^{\mathrm{w}})^{-1}(1 - \nabla \cdot \kappa \nabla)^M \eta(\mathbf{x}, M) = \widetilde{\eta}(\mathbf{x})$$

where

$$\gamma_d^{w} = 2^d \pi^{d/2} \frac{\Gamma(M)}{\Gamma(M - d/2)} |\kappa|^{1/2}.$$
 (54)

The associated correlation functions are given by

$$c_d^{\rm w}(\tilde{r}) = \frac{2^{1-M+d/2}}{\Gamma(M-d/2)} \tilde{r}^{M-d/2} K_{M-d/2}(\tilde{r}), \qquad (55)$$

with  $\tilde{r}$  defined by Eq. (48). As for the Gaussian, we can derive the following relationships between the Hessian tensor of  $c_d^{\kappa}(\tilde{r})$  and diffusion tensor  $\kappa$  (see Appendix A):

$$\begin{aligned}
 H_d^{\mathsf{w}} &= -\nabla \nabla^{\mathsf{T}} c_d^{\mathsf{w}} |_{\widetilde{r}=0}, \\
 D_d^{\mathsf{w}} &:= (H_d^{\mathsf{w}})^{-1}, \\
 \kappa &= \frac{1}{2M - d - 2} D_d^{\mathsf{w}}.
 \end{aligned}$$
(56)

The solution described in section 2.4 corresponds to the isotropic case  $\kappa = L^2 \mathbf{I}$  with  $L^2 = (D_d^w)^2 / (2M - d - 2)$ .

#### 3.2. Inhomogeneity and anisotropy

Analytical expressions for the correlation kernels of the anisotropic diffusion operators in  $\mathbb{R}^d$  with spatially varying diffusion tensors  $\kappa(\mathbf{x})$  are not known in general. Paciorek and Schervish (2006) describe a family of anisotropic and inhomogeneous correlation functions that generalize the standard isotropic and homogeneous Gaussian and Matérn family. These correlation functions have the form

$$C^{g}(\mathbf{x}, \mathbf{x}') = \beta(\mathbf{x}, \mathbf{x}') \exp(-\tilde{r}^{2}/2)$$
(57)

for the Gaussian-like function, and

$$C_d^{\mathsf{w}}(\mathbf{x}, \mathbf{x}') = \beta(\mathbf{x}, \mathbf{x}') \\ \times \frac{2^{1-M+d/2}}{\Gamma(M-d/2)} \widetilde{r}^{M-d/2} K_{M-d/2}(\widetilde{r})$$
(58)

for the Matérn-like functions ( $\nu = M - d/2$ ) where

$$\widetilde{r} = \sqrt{\left(\mathbf{x} - \mathbf{x}'\right)^{\mathrm{T}} \left(\frac{A(\mathbf{x}) + A(\mathbf{x}')}{2}\right)^{-1} \left(\mathbf{x} - \mathbf{x}'\right)^{\mathrm{T}}}$$

and

$$\beta(\mathbf{x}, \mathbf{x}') = |A(\mathbf{x})|^{1/4} |A(\mathbf{x}')|^{1/4}$$
$$\times \left|\frac{1}{2} (A(\mathbf{x}) + A(\mathbf{x}'))\right|^{-1/2},$$

 $A(\mathbf{x})$  and  $A(\mathbf{x}')$  denoting the (symmetric and positive definite) aspect tensors at points  $\mathbf{x}$  and  $\mathbf{x}'$ , respectively. Equations (57) and (58) with  $A(\mathbf{x}) \approx 2s \kappa(\mathbf{x})$  for the Gaussian-like function and  $A(\mathbf{x}) \approx \kappa(\mathbf{x})$  for the Matérn-like functions can be considered as the approximate kernels of the explicit and implicit forms of the anisotropic diffusion operator when the diffusion tensors  $\kappa(\mathbf{x})$  vary slowly and smoothly in space. This is illustrated in MW10 who provide examples in 1D comparing a two-step implicit-diffusion kernel and an inhomogeneous version of the SOAR function for different spatial distributions of the length-scale parameter.

#### 4. Specifying the anisotropic tensor

The elements  $\kappa_{xz}$ ,  $\kappa_{zx}$ ,  $\kappa_{yz}$  and  $\kappa_{zy}$  of the 3D diffusion tensor account for anisotropy between the horizontal and vertical directions. The importance of these terms compared to the diagonal terms is related to the choice of vertical coordinate in the correlation model. In an ocean model, for example, a natural vertical coordinate is a hybrid coordinate involving a standard geopotential (z)coordinate in unstratified regions such as the mixed layer, an isopycnal  $(\rho)$  coordinate in strongly stratified regions, and a terrain-following (s) coordinate near the ocean bottom, the latter being particularly important in shallow coastal regions (Haidvogel and Beckmann, 1999). In this hybrid coordinate system, the flow is more naturally decoupled into 'horizontal' and 'vertical' processes. If the same coordinate system is adopted for a background-error correlation model then it is reasonable to assume, at least from a physical viewpoint, that the non-diagonal tensor elements  $\kappa_{xz}$ ,  $\kappa_{zx}$ ,  $\kappa_{yz}$  and  $\kappa_{zy}$ , and possibly  $\kappa_{xy}$  and  $\kappa_{yx}$ , are small and can be neglected. However, anisotropy in background-error correlations can also arise from the assimilation of data, especially when the data coverage is irregular. In general, the relative importance of the diagonal and non-diagonal terms of the tensor can only be determined after a thorough diagnostic study involving, for instance, the direct estimation of the elements of the Daley tensor.

Many ocean models used for global- and basinscale circulation studies employ a z coordinate. WC01 illustrated how a standard isopycnal diffusion tensor used to parametrize mixing of unresolved processes in a z-coordinate ocean model could also be used to transform the coordinates of a background-error correlation model formulated as an explicit 3D diffusion operator. An analogous coordinate transformation was proposed within the framework of Optimal Interpolation by Balmaseda *et al.* (2008). While the isopycnal correlation model has appealing features, the implementation based on the explicit scheme proposed by WC01 is too expensive for routine applications since a prohibitively high number of iterations is required to maintain numerical stability in regions of strong isopycnal gradients. Moreover, the specification of the Daley length-scales must be performed in isopycnal space, which makes estimating them more difficult in a *z*-coordinate model. In the remainder of this section we explore alternative methods for defining anisotropic and inhomogeneous correlations, which involve estimating the Daley tensor directly in the model coordinate system.

#### 4.1. Ensemble estimation methods

Given an estimate of the Daley tensor, the anisotropic response of the explicit diffusion operator can be calibrated using Eq. (53), which relates the Daley tensor of the Gaussian function to the diffusion tensor. Alternatively, the anisotropic response of the implicit diffusion operator can be calibrated using the third expression in (56), which relates the Daley tensors of the v = M - d/2 Matérn functions to the diffusion tensor. Several authors have proposed methods for estimating the Daley tensor using perturbations from an ensemble of model states (Belo Pereira and Berre, 2006; Pannekoucke and Massart, 2008; Pannekoucke, 2009; Sato *et al.*, 2009). The basic procedure is outlined below. Two of the methods will then be compared in idealized experiments using the diffusion equation. For simplicity, we focus on the 2D case.

Assume that an ensemble of  $N_e$  model states is available and that the distribution of these states about their mean is a good approximation of the *true* probability distribution function (pdf) of the model-state (background) error  $\epsilon$ . In variational assimilation, this pdf is assumed to be Gaussian and thus fully described by its mean ( $E[\epsilon(\mathbf{x})] = 0$ ) and covariance function

$$B(\mathbf{x}, \mathbf{x}') = E[\epsilon(\mathbf{x}) \epsilon(\mathbf{x}')], \qquad (59)$$

where  $E[\cdot]$  denotes the expectation operator. The associated correlation function  $C(\mathbf{x}, \mathbf{x}')$  can be determined from the factorization

$$C(\mathbf{x}, \mathbf{x}') = \frac{B(\mathbf{x}, \mathbf{x}')}{\sigma(\mathbf{x}) \sigma(\mathbf{x}')},$$
(60)

where

$$\sigma(\mathbf{x}) = \sqrt{E[\epsilon(\mathbf{x})^2]}$$

is the standard deviation of  $\epsilon$  at **x**. Assuming that  $C(\mathbf{x}, \mathbf{x}')$  is at least twice differentiable then we can define the symmetric tensor

$$T(\mathbf{x},\mathbf{x}') = -\nabla \nabla'^{\mathrm{T}} C(\mathbf{x},\mathbf{x}'),$$

where  $\nabla = (\partial/\partial x \ \partial/\partial y)^{\mathrm{T}}$  and  $\nabla' = (\partial/\partial x' \ \partial/\partial y')^{\mathrm{T}}$ . The local correlation Hessian tensor is the value of T at  $\mathbf{x} = \mathbf{x}'$ . Assume further that, in a neighbourhood of  $\mathbf{x}$ ,  $C(\mathbf{x},\mathbf{x}')$  can be well approximated by a homogeneous function  $\widetilde{C}(\mathbf{r})$  where  $\mathbf{r} = \mathbf{x} - \mathbf{x}' = (x - x', y - y')^{\mathrm{T}} = (\widetilde{x}, \widetilde{y})^{\mathrm{T}}$ . Letting  $\widetilde{\nabla} = (\partial/\partial \widetilde{x} \ \partial/\partial \widetilde{y})^{\mathrm{T}}$ , then we can define

$$\widetilde{T}(\mathbf{r}) = -\widetilde{\nabla}\,\widetilde{\nabla}^{\mathrm{T}}\widetilde{C}(\mathbf{r})\,,\tag{61}$$

such that  $T(0) \approx T(\mathbf{x}, \mathbf{x})$  (see Appendix B).

