



Uncertainty-aware surrogate modeling for urban air pollutant dispersion prediction

Eliott Lumet^{a,*}, Mélanie C. Rochoux^a, Thomas Jaravel^a, Simon Lacroix^b

^a CECI, Université de Toulouse, CNRS, CERFACS, 42 Avenue Gaspard Coriolis, 31057 Toulouse cedex 1, France

^b LAAS-CNRS, Université de Toulouse, CNRS, 7 Avenue du Colonel Roche, BP54200, 31031 Toulouse cedex 4, France

ARTICLE INFO

Dataset link: PPMLES, https://github.com/eliott-lumet/pod_gpr_ppmles

Keywords:

Surrogate modeling
Uncertainty quantification
Microscale pollutant dispersion
Urban flow
Large-eddy simulation
Internal variability

ABSTRACT

This study evaluates a surrogate modeling approach that provides rapid ensemble predictions of air pollutant dispersion in urban environments for varying meteorological forcing, while estimating irreducible and modeling uncertainties. The POD–GPR approach combining Proper Orthogonal Decomposition (POD) and Gaussian Process Regression (GPR) is applied to emulate the response surface of a Large-Eddy Simulation (LES) model of the Mock Urban Setting Test (MUST) field-scale experiment. We design and validate new methods for (i) selecting the POD-latent space dimension to avoid overfitting noisy structures due to atmospheric internal variability, and (ii) estimating the uncertainty in POD–GPR predictions. To train and validate the POD–GPR surrogate in an offline phase, we build a large dataset of 200 LES 3-D time-averaged concentration fields, which are subject to substantial spatial variability from near-source to background concentration and have a very large dimension of several million grid cells. The results show that POD–GPR reaches the best achievable accuracy levels, except for the highest concentration near the source, while predicting full fields at a computational cost five orders of magnitude lower than an LES. The results also show that the proposed mode selection criterion avoids perturbing the surrogate response surface, and that the uncertainty estimate explains a large part of the surrogate error and is spatially consistent with the observed internal variability. Finally, POD–GPR can be robustly trained with much smaller datasets, paving the way for application to realistic urban configurations.

1. Introduction

Accidental releases of pollutants into the atmosphere, such as from industrial accidents, can degrade air quality and have significant short- and long-term health impacts [1,2]. In urban environments, these risks are exacerbated by high population density and reduced ventilation due to the urban canopy, leading to local pollution peaks [3–5]. For accurate mapping of these peaks and associated exposures, it is necessary to develop microscale dispersion models that take into account (i) the effect of urban buildings on the local flow, and (ii) the inherently multiscale and turbulent nature of the Atmospheric Boundary Layer (ABL).

To gain relevant insight into these processes, there is a growing consensus in the research community for the use of Computational Fluid Dynamics (CFD) [6,7]. Advanced models based on Reynolds-Averaged Navier–Stokes (RANS) and Large-Eddy Simulation (LES) are able to represent complex flow structures, in particular due to the interactions between the atmosphere and the built environment. However, their use in operational applications remains limited because their high computational cost prevents them from being used in real time, for

example in emergency response. Moreover, they still suffer from a lack of accuracy compared to field and wind tunnel measurements due to the large uncertainties involved [8–10]. These uncertainties can be classified into three different groups:

- *boundary condition uncertainties* due to measurement and representativeness errors in calibration data, and to boundary condition modeling assumptions, in relation to: (i) the meteorological forcing [11–13], (ii) the urban geometry representation [14–16], and (iii) the pollutant source [17,18];
- *structural modeling uncertainties*, inherent to the model solver and its underlying modeling assumptions, mainly related to turbulence modeling [19–25];
- *aleatory uncertainties*, mostly due to the turbulent and therefore stochastic nature of the ABL, and referred to as internal variability, which results in an irreducible uncertainty and is largely responsible for the discrepancies between field measurements and CFD model predictions [8,10,26–29].

In this work, we focus on atmospheric uncertainties, i.e. in how to represent the impact of large-scale atmospheric forcing uncertainties

* Corresponding author.

E-mail address: eliott.lumet@cerfacs.fr (E. Lumet).

and internal variability on microscale LES field predictions. We have chosen not to consider structural modeling uncertainties, as these have been extensively studied and remain small in the LES context. Instead, we have chosen to investigate how to design a surrogate modeling approach to quantify boundary condition uncertainties in LES, while accounting for internal variability. To our knowledge, the coupling between these two sources of uncertainty has not yet been studied, while this is one challenge expressed by Dauxois et al. [10] and Wu and Quan [30].

Surrogate modeling, also known as reduced-order modeling, aims at accurately emulating the response surface of complex and expensive numerical models, while significantly reducing computational time. By enabling real-time and large ensemble predictions, surrogate modeling is well suited to address the dual challenges of high cost and uncertainty in LES models, making it a hot topic of research in the CFD field [31,32]. For parametric studies, surrogate models are mostly based on fully data-driven approaches, which consist of learning the response surface of the CFD model from a dataset of reference simulations precomputed during an offline phase, to then provide fast predictions during an online phase. They have been successfully used to emulate urban wind and/or pollutant dispersion, with respect to urban geometry [33–36], or meteorological forcing and pollutant source [37–39]. Surrogate models are therefore valuable for ensemble prediction in more complex frameworks such as urban design optimization [30,33], sensitivity analysis [40,41], uncertainty quantification [11,42], and data assimilation [43–46].

While surrogate models have proven to be valuable tools for dealing with uncertainties related to CFD model boundary conditions, few studies have addressed the representation of internal variability, which is at least as important [26,27,29]. Moreover, surrogate models introduce a new form of structural uncertainty: the model reduction error, i.e. the error of the surrogate model relative to the full-order model. Our aim is to evaluate the model reduction error in a comprehensive and robust way, and to assess the ability of the surrogate model to retrieve reliable information on internal variability from the LES dataset and compare it with the model reduction error.

To this end, we adopt a surrogate modeling approach called POD–GPR [47], which combines Proper Orthogonal Decomposition [48,49] and Gaussian Process Regression [50]. It is a robust and standard method that has already been used for urban wind and pollutant dispersion prediction [38,39,41,51–53]. In this study, we construct a POD–GPR model for the MUST experiment of propylene dispersion in a simplified urban canopy [54]. For this purpose, we generate a large dataset of 200 LES using the model validated by Lumet et al. [29] by varying the wind boundary forcing. We choose LES over the more common and less expensive RANS approach because: (i) LES is expected to reduce structural uncertainties due to turbulence modeling compared to RANS [28,55], and (ii) LES provides instantaneous snapshots of the most energetic atmospheric eddies and can thus be used to estimate the effect of the microscale internal variability of the ABL on tracer dispersion [29], which is central to the objective of this study.

The novelty of the proposed surrogate modeling approach is related to the POD latent space, i.e. the reduced space compressing the LES information, and is twofold. First, we define a method to choose a priori the POD-latent space dimension, based on the projection of the internal variability into the latent space. Secondly, knowing the internal variability in the LES data and using regression uncertainty estimates from Gaussian processes, we develop a mathematical framework for propagating these uncertainty estimates from the POD latent space to the physical space to help interpret the uncertainty results, which to our knowledge has been little studied in physical applications.

This article is structured as follows: Section 2 briefly introduces the learning dataset of LES simulations. Section 3 describes the POD–GPR surrogate modeling approach and introduces our methods to estimate prediction uncertainty and select the latent space dimension. Finally, Section 4 provides a comprehensive validation of the POD–GPR predictions, uncertainty estimates, and ability to handle reduced-size training datasets.

2. Learning dataset of large-eddy simulations

This section summarizes the key points of the LES model for the MUST experiment, which has been extensively validated in previous work [29] and which is used here to build the surrogate learning dataset. Details are given on the choice of the parameter space, the field quantities of interest and the associated internal variability.

2.1. The MUST field campaign

MUST is a field-scale experiment conducted in September 2001 at the US Army Dugway Proving Ground test site in Utah’s desert to collect extensive measurements of urban pollutant dispersion [54,56]. During the field campaign, a series of trials were carried out by releasing a passive tracer, propylene, at different locations within an urban-like canopy consisting of 120 regularly-spaced shipping containers. It is a canonical experiment for dispersion model validation: (i) it was selected as one of the reference case studies for the COST Action 732 CFD dispersion model intercomparison [57], and (ii) it has been used in a large number of CFD studies involving RANS [58–64] or LES [27,65–67]. In this study, we focus on the trial 2681829 corresponding to neutral atmospheric conditions.

2.2. LES model of the MUST field experiment

We use the AVBP¹ [68,69] code to build the LES model. AVBP solves the LES-filtered Navier–Stokes equations on unstructured mesh using a second-order Lax–Wendroff finite-volume centered numerical scheme [68] and using pressure gradient scaling since the atmospheric flow features a low Mach number [70]. Tracer dispersion is modeled by the LES-filtered advection-diffusion equation using an Eulerian approach. Subgrid-scale turbulence is modeled using the Wall-Adaptative Local Eddy-Viscosity (WALE) model [71] for subgrid momentum transport, and a gradient-diffusion hypothesis for subgrid tracer transport (with the turbulent Schmidt number equal to $S'_t = 0.6$).

The computational domain is a rectangular box with dimensions of 420 m by 420 m by 50 m, discretized with a boundary-fitted mesh of 91 million tetrahedra, with a resolution ranging from 0.3 m in the canopy to 5 m at the top of the domain.