Let  $H(\mathbf{x}) = T(\mathbf{x}, \mathbf{x})$  denote the local correlation Hessian tensor at  $\mathbf{x}$ . Pannekoucke and Massart (2008) and Pannekoucke (2009) assume a Gaussian form for the correlation function and then invert this function at each point  $\mathbf{x}$  to estimate  $H(\mathbf{x})$  in terms of sample correlation estimates with neighbouring points. Belo Pereira and Berre (2006) propose an alternative method for estimating  $H(\mathbf{x})$ , which does not require a prior assumption on the functional form of the correlations. Their method leads to the expression

$$\widehat{H}(\mathbf{x}) = \frac{\overline{\nabla \epsilon(\mathbf{x}) (\nabla \epsilon(\mathbf{x}))^{\mathrm{T}} - \nabla \widehat{\sigma}(\mathbf{x}) (\nabla \widehat{\sigma}(\mathbf{x}))^{\mathrm{T}}}}{(\widehat{\sigma}(\mathbf{x}))^{2}}, \qquad (62)$$

where

$$\overline{\nabla \epsilon(\mathbf{x}) (\nabla \epsilon(\mathbf{x}))^{\mathrm{T}}} = \frac{1}{N_{\mathrm{e}} - 1} \sum_{l=1}^{N_{\mathrm{e}}} \nabla \epsilon_{l}(\mathbf{x}) (\nabla \epsilon_{l}(\mathbf{x}))^{\mathrm{T}}$$
  
and  $(\widehat{\sigma}(\mathbf{x}))^{2} = \overline{(\epsilon(\mathbf{x}))^{2}} = \frac{1}{N_{\mathrm{e}} - 1} \sum_{l=1}^{N_{\mathrm{e}}} (\epsilon_{l}(\mathbf{x}))^{2}.$ 

The gradient terms can be estimated numerically using finite differences. When the correlation function is strictly homogeneous,  $\hat{H}$  is equivalent to the constant tensor  $\tilde{T}(\mathbf{0})$  (Eq. (61)) if the sampling operator is replaced by the expectation operator (see Appendix B). Note that the derivation of Eq. (62) is based on the rather general assumptions of differentiability and local homogeneity of the correlation function. A specific assumption about the actual form of the correlation function is implied only when  $\hat{H}$  is employed with a particular diffusion operator (explicit or *M*-step implicit scheme). Hristopulos (2002) and Chorti and Hristopulos (2008) describe a related approach in geostatistics which involves estimating the aspect ratios and orientation angle required to transform the local covariance Hessian tensor into isotropic form.

Multidimensional Gaussian correlation operators can be applied efficiently using a combination of 1D recursive filters or 1D diffusion operators (Purser et al., 2003a, 2003b; MW10). For anisotropic Gaussian operators, the so-called triad (hexad) algorithm (Purser et al., 2003b; Purser, 2005) allows one to determine from the aspect tensor of the 2D (3D) Gaussian function, the three (six) generalized gridlines along which the 1D filters should be applied. Within that framework, various flow-dependent formulations of the aspect tensor have been proposed (De Pondeca et al., 2006; Liu et al., 2007, 2009; Sato et al., 2009). Of particular interest here is the hybrid formulation of Sato et al. (2009) where the inverse of the aspect tensor of the Gaussian function is defined as a linear combination of a 'conventional term' based on a quasi-isotropic, static formulation  $(A_{iso}^{-1})$  and an 'ensemble term' formed from the sample covariance of the gradient of the ensemble-generated perturbations, normalized by the sample variance of the perturbations:

$$\widehat{A}^{-1}(\mathbf{x}) = \alpha A_{\rm iso}^{-1}(\mathbf{x}) + \beta \widetilde{H}(\mathbf{x})$$
(63)

where

$$\widetilde{H}(\mathbf{x}) = \frac{\overline{\nabla \epsilon(\mathbf{x}) (\nabla \epsilon(\mathbf{x}))^{\mathrm{T}}}}{(\widehat{\sigma}(\mathbf{x}))^{2}}, \qquad (64)$$

Copyright © 2012 Royal Meteorological Society

and  $\alpha$  and  $\beta$  are weighting coefficients. Sato *et al.* (2009) provide a heuristic derivation of Eqs (63)–(64). They are equivalent to Eq. (62) when  $\alpha = 0$  and  $\beta = 1$ , and when the standard deviations are constant. As with the hybrid covariance formulations that involve combining static and ensemble-based expressions of the *full* covariance matrix (Wang *et al.*, 2008), the static term in the aspect tensor is intended to give the estimate more robustness especially when the ensemble size is small, although accounting for it requires extra parameters that must be tuned empirically.

Finally, with small ensemble sizes it can be advantageous to apply a local spatial averaging or filtering operator to the estimated variances and covariances in order to reduce the effects of sampling error (see, for example, Berre and Desroziers (2010) for a thorough review of recent work in this area). Letting F denote a particular filtering operator then the expressions for the filtered estimate of the inverse tensor are given by Eqs (62) and (63)–(64) with

$$\overline{\nabla \epsilon(\mathbf{x}) (\nabla \epsilon(\mathbf{x}))^{\mathrm{T}}} \mapsto F\left(\overline{\nabla \epsilon(\mathbf{x}) (\nabla \epsilon(\mathbf{x}))^{\mathrm{T}}}\right),$$
$$(\widehat{\sigma}(\mathbf{x}))^{2} \mapsto F\left((\widehat{\sigma}(\mathbf{x}))^{2}\right),$$
$$\widehat{\sigma}(\mathbf{x}) \mapsto \sqrt{F\left((\widehat{\sigma}(\mathbf{x}))^{2}\right)}.$$

and

## 4.2. Numerical experiments

In this section we perform idealized experiments to evaluate and compare the effectiveness of Eq. (62) and Eqs (63)–(64) for estimating the parameters of an anisotropic tensor. For simplicity, we focus on the 2D anisotropic diffusion problem and the solution algorithm based on the explicit scheme. Furthermore, for the tensor estimated using Eqs (63)–(64), only the special case  $\alpha = 0$  and  $\beta = 1$  is considered.

The experimental design is as follows. First, we define the 'true' covariance matrix of the problem as

$$\mathbf{B} = \boldsymbol{\Sigma} \, \boldsymbol{\Gamma}^{1/2} \, \mathbf{L} \, \boldsymbol{\Gamma}^{1/2} \, \boldsymbol{\Sigma},$$

where **L** is the *M*-step explicit diffusion operator  $(1 + \nabla \cdot \kappa \nabla)^M$  discretized using a standard centred finitedifference scheme on a uniform grid,  $\Gamma = \Gamma^{1/2} \Gamma^{1/2}$  is a diagonal matrix of normalization factors, and  $\Sigma$  is a diagonal matrix of standard deviations  $\sigma$ . With constant parameters and ignoring the influence of boundaries, **B** defines a Gaussian covariance matrix.

Next, a sample of  $N_e$  spatially uncorrelated random vectors  $\hat{\epsilon}_l, l = 1, ..., N_e$ , are produced on the grid, where the distribution of each  $\hat{\epsilon}_l$  is taken to be Gaussian with  $E[\hat{\epsilon}_l] = \mathbf{0}$  and  $E[\hat{\epsilon}_l \hat{\epsilon}_l^T] = \mathbf{I}$ . Each vector  $\hat{\epsilon}_l$  is then transformed into a new vector  $\epsilon_l$  such that  $E[\epsilon_l \epsilon_l^T] = \mathbf{B}$ . This is done using the 'square-root' of the **B**-operator,

$$\boldsymbol{\epsilon}_l = \boldsymbol{\Sigma} \, \boldsymbol{\Gamma}^{1/2} \, \mathbf{L}^{1/2} \, \boldsymbol{\widehat{\epsilon}}_l,$$

where  $\mathbf{L} = \mathbf{L}^{1/2} \mathbf{L}^{1/2}$ , the exponent 1/2 implying *M*/2 iterations of the explicit diffusion operator (with *M* taken to be even). The sample covariance matrix constructed from the ensemble of  $\epsilon_l$  vectors provides an estimate of the true covariance matrix:

$$\mathbf{B} \approx \frac{1}{N_{\rm e} - 1} \sum_{l=1}^{N_{\rm e}} \boldsymbol{\epsilon}_l' \boldsymbol{\epsilon}_l'^{\rm T},\tag{65}$$

where  $\boldsymbol{\epsilon}_l' = \boldsymbol{\epsilon}_l - \frac{1}{N_{\rm e}} \sum_{k=1}^{N_{\rm e}} \boldsymbol{\epsilon}_k.$ 

Our interest here is not to try to reconstruct the full covariance matrix from (65) but rather, with the help of Eq. (62) or Eq. (64), to try to reconstruct the anisotropic tensor used in L to generate the  $\epsilon_l$ . Indeed, the sampling errors resulting from estimating the local anisotropic tensor can be expected to be much smaller than those resulting from estimating the full covariance matrix. For the 2D problem, the tensor estimation requires, at each grid point, sample estimates of the standard deviation of  $\epsilon$  and of the three independent tensor elements involving the gradient of  $\epsilon$ , i.e. a total of 4N elements where N is the total number of grid points. This is much smaller than the  $(N^2 + N)/2$  independent elements required to determine the full covariance matrix.

The numerical experiments are performed in a square domain on a 2D grid  $\mathbf{x}_{i,i} = (x_i, y_i)$  where  $x_i =$  $i\Delta x$ ,  $i = 1, \dots, \sqrt{N}$ , and  $y_j = j\Delta y$ ,  $j = 1, \dots, \sqrt{N}$ . Here,  $\Delta x = \Delta y = 1$  unit and  $N = 200 \times 60$ , and thus the effective size of the **B** matrix is  $(1.2 \times 10^4) \times (1.2 \times 10^4)$ . Neumann boundary conditions are employed at the solid walls located at the domain edges. As a result, the implied correlation function near the boundary is slightly modified from the target Gaussian (MW10).

At each grid point  $\mathbf{x}_{i,j}$ , the 'true' diffusion tensor  $\kappa$  for L is defined according to (cf. Eq. (53))

$$\boldsymbol{\kappa}_{i,j} = \frac{1}{2M} \boldsymbol{D}_{i,j},\tag{66}$$

where  $D_{i,j}$  is the local Daley tensor which is formulated as

$$D_{i,j} = R_{i,j}\overline{D}_{i,j}R_{i,j}^{-1}$$

 $\overline{D}_{i,j}$  being a diagonal matrix and  $R_{i,j}$  a rotation matrix ( $R_{i,j}^{-1} =$  $\boldsymbol{R}_{i,j}^{\mathrm{T}}$ ). The elements  $(D_{xx})_{i,j}$ ,  $(D_{yy})_{i,j}$  and  $(D_{xy})_{i,j} = (D_{yx})_{i,j}$  of  $D_{i,j}$  are thus determined by the diagonal elements  $(\overline{D}_{xx})_{i,j}$ and  $(\overline{D}_{yy})_{i,j}$  of  $\overline{D}_{i,j}$  and by the rotation angle  $\theta_{i,j}$  of  $R_{i,j}$ . For the experiments described here,  $\theta_{ij} = \theta$  is constant, while the parameters  $(\overline{D}_{xx})_{ij}$  and  $(\overline{D}_{yy})_{ij}$  are specified as a simple oscillatory function of the spatial coordinates  $\mathbf{x}_{i,j}$ :

$$(\overline{D}_{xx})_{i,j} = A_1 f(\mathbf{x}_{i,j}) + B_1,$$
  
$$(\overline{D}_{yy})_{i,j} = -A_1 f(\mathbf{x}_{i,j}) + B_1,$$

where

$$f(\mathbf{x}_{i,j}) = \cos\left(\frac{2\pi x_i}{X}\right)\cos\left(\frac{2\pi y_j}{Y}\right),$$

 $\begin{array}{l} A_1 = \frac{1}{2} (\overline{D}_{\max} - \overline{D}_{\min}), \qquad B_1 = \frac{1}{2} (\overline{D}_{\max} + \overline{D}_{\min}), \\ X = Y = 20. \text{ Similarly, the variances are specified as} \end{array}$ and

$$\sigma_{i,j}^2 = A_2 f(\mathbf{x}_{i,j}) + B_2,$$

where  $A_2 = \frac{1}{2}(\sigma_{\max}^2 - \sigma_{\min}^2)$  and  $B_2 = \frac{1}{2}(\sigma_{\max}^2 + \sigma_{\min}^2)$ . Experiments are performed with different values of the parameters  $\theta$ ,  $\overline{D}_{\min}$ ,  $\overline{D}_{\max}$ ,  $\sigma_{\min}^2$  and  $\sigma_{\max}^2$  (see Table 2). The normalization factors  $\gamma_{i,j}$  of the diagonal matrix  $\Gamma$ 

are approximately given by the expression

$$\gamma_{i,j} \approx 2\pi |\boldsymbol{D}_{i,j}|^{1/2}.$$
 (67)

This approximation was used by Pannekoucke and Massart (2008), for example, and is reasonable if the

diffusion tensor varies in space on a scale much larger than the local correlation scale and with a proper treatment of the boundary conditions (MW10). The factors can be estimated to a higher accuracy using more refined analytical approximations (Purser et al., 2003b; Purser, 2008a, 2008b; MW10; Yaremchuk and Carrier, 2012) or randomization methods (WC01; Yaremchuk and Carrier, 2012). They can also be computed exactly using the  $\delta$ -function method (WC01; MW10). In this idealized study, we employ the exact normalization method in order to avoid introducing a bias in the ensemble perturbations and thus complicating the interpretation of the results. In practice, however, the exact computation is generally not affordable and hence the representation of covariances using the diffusion equation will also be affected by approximations in the normalization factors. The errors that can result from using the approximate expression (67) are illustrated in the experiments below.

The estimation of the tensor via the statistical relationships (62) and (64) is achieved using centred finite-differences. Estimates of the first derivatives of the error and its standard deviation produce values at the interface of the grid cells, i.e. at the half-integer points (i + 1/2, j) for the x-component and (i, j + 1/2) for the y-component. The sample variance of these quantities is computed directly at these points to evaluate the numerator in the expressions for the diagonal elements of the tensor. The off-diagonal elements involve estimates of the cross-product of the x- and y-components of the derivatives. This requires interpolation of one of the component derivatives to the point where the other component derivative is defined. To estimate the cross-product at (i + 1/2, i), the x-component derivative that is defined there is multiplied with an estimate of the y-component derivative obtained by averaging its values from the four surrounding points (i, j + 1/2), (i+1, j+1/2), (i, j-1/2) and (i+1, j-1/2), and viceversa for estimating the cross-product at (i, j + 1/2). To compute the denominator in the expressions for the tensor elements, the sample variance of the error is interpolated from (i, j) points to (i + 1/2, j) or (i, j + 1/2) points. In order to use the estimated tensor elements in the diffusion equation, the elements are first averaged to the (i, j) points and then the off-diagonal elements averaged to force symmetry. The estimated tensor is then inverted at each point and used with the relation (66) to define the diffusion tensor at each point. Finally, interpolation is used to define the values of the tensor elements at the half-integer points (i + 1/2, j) or (i, j + 1/2) where they are required with the centred-difference formulation of  $\nabla \cdot \kappa \nabla$ .