In terms of boundary conditions, the wind velocity vector \mathbf{u} imposed at the inlet and expressed in the fixed coordinate system aligned with the containers defined by Yee and Biltoft [54] reads:

$$\mathbf{u} = \bar{\mathbf{u}} + \mathbf{u}' \quad \text{with} \quad \bar{\mathbf{u}}(z) = \frac{u_*}{\kappa} \ln \left(\frac{z + z_0}{z_0} \right) \begin{pmatrix} \cos(\alpha_{inlet}) \\ \sin(\alpha_{inlet}) \\ 0 \end{pmatrix}, \quad (1)$$

where the logarithmic wind profile used for the time-averaged part $\bar{\mathbf{u}}$ is set with κ the von Kármán constant equal to 0.4, z_0 the aerodynamic roughness length equal to 0.045 m [54], and u_* the friction velocity. The inlet wind velocity fluctuations \mathbf{u}' are prescribed using the synthetic turbulence injection method from Smirnov et al. [72], which allows to impose a level of wind speed fluctuations that is anisotropic and height-dependent. The prescribed Reynolds stress tensor is computed from a precursor run (corresponding to a simulation with the same surface roughness but without obstacles, and with periodic boundary conditions at the inlet and outlet inspired by Vasaturo et al. [73]) and is further described in Lumet et al. [29]. At the lateral boundaries, symmetry boundary conditions are used. Static pressure is imposed at the outlet and top boundaries. Standard laws of the wall are imposed for the ground and obstacle boundaries. The pollutant source is modeled by a local source term in the advection-diffusion equation to match the experimental volumetric flow rate. A comprehensive description of the boundary conditions is given in Lumet [46].

¹ AVBP documentation, see <https://www.cerfacs.fr/avbp7x/>

To be comparable to the MUST observational time series, we need to simulate a 200-s time sequence for each snapshot of the learning dataset. Before running this time sequence, we need to initialize each simulation until first- and second-order statistics of the flow and tracer variables reach a stationary state. For this initialization, a spin-up time $t_{spin-up}$ of 1.5 times the convective time scale is used:

$$t_{spin-up} = 1.5 \times \left(\frac{L}{U_{bulk}} \right) = 1.5 \times \frac{\kappa H L}{u_* \left[(H + z_0) \ln \left(\frac{H+z_0}{z_0} \right) - H \right]}, \quad (2)$$

with $L = 420$ m the domain length and $H = 50$ m the domain height. This spin-up time is specific to each snapshot as the bulk velocity U_{bulk} is an uncertain quantity (Section 2.3). Note that the average computational cost for a given simulation of 200 s is around 15,000 core hours, which motivates the development of a surrogate model to speed up predictions.

2.3. Definition of the input parameter space

2.3.1. Choice of input parameters

In this work, we focus on atmospheric parametric uncertainties. For the surrogate model to be useful, it must capture the dependence of the tracer dispersion on the most influential and uncertain atmospheric parameters of the LES model. In preliminary work (Lumet [46], Chapter III), we carried out one-at-a-time sensitivity analysis and showed that the inlet wind direction α_{inlet} and the friction velocity u_* are the two parameters that most significantly affect the LES mean concentration predictions. In particular, the aerodynamic roughness length z_0 is well identified in the MUST experiment (z_0 is equal to 0.045 ± 0.005 m according to observations, Yee and Biltoft [54]) and was found to have a negligible impact. For these reasons, we consider only two uncertain parameters:

$$\theta = (\alpha_{inlet}, u_*), \quad (3)$$

to define the input space of the surrogate model. Note that this choice is quite common in urban flow surrogate modeling [11,37,42]. Note also that, under neutral conditions, the mean concentration is inversely proportional to the friction velocity and the reduction problem could thus be simplified by predicting dimensionless quantities, as done by Sousa et al. [44] and Lamberti and Górlé [74]. This normalization was investigated in Lumet [46], Chapter IV, but we choose to present results with multiple input dimensions here for generalization purposes.

2.3.2. Parameter variation ranges

The surrogate model must cover a wide, but plausible and feasible, range of variation in the input parameters (Eq. (3)). Based on a microclimatology constructed using all available data from the closest micrometeorological station to the MUST site (Lumet [46], Chapter IV), all wind directions are likely to occur and more than 99% of the horizontal wind speed measurements at $z = 10$ m are below 12 m s^{-1} , which corresponds to a friction velocity u_* of 0.89 m s^{-1} and which is therefore chosen as the maximum friction velocity here. We limit the minimum friction velocity to 0.07 m s^{-1} , which corresponds to a wind speed of about 1 m s^{-1} at an altitude of 10 m, since we are interested in windy conditions. To reduce the number of LES, we also restrict the range of variation for the inlet wind direction to wind directions for which the plume crosses the array of containers. In the end, the input parameter space reads

$$\Omega_\theta = [-90^\circ, 30^\circ] \times [0.07 \text{ m s}^{-1}, 0.89 \text{ m s}^{-1}]. \quad (4)$$

2.3.3. Parameter space sampling

To sample the input parameter space (Eq. (4)), we use Halton's sequence (1964). As a low-discrepancy sequence, it samples the space uniformly and more efficiently than a purely random sequence for a limited number of samples, avoiding redundant sampling in the same areas and it is well adapted to a small number of parameters. Fig. 1 shows the location of the 200 samples thus obtained in the uncertain parameter space.