Table 2 summarizes the results from several experiments with different parameter settings  $\mathcal{P} = (\theta, \overline{D}_{\min}, \overline{D}_{\max}, \sigma_{\min}, \sigma_{\max})$ . Three cases are considered. In the first case, the principal axes of the anisotropic correlations are aligned with the grid-lines, and the variance is constant:  $\mathcal{P}_1 = (0, 3, 6, 1, 1)$ . The second case extends the first case by allowing the variances to vary in space:  $\mathcal{P}_2 = (0, 3, 6, 1, 5)$ . Finally, the third case extends the second case by rotating the principal axes of the anisotropic correlations relative to the grid-lines:  $\mathcal{P}_3 = (\pi/4, 3, 6, 1, 5)$ . The quality of the estimation is measured in terms of the domainaveraged bias and root-mean-square error (RMSE) of the estimates of the elements  $(H_{xx}, H_{yy}, H_{xy})$ . (The results for the estimates of  $H_{yx}$  are not given since they are almost identical to those of  $H_{xy}$ ). For reference, the domain-averaged RMS

Copyright © 2012 Royal Meteorological Society

Exp	$(\theta, D_{\min}, D_{\max}, \sigma_{\min}, \sigma_{\max})$	$(N_{\rm e},N_{\rm avg})$	Method	$(H_{xx}, H_{yy}, H_{xy})$	
				Bias $\times 10^{-2}$	RMSE $\times 10^{-2}$
1	(0, 3, 6, 1, 1)	(100, 0)	$\widehat{H}$	(0.06, 0.20, 0.0)	(1.2, 1.6, 0.46)
2	(0, 3, 6, 1, 1)	(100, 0)	$\widetilde{H}$	(-0.01, 0.15, 0.0)	(1.2, 1.6, 0.47)
3	(0, 3, 6, 1, 1)	(10, 0)	$\widehat{H}$	(0.33, 0.39, 0.07)	(4.6, 4.8, 2.1)
4	(0, 3, 6, 1, 1)	(10, 0)	$\widetilde{H}$	(0.92, 0.99, 0.10)	(4.9, 5.2, 2.2)
5	(0, 3, 6, 1, 5)	(100, 0)	$\widehat{H}$	(-0.07, -0.22, 0.0)	(1.2, 1.6, 0.47)
6	(0, 3, 6, 1, 5)	(100, 0)	$\widetilde{H}$	(1.2, 1.1, -0.01)	(2.5, 2.8, 1.2)
7	$(\pi/4, 3, 6, 1, 5)$	(100, 0)	$\widehat{H}$	(-0.22, -0.40, 0.01)	(1.4, 1.7, 0.85)
8	$(\pi/4, 3, 6, 1, 5)$	(100, 0)	$\widetilde{H}$	(1.1, 0.88, 0.0)	(2.6, 2.7, 1.4)
9	$(\pi/4, 3, 6, 1, 5)$	(10, 0)	$\widehat{H}$	(0.17, 0.10, -0.02)	(4.0, 4.5, 2.3)
10	$(\pi/4, 3, 6, 1, 5)$	(10, 1)	$\widehat{H}$	(0.33, 0.28, 1.9)	(3.3, 3.6, 1.9)
11	$(\pi/4, 3, 6, 1, 5)$	(10, 3)	$\widehat{H}$	(0.51, 0.45, -0.36)	(2.0, 2.2, 2.8)
12	$(\pi/4, 3, 6, 1, 5)$	(100, 1)	$\widehat{H}$	(0.02, 0.14, -0.07)	(1.2, 1.4, 0.86)
13	$(\pi/4, 3, 6, 1, 5)$	(100, 3)	$\widehat{H}$	(0.31, 0.19, -0.32)	(0.98, 1.0, 1.1)

Table 2. Bias and RMSE of the estimates of the correlation Hessian tensor elements  $(H_{xx}, H_{yy}, H_{xy})$  using expressions (62)  $(\widehat{H})$  and (64)  $(\widetilde{H})$ . The second column lists the parameter settings in the 'true' covariance model and the third column indicates the choice of ensemble size  $(N_e)$  and spatial filtering scale  $(N_{avg})$  used in the estimation process. The RMS of the true values of  $(H_{xx}, H_{yy}, H_{xy})$  are (5.3, 5.3, 0) × 10<sup>-2</sup> when  $\theta = 0$ , and (5.0, 5.0, 1.7) × 10<sup>-2</sup> when  $\theta = \pi/4$ .

of the true values of  $(H_{xx}, H_{yy}, H_{xy})$  are  $(5.3, 5.3, 0) \times 10^{-2}$ when  $\theta = 0$ , and  $(5.0, 5.0, 1.7) \times 10^{-2}$  when  $\theta = \pi/4$ .

With  $\mathcal{P}_1$ ,  $\widehat{H}$  (Eq. (62)) and  $\widetilde{H}$  (Eq. (64)) produce similar results with a relatively large ensemble ( $N_e = 100$ ) as one might expect since the true variances are constant (Exps 1–2 in Table 2). Interestingly, however,  $\widehat{H}$  is noticeably more accurate than  $\widetilde{H}$  with a small ensemble size ( $N_e = 10$ ; Exps 3–4). When the variances are spatially varying ( $\mathcal{P}_2$  and  $\mathcal{P}_3$ ), the errors for  $\widetilde{H}$  become significantly larger, whereas those for  $\widehat{H}$  are similar to the constant variance case (Exps 5-8). This illustrates the importance of the second term in Eq. (62).

Local spatial filtering is beneficial for reducing the RMSE especially when the ensemble size is small ( $N_e = 10$ ; Exps 9–11). With the raw ensemble estimates  $(N_{\text{avg}} = 0)$ , the RMSE is comparable to the RMS of the true signal (Exps 3-4, 9). Here, a very simple filtering procedure has been used in which the estimate at  $x_{i,j}$  is obtained by averaging estimates at points within  $N_{\text{avg}}$  grid points of (i, j) where in the examples considered  $N_{avg} = 1$  or 3. This increases the size of the averaging sample at each point to  $N_{\rm eff} = (2N_{\rm avg} + 1)^2 \times N_{\rm e}$ , except near the boundary where fewer points are used in the averaging process. While increasing the value of  $N_{avg}$  reduces the RMSE, it does so at the expense of increasing the bias in the estimates. With a larger ensemble ( $N_e = 100$ ), good results are obtained when a 'light' filtering is applied  $(N_{\text{avg}} = 1)$ , with both the bias and RMSE being reduced relative to the no filtering case for all but the off-diagonal elements which are very slightly degraded (Exps 7, 12–13). The filter in this example is very simple and the choice of filtering scale is somewhat ad hoc. More sophisticated (objective) filters could be expected to perform better as discussed by Raynaud et al. (2009) and Berre and Desroziers (2010), and recently by Raynaud and Pannekoucke (2012) within the context of filters based on diffusion.

The correlations obtained using the 'true' tensor with the parameter settings  $\mathcal{P}_3$  are illustrated at selected points in Figure 5(a). Figure 5(b) shows the corresponding correlations estimated directly from the sample covariance matrix (Eq. (65)) with a 100-member ensemble. Sampling errors are large and manifest themselves as spurious non-local correlations. In contrast, the diffusion-based correlation model is localized by construction. The correlations resulting from estimating the diffusion tensor from the 100-member ensemble are shown in Figure 5(c). The estimated correlations are in good agreement with the true correlations and notably capture prominent anisotropic features such as the rotation of the principal axes relative to the grid lines. The third correlation pattern from the left boundary is computed with respect to a point that is located midway between maximum and minimum values of  $D_{xx}$ ,  $D_{yy}$  and  $\sigma^2$  where the spatial derivative of these parameters is maximum and thus where one would expect the local homogeneous assumption to be least valid. At this location the estimated errors are largest and up to 20% (Figure 5d). The breakdown of homogeneity also affects the accuracy of the approximate expression (67) for the normalization factors. This can be seen in Figure 5 parts (e) and (f), which show the estimated correlations and associated error when approximate normalization factors from Eq. (67) are used in place of the exact factors that were used to produce Figure 5(c). The amplitude of the error now reaches 50% for the third correlation pattern (the colour bar is truncated at 30%) and is noticeably larger for the other correlation patterns as well. Finally, Figure 5 parts (g) and (h) show the correlations and associated errors obtained using the tensor estimated with  $N_e = 10$ combined with local spatial averaging. While the correlations are not as accurate as those with  $N_e = 100$ , they are still reasonable approximations. The maximum error for the third correlation pattern is approximately 25% and reaches 36% when the approximate normalization factors are used (not shown).

#### 5. Summary and conclusions

Accounting for general background-error correlations effectively and efficiently is a considerable challenge in geophysical data assimilation. In VDA, general backgrounderror correlation models can be defined using differential operators constructed numerically from the explicit or implicit solution of a diffusion equation. Theoretical results underpinning the diffusion approach to correlation modelling were reviewed in this paper. First, the isotropic, constant-coefficient diffusion problem was considered both



**Figure 5.** (a) The 'true' correlations at selected points in the domain. (b) The correlations estimated directly from the sample covariance matrix (Eq. (65)) with 100 ensemble members. The correlations produced using the diffusion equation with the diffusion tensor estimated with (c) 100 members, no local spatial averaging, and exact normalization; (e) 100 members, no local spatial averaging, and normalization factors approximated using Eq. (67); and (g) 10 members, local spatial averaging with a 7 × 7 grid-point window, and exact normalization. The differences between the correlations obtained using the estimated diffusion tensor and the true correlations are illustrated in panels (d), (f) and (h).

on the sphere and in the *d*-dimensional Euclidean space  $\mathbb{R}^d$ . The covariance functions (kernels) of the integral solution operators implied by explicit and implicit diffusion in these spaces were identified. The solutions on the sphere were shown to be well approximated by the solutions on  $\mathbb{R}^2$  for scales of interest in ocean and meteorological data assimilation. Expressions relating the diffusion model parameters to the parameters that control the length-scale and amplitude (normalization factor) of the covariance function were also given. These results provided the basis for constructing more general correlation operators via anisotropic diffusion, which was the focus of the second part of the paper.