2.4. Generation of the LES dataset

We run an LES for each of the 200 input parameter samples (Fig. 1) to provide the learning dataset for the surrogate model. The main quantity of interest for the surrogate modeling approach is the 3-D mean (time-averaged) concentration field averaged over the 200-s analysis time period of the MUST experiment.

To generate this ensemble, the computational domain is rotated to align with the mean wind direction α_{inlet} to avoid inducing lateral confinement and numerical instabilities due to the shear-free boundary conditions at the domain sides. The spin-up time before collecting LES statistics is scaled by the friction velocity according to Eq. (2) to account for the slowing down of the flow establishment with decreasing u_* . Finally, the Reynolds stress tensor prescribed for the turbulent injection method is rescaled by u_*^2 following similarity theory.

The total cost of generating this LES ensemble is about 5.7 million core hours. Note that a subset of the most relevant data from these simulations, including all the data used in this study, is available in open access [76].

Fig. 2a shows the topology of the LES ensemble with the example of the mean concentration c at one specific location within the canopy (the green square in Fig. 2b, c corresponding to the tower B in the MUST experiment). The mean concentration increases linearly with decreasing friction velocity. The dependence on the wind direction is more complex with a concentration maximum obtained for $\alpha_{inlet} \approx 30^\circ$ and a rapid decay in both directions down to 0 ppm as the plume no longer crosses the probe location. The two examples of horizontal cuts of the LES mean concentration fields (Fig. 2b, c) obtained for two different wind conditions highlight the high spatial variability of the fields, especially within the plumes, which is a challenge for the surrogate modeling problem.

2.5. Noise in the learning dataset

Atmospheric flows are naturally unsteady with strong variations occurring over a wide range of frequencies corresponding to the time scales of the atmospheric eddies. When considering statistics over finite temporal periods, this internal variability yields sampling errors and is therefore a source of aleatory uncertainty, which is inherent to the physical system under study and thereby irreducible. For the MUST case, internal variability has a significant impact on the tracer concentration statistics when computed over the standard 200-s analysis period [29,77]. One of the challenges of this study is to build a surrogate model that explicitly estimates this uncertainty when emulating the mean concentration fields.

To quantify the effect of internal variability on the LES predictions, we use the stationary bootstrap approach from Lumet et al. [29], which relies on resampling of the sub-averages of the physical fields using the algorithm of Politis and Romano [78] and which involves a mean bootstrap block length to account for temporal correlation between sub-averages. This approach is applied separately for each snapshot in the dataset (Fig. 1) using 1,000 bootstrap replicates to estimate the internal variability.

Fig. 3 confirms that the internal microscale variability of the ABL significantly affects the LES learning dataset, with spatially-averaged relative standard deviations of up to just over 20% for a few samples of the LES dataset. Looking at the mean concentration fields, these deviations can be even larger locally, especially in areas of strong gradients or close to the source. We note in Fig. 3 that the noise induced by internal variability is not homogeneous in the input parameter space, as it increases as the friction velocity decreases. This is because as advection decreases, the temporal correlation of concentration increases, which increases the uncertainty of the mean over the 200-s analysis period (less independent information to estimate the mean). We also note that the noise decreases as α_{inlet} moves away from the median value of -30° , due to a zoning bias: the plume moves further outside the domain at the boundary angles (Eq. (4)), and there is

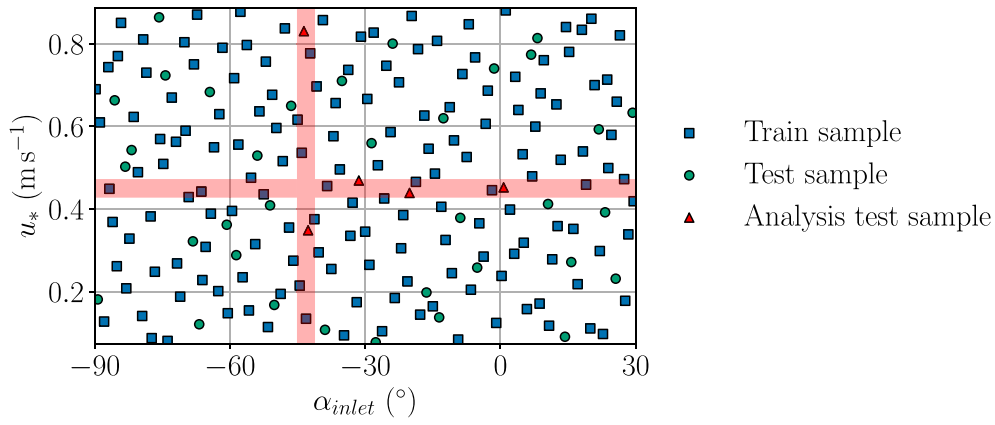


Fig. 1. Input parameter space sampling obtained with Halton's sequence. Each point is a pair of parameters for which we perform an LES prediction. The training (80%) and test (20%) sets are represented as blue squares and green circles, respectively. The horizontal red shaded area corresponds to the parameter space sub-section scanned by taking a margin of $\pm 5\%$ around the constant friction speed $u_*^{plat} = 0.45 \text{ m s}^{-1}$. The vertical shaded area is similarly defined around the constant inlet wind direction $\alpha_{inlet}^{plat} = -43^\circ$ with a margin of $\pm 2^\circ$. The test samples within these ranges (red triangles) are used in Section 4.3 to evaluate the surrogate model.

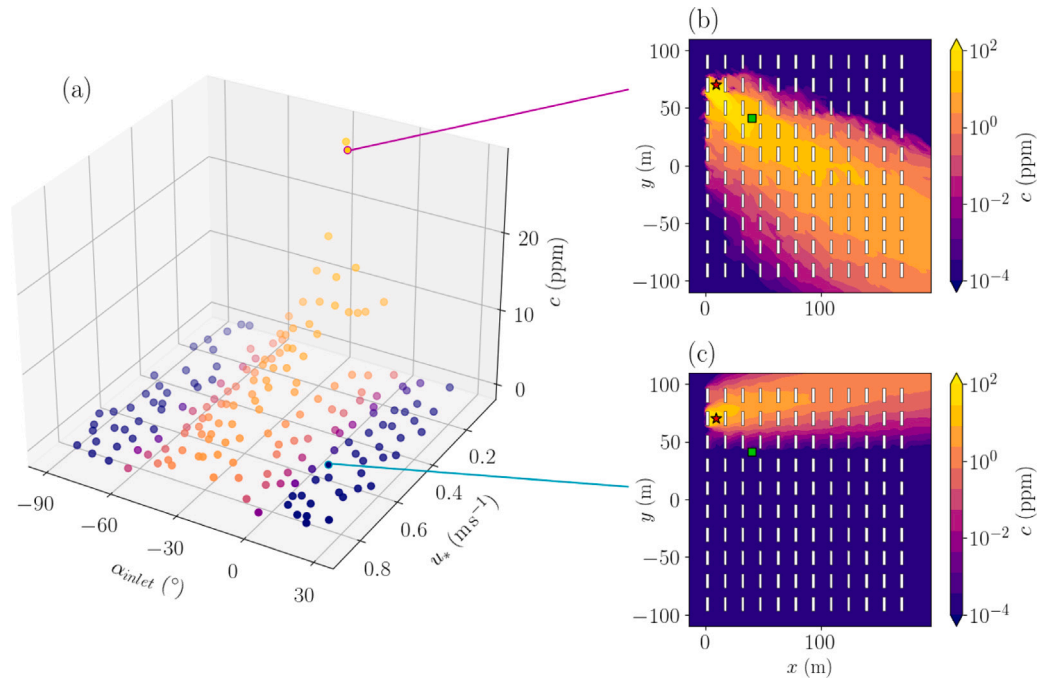


Fig. 2. (a) LES prediction of the local mean (time-averaged) concentration c at tower B at $z = 2 \text{ m}$ for each sample of parameters $\theta = (\alpha_{inlet}, u_*)$ from Fig. 1. (b, c) Horizontal cuts of the mean concentration at $z = 1.6 \text{ m}$ for the two samples $(\alpha_{inlet}^{(81)}, u_*^{(81)}) = (-27.7^\circ, 0.08 \text{ m s}^{-1})$ and $(\alpha_{inlet}^{(133)}, u_*^{(133)}) = (7.73^\circ, 0.60 \text{ m s}^{-1})$ in (a). The green square corresponds to the tower B, and the red star corresponds to the tracer source.

therefore a larger proportion of the domain where the concentration is zero at these angles.

This quantification of the noise in the learning dataset is of paramount importance for the construction and validation of surrogate models. In particular, this information can be used to select the dimension of the latent space to prevent the surrogate model from overfitting the noise associated with internal variability (Section 3.4). Internal variability estimates can also be used as a reference to check that the surrogate model uncertainty is not underestimated (Section 3.3), and as a performance target for the surrogate model (Section 3.5).