Anisotropic diffusion was considered in  $\mathbb{R}^d$ . The anisotropic diffusion problem is characterized by a diffusion tensor that controls the direction of the covariance response, as well as its scale and amplitude. Solutions to the

anisotropic, constant-tensor diffusion problem are integral operators that involve covariance kernels with the same basic form as those of the isotropic, constant-coefficient problem. With the explicit scheme, these functions are approximately Gaussian, whereas with the implicit scheme they are members of the larger Matérn family (e.g., in  $\mathbb{R}^3$  they are AR functions). For the anisotropic functions, distance is defined by a norm whose metric is given by the inverse of the diffusion tensor. This metric can in turn be related to the correlation Hessian tensor which is defined by the tensor of second-derivatives of the correlation function evaluated at zero separation. The importance of this tensor is that it can be related to quantities that can be estimated directly from ensemble statistics. The inverse of the correlation Hessian tensor was referred to as the Daley tensor in this paper in view of its close connection to the conventional Daley length-scale in the isotropic case.

Ensemble data assimilation methods can be used to provide flow-dependent estimates of the backgrounderror covariances. In realistic applications, the number of independent background-error covariances that need to be estimated is huge and the number of ensemble members that can be affordably run is very limited. Methods are then required to synthesize the ensemble-covariance information to avoid manipulating huge covariance matrices, on the one hand, and to reduce the effects of sampling error, on the other.

The correlation information in the ensemble can be synthesized using a diffusion model with an anisotropic and spatially varying tensor.<sup>†</sup> Procedures for estimating the local Hessian tensor (which in turn can be related to the diffusion tensor) from ensemble perturbations were described and compared in idealized numerical experiments. The method of Belo Pereira and Berre (2006), which assumes local homogeneity of the correlation function but accounts for spatially varying variances, was shown to work particularly well, and is well suited for the automated computations required in a cycled ensemble data assimilation system. Local spatial filtering of the tensor was critical with small ensemble sizes (order 10), but the raw ensemble with 100 members gave good results without spatial filtering in our example. In general, a carefully designed objective filter would be beneficial in order to maximize the signal-tonoise ratio of the ensemble-estimates of the tensor elements in a similar way that it has been shown to be beneficial to the ensemble estimation of background-error variances (Raynaud et al., 2009; Berre and Desroziers, 2010).

In realistic applications, the numerical stability condition associated with explicit diffusion schemes can severely limit their computational efficiency. In particular, many iterations are likely to be needed with general anisotropic and inhomogeneous diffusion models that employ ensembleestimated tensors. Implicit diffusion schemes are more robust but require solving a large linear system for which efficient methods that are well-suited to massively parallel machines are required. This important practical aspect of the problem was not addressed in this paper and should be the subject of further research.

#### Acknowledgements

Financial support from the French National Research Agency (ANR) COSINUS programme (VODA project, no. ANR-08-COSI-016), the RTRA STAE foundation (ADTAO project), the European Framework Programme 7 (COMBINE project, GA 226520), and the French LEFE-ASSIM programme is gratefully acknowledged. This work benefited from discussions with Loik Berre, Serge Gratton, Sébastien Massart, Olivier Pannekoucke and Andrea Piacentini. The anonymous reviewers and the Associate Editor Martin Leutbecher provided many helpful remarks for improving the presentation of the paper.

#### Appendix A

## The Daley tensor of the 2D implicit-diffusion kernels

In this appendix we show that the expression for the Daley tensor of the 2D implicit-diffusion kernels  $c_2^{w}(\tilde{r})$  (Eq. (55) with d = 2) is related to the diffusion tensor by Eq. (56). For clarity of notation the subscript and superscript of  $c_2^{w}(\tilde{r})$  will be dropped hereafter. The relationships between the Daley and diffusion tensors for the implicit-diffusion kernels in higher dimensions and for the Gaussian function (Eqs (49)–(51)) are straightforward to verify following the basic procedure outlined here.

From the chain rule, the three independent elements of the outer product in the first equation of (56) can be written as

$$\frac{\partial^{2}c}{\partial x^{2}} = \frac{\partial^{2}c}{\partial \widetilde{r}^{2}} \left(\frac{\partial \widetilde{r}}{\partial x}\right)^{2} + \frac{\partial c}{\partial \widetilde{r}} \left(\frac{\partial^{2} \widetilde{r}}{\partial x^{2}}\right), \\
\frac{\partial^{2}c}{\partial y^{2}} = \frac{\partial^{2}c}{\partial \widetilde{r}^{2}} \left(\frac{\partial \widetilde{r}}{\partial y}\right)^{2} + \frac{\partial c}{\partial \widetilde{r}} \left(\frac{\partial^{2} \widetilde{r}}{\partial y^{2}}\right), \\
\frac{\partial^{2}c}{\partial x \partial y} = \frac{\partial^{2}c}{\partial \widetilde{r}^{2}} \left(\frac{\partial \widetilde{r}}{\partial x}\right) \left(\frac{\partial \widetilde{r}}{\partial y}\right) + \frac{\partial c}{\partial \widetilde{r}} \left(\frac{\partial^{2} \widetilde{r}}{\partial x \partial y}\right).$$
(68)

Expressing Eq. (55) with d = 2 as

$$c(\widetilde{r}) = \alpha_M \widetilde{r}^{M-1} K_{M-1}(\widetilde{r}), \qquad (69)$$

where  $\alpha_M = 2^{2-M}/(M-2)!$  and M > 2, and using the following recurrence relation for the modified Bessel functions of the second kind of integer order *n* (Eq. 9.6.26 of Abramowitz and Stegun, 1970),

$$\frac{\partial K_n}{\partial \widetilde{r}} = -\frac{n}{\widetilde{r}} K_n - K_{n-1},$$

where  $K_n = K_n(\tilde{r})$ , allows us to write

$$\frac{\partial c}{\partial \widetilde{r}} = -\alpha_M \widetilde{r}^{M-1} K_{M-2}, 
\frac{\partial^2 c}{\partial \widetilde{r}^2} = -\alpha_M \left( \widetilde{r}^{M-2} K_{M-2} - \widetilde{r}^{M-1} K_{M-3} \right).$$
(70)

The inverse of the symmetric diffusion tensor (46) can be written as

$$\boldsymbol{\kappa}^{-1} = \begin{pmatrix} \kappa_{xx}^{-1} & -\tau \kappa_{xy}^{-1} \\ -\tau \kappa_{xy}^{-1} & \kappa_{yy}^{-1} \end{pmatrix},$$

where  $\tau = 1/(\mu - 1)$  and  $\mu = \kappa_{xx}\kappa_{yy}/\kappa_{xy}^2$ . In expanded form the nondimensional distance measure (48) then reads

$$\widetilde{r}^{2} = \kappa_{xx}^{-1} (x - x')^{2} + \kappa_{yy}^{-1} (y - y')^{2} - 2\tau \kappa_{xy}^{-1} (x - x') (y - y').$$
(71)

From Eq. (71) we can derive the following relations

$$\frac{\partial \widetilde{r}}{\partial x} = \widetilde{r}^{-1}X,$$

$$\frac{\partial \widetilde{r}}{\partial y} = \widetilde{r}^{-1}Y,$$

$$\frac{\partial^{2} \widetilde{r}}{\partial x^{2}} = -\widetilde{r}^{-3}X^{2} + \widetilde{r}^{-1}\kappa_{xx}^{-1},$$

$$\frac{\partial^{2} \widetilde{r}}{\partial y^{2}} = -\widetilde{r}^{-3}Y^{2} + \widetilde{r}^{-1}\kappa_{yy}^{-1},$$

$$\frac{\partial^{2} \widetilde{r}}{\partial x \partial y} = -\widetilde{r}^{-3}XY - \widetilde{r}^{-1}\tau\kappa_{xy}^{-1},$$
(72)

<sup>&</sup>lt;sup>†</sup>It is common in VDA to use multivariate balance operators to transform the background-error variables into a set of new variables whose cross-covariances can be effectively neglected. The background-error covariances of the transformed variables can then be treated as univariate functions and represented via a diffusion model.