3. Surrogate modeling approach

This section presents the POD–GPR surrogate modeling approach and specifies the inputs/outputs and metrics used for validation. The focus is on two points. The first point is how to estimate the uncertainty associated with POD–GPR predictions and relate it to internal

variability. The second point is how to make an informed choice about the surrogate latent space dimension.

3.1. Problem statement

The goal of the surrogate model is to emulate as closely as possible the response surface of the LES model (Section 2.2) with respect to the input parameters $\theta = (\alpha_{inlet}, u_*)$ defined over the space Ω_θ (Eq. (4), Section 2.3). This means finding a function:

$$\begin{aligned} \mathcal{M}_{\text{surrogate}} : \Omega_\theta &\longrightarrow \mathbb{R}^N, \\ \theta &\longmapsto \mathbf{y}_{\text{surrogate}}, \end{aligned} \quad (5)$$

that minimizes $\int_{\Omega_\theta} \|\mathbf{y}_{\text{surrogate}}(\theta) - \mathbf{y}_{\text{LES}}(\theta)\| d\theta$, where $\mathbf{y}_{\text{LES}} \in \mathbb{R}^N$ is the field to be emulated, discretized on a grid of N nodes, and where $\mathbf{y}_{\text{surrogate}}$ is its counterpart predicted by the surrogate. This function is

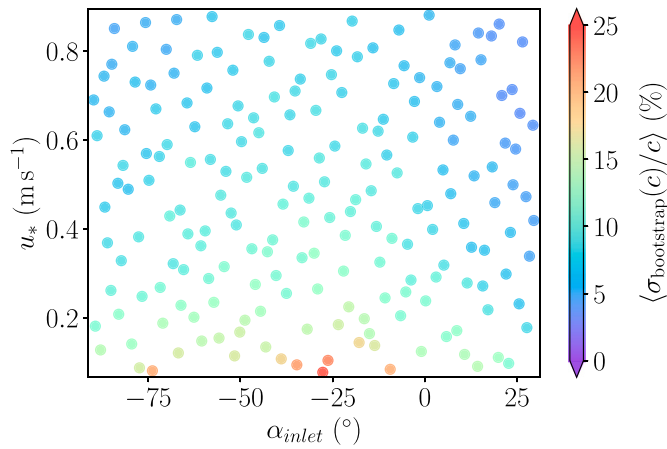


Fig. 3. Relative uncertainty of the mean concentration in the parameter space estimated using stationary bootstrap [29] and averaged over the whole spatial domain. Each circle corresponds to the averaged uncertainty of one LES sample of the learning dataset obtained from Halton's sequence (Fig. 1).

obtained here by learning from the training set $\{\theta^{(i)}, \mathbf{y}_{\text{LES}}(\theta^{(i)})\}_{i=1}^{N_{\text{train}}}$ with $N_{\text{train}} = 160$ (80% of the full LES dataset, see Fig. 1).

In this study, we focus on the emulation of the time-averaged tracer concentration fields, which are noisy due to the internal variability of the ABL (Section 2.5). Taking into account this aleatory uncertainty in the construction and validation of the surrogate model is a key challenge we address here.

To reduce the computational cost associated with the high dimension N of the solver grid on which the fields of interest are expressed, we interpolate all the fields on an analysis mesh twice as coarse, centered around the container array, and with a height limited to 20 m as most of the tracer is located in this area. This leads to an analysis mesh of $N = 1.88 \times 10^6$ nodes, with characteristic cell sizes ranging from 0.6 m to 4 m, which facilitates efficient model reduction. We have checked that using a coarser-resolution mesh has a negligible effect on the surrogate model accuracy (not shown here).

3.2. The POD–GPR surrogate model

3.2.1. Principle

We choose to use a POD–GPR surrogate model because it has proven to be efficient, relatively inexpensive and robust [39,47,79]. The fundamental principle of the POD–GPR approach is to combine:

- (i) a *reduction step* using Proper Orthogonal Decomposition (POD) [48,49], which is very popular in fluid mechanics [32,80,81] and consists in finding a low-dimensional space, called *latent space*, of dimension $L \ll N$, on which the fields to be emulated $\mathbf{y}(\theta)$ are projected;
- (ii) and a *regression step* using standard Gaussian Process Regression (GPR) [50], which consists in learning from the training set, the relationship between the LES model input parameters θ and the latent coefficients $\{k_\ell(\theta)\}_{\ell=1}^L$ resulting from the field projection onto the latent space.

This reduction-regression approach allows (i) to reduce the dimension of the regression problem to L latent variables ($L \ll N$) and thereby drastically reduce the computational burden of the learning task; and (ii) to separate the parametric dependence of the field from the spatial variability.