Copyright © 2012 Royal Meteorological Society

where

$$\begin{split} X &= \kappa_{xx}^{-1} \left( x - x' \right) \, - \, \tau \kappa_{xy}^{-1} \left( y - y' \right), \\ Y &= \kappa_{yy}^{-1} \left( y - y' \right) \, - \, \tau \kappa_{xy}^{-1} \left( x - x' \right). \end{split}$$

Substituting Eqs (70) and (72) in Eq. (68) yields

$$\frac{\partial^2 c}{\partial x^2} = \alpha_M \left( \widetilde{r}^{M-3} K_{M-3} X^2 - \widetilde{r}^{M-2} K_{M-2} \kappa_{xx}^{-1} \right),$$
  

$$\frac{\partial^2 c}{\partial y^2} = \alpha_M \left( \widetilde{r}^{M-3} K_{M-3} Y^2 - \widetilde{r}^{M-2} K_{M-2} \kappa_{yy}^{-1} \right),$$
  

$$\frac{\partial^2 c}{\partial x \partial y} = \alpha_M \left( \widetilde{r}^{M-3} K_{M-3} X Y + \widetilde{r}^{M-2} K_{M-2} \tau \kappa_{xy}^{-1} \right).$$

Since c(0) = 1, for all allowable *M*, we have from Eq. (69) the general relation

$$\widetilde{r}^n K_n(\widetilde{r})\Big|_{\widetilde{r}=0} = \frac{1}{\alpha_{n+1}}.$$

Hence

$$\frac{\partial^2 c}{\partial x^2}\Big|_{\widetilde{r}=0} = \frac{\alpha_M}{\alpha_{M-2}} X^2 \Big|_{\widetilde{r}=0} - \frac{\alpha_M}{\alpha_{M-1}} \kappa_{xx}^{-1} \Big|_{\widetilde{r}=0}, \quad (73)$$

$$\frac{\partial^2 c}{\partial y^2}\Big|_{\widetilde{r}=0} = \frac{\alpha_M}{\alpha_{M-2}} Y^2 \Big|_{\widetilde{r}=0} - \frac{\alpha_M}{\alpha_{M-1}} \kappa_{yy}^{-1} \Big|_{\widetilde{r}=0}, \quad (74)$$

$$\frac{\partial^2 c}{\partial x \partial y}\Big|_{\widetilde{r}=0} = \frac{\alpha_M}{\alpha_{M-2}} X Y \Big|_{\widetilde{r}=0} + \frac{\alpha_M}{\alpha_{M-1}} \tau \kappa_{xy}^{-1} \Big|_{\widetilde{r}=0}.$$

The first term on the right-hand side of each of the above equations vanishes since X = Y = 0 at  $\tilde{r} = 0$ , while the common coefficient of the second term in each equation is

$$\frac{\alpha_M}{\alpha_{M-1}} = \frac{1}{2M-4}.$$

Thus we obtain the relationship governed by (56) with d = 2.

In the isotropic case,  $\kappa_{xx} = \kappa_{yy} = L^2$  and  $\tau \kappa_{xy}^{-1} = 0$ . Equations (73) and (74) can then be averaged and inverted to yield the standard definition of the square of the 2D Daley length-scale involving the Laplacian operator (Eq. (44) with d = 2).

### Appendix B

## Estimating the Hessian tensor from an ensemble of simulated errors

In this appendix we provide a derivation of Eq. (62) for the special case of a homogeneous correlation function. The derivation is similar to the one given in the Appendix of Belo Pereira and Berre (2006) except for notational changes and greater emphasis here on some of the underlying assumptions.

The starting point is the general expression (59) for the covariance function  $B(\mathbf{x}, \mathbf{x}')$  of the ensemble of model-state errors. We consider here the 2D case where  $\mathbf{x} = (x, y)^{\mathrm{T}}$  and  $\mathbf{x}' = (x', y')^{\mathrm{T}}$ , and assume that  $B(\mathbf{x}, \mathbf{x}')$  is at least twice differentiable. We can express the covariance function of

Copyright © 2012 Royal Meteorological Society

the derivatives of the ensemble errors as follows (Swerling, 1962; Daley, 1991, p. 156):

$$B_{xx'}(\mathbf{x}, \mathbf{x}') = E\left[\left(\frac{\partial \epsilon(\mathbf{x})}{\partial x}\right) \left(\frac{\partial \epsilon(\mathbf{x}')}{\partial x'}\right)\right] = \frac{\partial^2 B(\mathbf{x}, \mathbf{x}')}{\partial x \, \partial x'},$$
  

$$B_{yy'}(\mathbf{x}, \mathbf{x}') = E\left[\left(\frac{\partial \epsilon(\mathbf{x})}{\partial y}\right) \left(\frac{\partial \epsilon(\mathbf{x}')}{\partial y'}\right)\right] = \frac{\partial^2 B(\mathbf{x}, \mathbf{x}')}{\partial y \, \partial y'},$$
  

$$B_{xy'}(\mathbf{x}, \mathbf{x}') = E\left[\left(\frac{\partial \epsilon(\mathbf{x})}{\partial x}\right) \left(\frac{\partial \epsilon(\mathbf{x}')}{\partial y'}\right)\right] = \frac{\partial^2 B(\mathbf{x}, \mathbf{x}')}{\partial x \, \partial y'},$$
  

$$B_{yx'}(\mathbf{x}, \mathbf{x}') = E\left[\left(\frac{\partial \epsilon(\mathbf{x})}{\partial y}\right) \left(\frac{\partial \epsilon(\mathbf{x}')}{\partial x'}\right)\right] = \frac{\partial^2 B(\mathbf{x}, \mathbf{x}')}{\partial y \, \partial x'}.$$

Using Eq. (60) the derivatives on the right-hand side of the above equations can be evaluated in terms of the standard deviations  $\sigma(\mathbf{x})$  and correlation function  $C(\mathbf{x}, \mathbf{x}')$ . Focusing on the first of these equations this yields

$$B_{xx'}(\mathbf{x}, \mathbf{x}') = \left(\frac{\partial \sigma(\mathbf{x})}{\partial x}\right) \left(\frac{\partial \sigma(\mathbf{x}')}{\partial x'}\right) C(\mathbf{x}, \mathbf{x}') + \sigma(\mathbf{x}) \left(\frac{\partial \sigma(\mathbf{x}')}{\partial x'}\right) \left(\frac{\partial C(\mathbf{x}, \mathbf{x}')}{\partial x}\right) + \sigma(\mathbf{x}') \left(\frac{\partial \sigma(\mathbf{x})}{\partial x}\right) \left(\frac{\partial C(\mathbf{x}, \mathbf{x}')}{\partial x'}\right) + \sigma(\mathbf{x}) \sigma(\mathbf{x}') \frac{\partial^2 C(\mathbf{x}, \mathbf{x}')}{\partial x \partial x'}.$$

Under the assumption of homogeneous correlations, we can write  $C(\mathbf{x}, \mathbf{x}') = c(\mathbf{r})$  where  $\mathbf{r} = \mathbf{x} - \mathbf{x}' = (x - x', y - y')^{\mathrm{T}}$  (Gaspari and Cohn, 1999). Using the chain rule, the derivatives of *C* with respect to *x*, *x'*, *y* and *y'* can be rewritten in terms of derivatives of *c* with respect to  $\tilde{x} = x - x'$  and  $\tilde{y} = y - y'$ . For  $B_{xx'}(\mathbf{x}, \mathbf{x}')$  this gives

$$B_{xx'}(\mathbf{x}, \mathbf{x}') = \left(\frac{\partial \sigma(\mathbf{x})}{\partial x}\right) \left(\frac{\partial \sigma(\mathbf{x}')}{\partial x'}\right) c(\mathbf{r}) + \sigma(\mathbf{x}) \left(\frac{\partial \sigma(\mathbf{x}')}{\partial x'}\right) \left(\frac{\partial c(\mathbf{r})}{\partial \widetilde{x}}\right) - \sigma(\mathbf{x}') \left(\frac{\partial \sigma(\mathbf{x})}{\partial x}\right) \left(\frac{\partial c(\mathbf{r})}{\partial \widetilde{x}}\right) - \sigma(\mathbf{x}) \sigma(\mathbf{x}') \frac{\partial^2 c(\mathbf{r})}{\partial \widetilde{x}^2}.$$