The POD–GPR model is implemented as a standard statistical learning approach, i.e. with an initial training phase consisting of (i) preprocessing the LES fields, (ii) building the POD reduced basis based on the training set, and (iii) optimizing the GPR models in the

latent space (Fig. 4a). This training phase is done offline and only once. The trained POD–GPR can then provide online field predictions for new inputs θ as follows: (i) the associated POD reduced coefficients are predicted by the fitted GPR models, and (ii) the inverse POD projection and inverse fields scaling are applied to these coefficients to recover the physical field $\mathbf{y}_{\text{surrogate}}$ (Fig. 4b). The following sections present the theoretical elements of the POD and GPR techniques required for this study.

3.2.2. Field preprocessing and dimension reduction using POD

With POD, the field dimension is reduced by linearly projecting the fields into the latent space generated by the basis of the POD modes $\{\psi_\ell\}_{\ell=1}^L \in \mathbb{R}^{N \times L}$. These modes are the eigenvectors obtained by diagonalizing the covariance matrix $\mathbf{C} = \frac{1}{N_{\text{train}} - 1} \mathcal{T}(\mathbf{S})\mathcal{T}(\mathbf{S})^T \in \mathbb{R}^{N \times N}$ of the snapshot matrix $\mathbf{S} = (\mathbf{y}_{\text{LES}}^{(1)} | \dots | \mathbf{y}_{\text{LES}}^{(N_{\text{train}})}) \in \mathbb{R}^{N \times N_{\text{train}}}$ composed of the training set of mean concentration fields, with \mathcal{T} a preprocessing that includes centering. The choice of the preprocessing is fundamental, and the one retained in this study is further defined in Eq. (7). The basis is then truncated to retain only the L eigenvectors $\{\psi_\ell\}_{\ell=1}^L$ associated with the L largest eigenvalues $\{\Lambda_\ell\}_{\ell=1}^L$ of the covariance matrix \mathbf{C} . These eigenvectors are the most informative about the coherent spatial structures emerging from variations in the wind conditions $\theta = (\alpha_{\text{inlet}}, u_*)$. The question of how to choose L is discussed in detail in Section 3.4.

The projection of one preprocessed field $\mathcal{T}(\mathbf{y})$ onto the POD latent space can be formulated as

$$\mathcal{T}(\mathbf{y}) = \sum_{\ell=1}^L \sqrt{\Lambda_\ell} k_\ell \psi_\ell, \quad (6)$$

where $\{k_\ell\}_{\ell=1}^L$ are the POD reduced coefficients defined as the coefficients in the projection of the given field $\mathbf{y}(\theta)$ normalized by $\sqrt{\Lambda_\ell}$. This scaling, called POD whitening [82], ensures that the set of reduced coefficients $\{k_\ell\}_{\ell=1}^L$ is centered and has unit component-wise variances on average, so that the regression problem is well posed for GPR.

The orthogonality of POD modes leads to some very useful properties [49,83]: (i) the POD decomposition (Eq. (6)) is the linear combination that reproduces the most variance of the original set, and (ii) POD reduced coefficients are uncorrelated, i.e. $\text{Cov}(k_i, k_j) = 0$, if $i \neq j$, which justifies why we build one GPR model per mode (Fig. 4).

For pollutant dispersion applications, a particular difficulty arises from the wide disparity of the concentration scale, which significantly limits POD approximation accuracy. This can be addressed by preprocessing the fields before building the POD, as this changes the meaning of the optimality and orthogonality properties of the POD modes [84], and thus conditions the POD ability to efficiently represent fields in a smaller dimension. Using a logarithmic preprocessing, which is a natural choice for concentrations since they can be assumed to follow a log-normal distribution [85], results in better overall projection performance for the MUST case study (not shown here — see [46], Chapter IV, for further discussion on preprocessing choices). The logarithmic preprocessing used in this study reads:

$$\mathcal{T} : \mathbb{R}^N \longrightarrow \mathbb{R}^N, \quad (7)$$

$$\mathbf{y}(\mathbf{x}_k) \longmapsto \sqrt{\frac{\omega(\mathbf{x}_k)}{\Omega}} \left[\ln(\mathbf{y}(\mathbf{x}_k) + y_t) - \langle \ln(\mathbf{y}_{\text{LES}}(\mathbf{x}_k) + y_t) \rangle_{\text{train}} \right], \quad 1 \leq k \leq N,$$

where $\omega(\mathbf{x}_k)/\Omega$ is the relative volume of the node \mathbf{x}_k , y_t is a threshold set to 10^{-4} ppm to avoid issues with concentration values close to zero, and $\langle \cdot \rangle_{\text{train}}$ denotes the mean over the training set. This choice provides an effective compromise that does not over-cut low concentrations and does not over-emphasize very low variances, which are mainly numerical noise. Note that this preprocessing also includes the centering required for POD [49], and volume node weighting to avoid over-weighting refined locations [84].