Evaluating the above equation at  $\mathbf{x} = \mathbf{x}'$  ( $\mathbf{r} = \mathbf{0}$ ), and noting that  $\partial c(\mathbf{r}) / \partial \widetilde{x}|_{\mathbf{r}=\mathbf{0}} = 0$  since  $c(\mathbf{r})$  is maximum at  $\mathbf{r} = \mathbf{0}$  and that  $c(\mathbf{0}) = 1$ , yields

$$B_{xx}(\mathbf{x},\mathbf{x}) = \left(\frac{\partial \sigma(\mathbf{x})}{\partial x}\right)^2 - \left(\sigma(\mathbf{x})\right)^2 \left.\frac{\partial^2 c(\mathbf{r})}{\partial \tilde{x}^2}\right|_{\mathbf{r}=\mathbf{0}}$$

which can be rearranged as

$$-\frac{\partial^2 c(\mathbf{r})}{\partial \tilde{x}^2}\Big|_{\mathbf{r}=\mathbf{0}} = \frac{1}{\sigma^2} \left( B_{xx}(\mathbf{x}, \mathbf{x}) - \left(\frac{\partial \sigma}{\partial x}\right)^2 \right), \quad (75)$$

where  $\sigma = \sigma(\mathbf{x})$ . A similar analysis for the other covariance functions yields

$$-\frac{\partial^2 c(\mathbf{r})}{\partial \tilde{\gamma}^2}\Big|_{\mathbf{r}=\mathbf{0}} = \frac{1}{\sigma^2} \left( B_{yy}(\mathbf{x}, \mathbf{x}) - \left(\frac{\partial \sigma}{\partial y}\right)^2 \right),\tag{76}$$

$$-\frac{\partial^2 c(\mathbf{r})}{\partial \widetilde{\mathbf{x}} \, \partial \widetilde{\mathbf{y}}} \bigg|_{\mathbf{r}=\mathbf{0}} = \frac{1}{\sigma^2} \bigg( B_{xy}(\mathbf{x}, \mathbf{x}) - \left(\frac{\partial \sigma}{\partial x}\right) \left(\frac{\partial \sigma}{\partial y}\right) \bigg), \qquad (77)$$

$$-\frac{\partial^2 c(\mathbf{r})}{\partial \widetilde{y} \, \partial \widetilde{x}} \bigg|_{\mathbf{r}=\mathbf{0}} = \frac{1}{\sigma^2} \bigg( B_{yx}(\mathbf{x}, \mathbf{x}) - \left(\frac{\partial \sigma}{\partial y}\right) \left(\frac{\partial \sigma}{\partial x}\right) \bigg).$$
(78)

In tensor notation, the left-hand side of Eqs (75)–(78) is equivalent to Eq. (61) evaluated at  $\mathbf{r} = \mathbf{0}$ , while the right-hand side of the equations can be identified with the right-hand side of Eq. (62).

#### References

- Balmaseda MA, Vidard A, Anderson DLT. 2008. The ECMWF Ocean Analysis System: ORA-S3. *Mon. Weather Rev.* **136**: 3018–3034.
- Bannister RN. 2008. A review of forecast error covariance statistics in atmospheric variational data assimilation. II: Modelling the forecast error covariance statistics. Q. J. R. Meteorol. Soc. 134: 1971–1996.
- Belo Pereira M, Berre L. 2006. The use of an ensemble approach to study the background error covariances in a global NWP model. *Mon. Weather Rev.* **134**: 2466–2489.
- Bennett AF, Chua BS, Leslie LM. 1996. Generalized inversion of a global numerical weather prediction model. *Meteorol. Atmos. Phys.* 60: 165–178.
- Bennett AF, Chua BS, Leslie LM. 1997. Generalized inversion of a global numerical weather prediction model. II: analysis and implementation. *Meteorol. Atmos. Phys.* 62: 129–140.
- Berre L, Desroziers G. 2010. Filtering of background error variances and correlations by local spatial averaging: a review. *Mon. Weather Rev.* 138: 3693–3720.
- Carrier MJ, Ngodock H. 2010. Background-error correlation model based on the implicit solution of a diffusion equation. *Ocean Model*. 35: 45–53.
- Chorti A, Hristopulos DT. 2008. Non-parametric identification of anisotropic (elliptic) correlations in spatially distributed data sets. *IEEE Trans. Signal Process.* **56**: 4738–4751.
- Chua BS, Bennett AF. 2001. An inverse ocean modeling system. Ocean Model. 3: 137–165.
- Daget N, Weaver AT, Balmaseda MA. 2009. Ensemble estimation of background-error variances in a three-dimensional variational data assimilation system for the global ocean. Q. J. R. Meteorol Soc. 135: 1071–1094.
- Daley R. 1991. Atmospheric Data Analysis. Cambridge University Press: Cambridge, UK.
- De Pondeca MSFV, Purser RJ, Parrish DF, Derber J. 2006. 'Comparison of strategies for the specification of anisotropies in the covariances of a three-dimensional atmospheric data assimilation system'. Office Note 452, National Centers for Environmental Prediction: Camp Springs, MD.
- Derber JC, Rosati A. 1989. A global oceanic data assimilation system. J. Phys. Oceanogr. 19: 1333–1347.
- Di Lorenzo E, Moore AM, Arango HG, Cornuelle BD, Miller AJ, Powell B, Chua BS, Bennett AF. 2007. Weak and strong constraint data assimilation in the inverse Regional Ocean Modeling System (ROMS): Development and application for a baroclinic coastal upwelling system. *Ocean Model.* **16**: 160–187.
- Dobricic S, Pinardi N. 2008. An oceanographic three-dimensional variational data assimilation scheme. *Ocean Model.* **22**: 89–105.
- Egbert G, Bennett A, Foreman M. 1994. Topex/Poseidon tides estimated using a global inverse model. *J. Geophys. Res.* **99**: 24821–24852.
- Elbern H, Schwinger J, Botchorishvili R. 2010. Chemical state estimation for the middle atmosphere by four-dimensional variational data assimilation: system configuration. J. Geophys. Res. 115: D06302, DOI: 10.1029/2009 JD011953.
- Gaspari G, Cohn S. 1999. Construction of correlation functions in two and three dimensions. Q. J. R. Meteorol. Soc. 125: 723–757.
- Geer AJ, Lahoz WA, Bekki S, Bormann N, Errera Q, Eskes HJ, Fonteyn D, Jackson DR, Juckes MN, Massart S, Peuch V-H, Rharmili S, Segers A. 2006. The ASSET intercomparison of ozone analyses: method and first results. *Atmos. Chem. Phys.* 6: 5445–5474.

- Gneiting T. 1999. Correlation functions for atmospheric data analysis. *Q. J. R. Meteorol. Soc.* **125**: 2449–2464.
- Gneiting T, Kleiber W, Schlather M. 2009. 'Matérn cross-covariance functions for multivariate random fields'. Technical Report No. 549, University of Washington, Department of Statistics. http://www.stat.washington.edu/research/reports/2009/tr549.pdf
- Gregori P, Porcu E, Mateu J, Sasvári Z. 2008. On potentially negative space time covariances obtained as sum of products of marginal ones. *Ann. Inst. Stat. Math.* **60**: 865–882.
- Guttorp P, Gneiting T. 2006. Miscellanea studies in the history of probability and statistics XLIX: on the Matérn correlation family. *Biometrika* 93: 989–995.
- Haidvogel DB, Beckmann A. 1999. Numerical Ocean Circulation Modeling. Imperial College Press: London.
- Hartman P, Watson GS. 1974. 'Normal' distribution functions on spheres and the modified Bessel functions. Ann. Prob. 2: 593–607.
- Hayden CM, Purser RJ. 1995. Recursive filter objective analysis of meteorological fields: applications to NESDIS operational processing. *J. Appl. Meteorol.* 34: 3–15.
- Hristopulos DT. 2002. New anisotropic covariance models and estimation of aniostropic parameters based on the covariance tensor identity. *Stoch. Environ. Res. Risk Assess.* **16**: 43–62.
- Hristopulos DT. 2003. Spartan Gibbs random field models for geostatistical applications. SIAM J. Sci. Comput. 24: 2125–2162.
- Hristopulos DT, Elogne SN. 2007. Analytic properties and covariance functions of a new class of generalized Gibbs random fields. *IEEE Trans. Inform. Theory.* **53**: 4467–4679.
- Kurapov AL, Egbert GD, Allen JS, Miller RN. 2009. Representer-based analyses in the coastal upwelling system. Dyn. Atmos. Oceans 48: 198–218.
- Liu H, Xue M, Purser RJ, Parrish DF. 2007. Retrieval of moisture from simulated GPS slant-path water vapor observations using 3DVAR with anisotropic recursive filter. *Mon. Weather Rev.* 135: 1506–1521.
- Liu Y, Zhu J, She J, Zhuang S, Fu W, Gao J. 2009. Assimilating temperature and salinity profile observations using an anisotropic recursive filter in a coastal ocean model. *Ocean Model.* **30**: 75–87.
- Lorenc AC. 1992. Iterative analysis using covariance functions and filters. *Q. J. R. Meteorol. Soc.* **118**: 569–591.
- Lorenc AC, Ballard SP, Bell RS, Ingleby NB, Andrews PLF, Barker DM, Bray JR, Clayton AM, Dalby T, Li D, Payne TJ, Saunders FW. 2000. The Met Office global three-dimensional variational data assimilation system. *Q. J. R. Meteorol. Soc.* **126**: 2991–3012.
- Martin MJ, Hines A, Bell MJ. 2007. Data assimilation in the FOAM operational short-range ocean forecasting system: a description of the scheme and its impact. *Q. J. R. Meteorol. Soc.* **133**: 981–995.
- Massart S, Piacentini P, Pannekouck O. 2012. Importance of using ensemble estimated background error covariances for the quality of atmospheric ozone analyses. Q. J. R. Meteorol. Soc., DOI: 10.1002/qj.971.
- McIntosh P. 1990. Oceanographic data interpolation: objective analysis and splines. J. Geophys. Res. 95: 13529–13541.
- Mirouze I, Weaver AT. 2010. Representation of correlation functions in variational assimilation using an implicit diffusion operator. *Q. J. R. Meteorol. Soc.* **136**: 1421–1443.
- Moore AM, Arango HG, Broquet G, Powell BS, Weaver AT, Zavala-Garay J. 2011. The Regional Ocean Modeling System (ROMS) 4dimensional variational data assimilation systems. I: System overview and formulation. *Prog. Oceanogr.* **91**: 50–73.
- Muccino JC, Arango HG, Bennett AF, Chua BS, Cornuelle BD, Di Lorenzo E, Egbert GD, Haidvogel D, Levin JC, Luo H, Miller AJ, Moore AM, Zaron ED. 2008. The Inverse Ocean Modelling System. Part II: applications. J. Atmos. Ocean. Tech. 25: 1623–1637.
- Ngodock HE. 2005. Efficient implementation of covariance multiplication for data assimilation with the representer method. *Ocean Model.* **8**: 237–251.
- Paciorek CJ, Schervish MJ. 2006. Spatial modelling using a new class of nonstationary covariance functions. *Environmetrics* 17: 483–506.
- Pannekoucke O. 2009. Heterogeneous correlation modeling based on the wavelet diagonal assumption and on the diffusion operator. *Mon. Weather Rev.* 131: 2995–3012.
- Pannekoucke O, Massart S. 2008. Estimation of the local diffusion tensor and normalization for heterogeneous correlation modelling using a diffusion equation. Q. J. R. Meteorol. Soc. 134: 1425–1438.
- Pannekoucke O, Berre L, Desroziers G. 2008. Background-error correlation length-scale estimates and their sampling statistics. Q. J. R. Meteorol. Soc. 134: 497–508.
  Purser RJ. 2005. 'A geometrical approach to the synthesis of smooth statemetrical approach to the synthesis of smooth statemetrical approach.
- Purser RJ. 2005. 'A geometrical approach to the synthesis of smooth anisotropic covariance operators for data assimilation'. Office Note 447, National Centers for Environmental Prediction: Camp Springs, MD.

Copyright © 2012 Royal Meteorological Society

- Purser RJ. 2008a. 'Normalization of the diffusive filters that represent the inhomogeneous covariance operators of variational assimilation, using asymptotic expansions and techniques of non-Euclidean geometry. Part I: Analytic solutions for symmetrical configurations and the validation of practical algorithms'. Office Note 456, National Centers for Environmental Prediction: Camp Springs, MD.
- Purser RJ. 2008b. 'Normalization of the diffusive filters that represent the inhomogeneous covariance operators of variational assimilation, using asymptotic expansions and techniques of non-Euclidean geometry. Part II: Riemannian geometry and the generic parametrix expansion method'. Office Note 457, National Centers for Environmental Prediction: Camp Springs, MD.
- Purser RJ, Wu WS, Parrish DF, Roberts NM. 2003a. Numerical aspects of the application of recursive filters to variational statistical analysis. Part I: Spatially homogeneous and isotropic Gaussian covariances. *Mon. Weather Rev.* 131: 1524–1535.
- Purser RJ, Wu WS, Parrish DF, Roberts NM. 2003b. Numerical aspects of the application of recursive filters to variational statistical analysis. Part II: Spatially inhomogeneous and anisotropic general covariances. *Mon. Weather Rev.* 131: 1536–1548.
- Raynaud L, Pannekoucke O. 2012. Heterogeneous filtering of ensemble-based background-error variances. Q. J. R. Meteorol. Soc., DOI:10.1002/qj.1890.
- Raynaud L, Berre L, Desroziers G. 2009. Objective filtering of ensemblebased background-error variances. Q. J. R. Meteorol. Soc. 134: 1003–1014.
- Roberts PH, Ursell HD. 1960. Random walk on a sphere and on a Riemannian manifold. *Philos. Trans. R. Soc. London A* **252**: 317–356.
- Sato Y, De Pondeca MSFV, Purser RJ, Parrish DF. 2009. 'Ensemble-based background error covariance implementations using spatial recursive filters in NCEP's grid-point statistical interpolation system'. Office Note 459, National Centers for Environmental Prediction: Camp Springs, MD.
- Sheinbaum J, Anderson DLT. 1990. Variational assimilation of XBT data. Part II: Sensitivity studies and use of smoothing constraints. J. Phys. Oceanogr. 20: 689–704.
- Stein ML. 1999. Interpolation of Spatial Data: Some Theory for Kriging. Springer: New-York.
- Swerling P. 1962. Statistical properties of the contours of random surfaces. IRE Trans. Inform. Theory IT-8: 315–321.
- Thacker WC. 1988. Fitting models to inadequate data by enforcing spatial and temporal smoothness. J. Geophys. Res. 93: 10655–10665.
- Wahba G. 1982. 'Vector splines on the sphere, with application to the estimation of vorticity and divergence from discrete, noisy data'. In

*Multivariate Approximation Theory II*, Schempp W, Zeller K (eds). Birkhäuser: Basel; 407–429.

- Wahba G, Wendelberger J. 1980. Some new mathematical methods for variational objective analysis using splines and cross-validation. *Mon. Weather Rev.* **108**: 36–57.
- Wang X, Barker DM, Synder C, Hamill TM. 2008. A hybrid ETKF-3DVAR data assimilation scheme for the WRF model. Part I: Observing system simulation experiment. *Mon. Weather Rev.* **136**: 5116–5131.
- Weaver AT, Courtier P. 2001. Correlation modelling on the sphere using a generalized diffusion equation. Q. J. R. Meteorol. Soc. 127: 1815–1846.
- Weaver AT, Ricci S. 2004. 'Constructing a background-error correlation model using generalized diffusion operators'. In ECMWF Proceedings of the Seminar on Recent Developments in Data Assimilation for Atmosphere and Ocean. ECMWF: Reading, UK; 327–339.
- Weaver AT, Vialard J, Anderson DLT. 2003. Three- and fourdimensional variational assimilation with an ocean general circulation model of the tropical Pacific Ocean. Part 1: Formulation, internal diagnostics and consistency checks. *Mon. Weather Rev.* 131: 1360–1378.
- Weber RO, Talkner P. 1993. Some remarks on spatial correlation function models. *Mon. Weather Rev.* **121**: 2611–2617.
- Whittle P. 1954. On stationary processes in the plane. *Biometrika* 41: 434-449.
- Whittle P. 1963. Stochastic processes in several dimensions. Bull. Inst. Int. Statist. 40: 974-994.
- Wu WS, Purser RJ, Parrish DF. 2002. Three-dimensional variational analysis with spatially inhomogeneous covariances. *Mon. Weather Rev.* **130**: 2905–2916.
- Xu Q. 2005. Representations of inverse covariances by differential operators. Adv. Atmos. Sci. 22: 181–198.
- Yaglom AM. 1987. Correlation Theory of Stationary and Related Random Functions. I: Basic Results. Springer: New York.
- Yaremchuk M, Carrier M. 2012. On the renormalization of the covariance operators. *Mon. Weather Rev.* 140: 637–649.
- Yaremchuk M, Smith S. 2011. On the correlation functions associated with polynomials of the diffusion operator. Q. J. R. Meteorol. Soc. 137: 1927–1932.
- Yaremchuk M, Nechaev D, Pan C. 2011. A hybrid background-error covariance model for assimilating glider data into a coastal ocean model. *Mon. Weather Rev.* 139: 1879–1890.
- Zaron ED, Chavanne C, Egbert GD, Flament P. 2009. Baroclinic tidal generation in the Kauai Channel inferred from high-frequency radio Doppler current meters. *Dyn. Atmos. Oceans* **48**: 93–120.