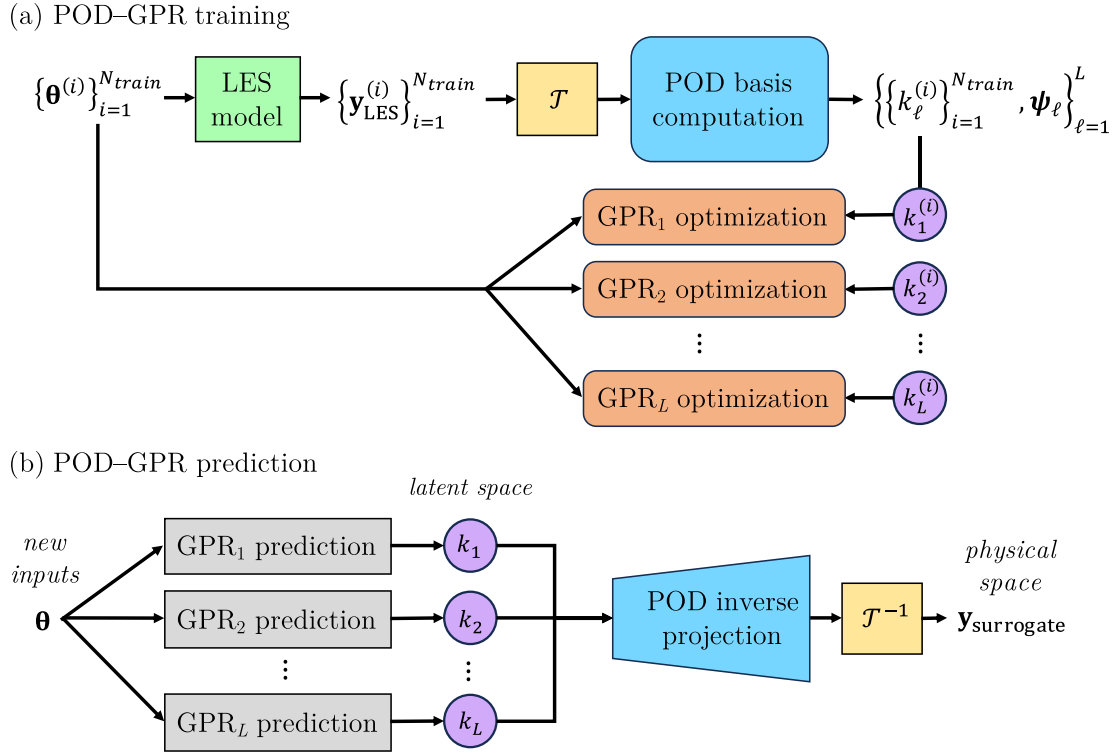


Fig. 4. Schematic of the POD-GPR surrogate model. Its operation is divided into two phases: the training phase (a) and the prediction phase (b). For the training phase, first, a preprocessing \mathcal{T} is applied to the LES training fields, which are then used to build the POD basis $\{\psi_l\}_{l=1}^L$; then L independent GPs are optimized to learn the dependence of the POD reduced coefficients $\{k_\ell\}_{\ell=1}^L$ on the input parameters θ . For the prediction phase, the optimized GPs predict the POD reduced coefficients associated with any set of wind conditions θ , then the inverse POD projection and inverse scaling \mathcal{T}^{-1} are applied to recover the associated physical field.

3.2.3. Latent coefficients estimation by Gaussian processes

Once the POD latent space is constructed, the next step is to predict the POD reduced coefficients $\{k_\ell(\theta)\}_{\ell=1}^L$ for any new wind conditions $\theta = (\alpha_{inlet}, u_*) \in \Omega_\theta \subset \mathbb{R}^2$ (Fig. 4b). Since POD coefficients are uncorrelated, we simplify this vector regression problem into L scalar regression problems solved by GPR [50]. There are three main reasons for this choice: (i) simple interpolation may fail to predict latent space components [86]; (ii) GPR was found to be one of the best machine learning regression methods for predicting POD-reduced coefficients of LES concentration fields [87]; and (iii) GPR models predict probability distributions and not just pointwise estimates, which is in line with our objective to quantify surrogate model uncertainties.

The principle of Gaussian processes (GP) is that the data distribution can be described by a Gaussian stochastic process, implying

$$k_\ell = f_\ell(\theta) + \epsilon_\ell \text{ with } \begin{cases} f_\ell(\theta) \sim \mathcal{GP}(\mathbf{0}, r_\ell(\theta, \theta^*)), \forall(\theta, \theta^*) \in \Omega_\theta^2 \\ \epsilon_\ell \sim \mathcal{N}(0, s_\ell^2), \end{cases} \quad (8)$$

where r_ℓ is the GP covariance function, or *kernel*, and where ϵ_ℓ is an additive Gaussian noise with variance s_ℓ^2 accounting for the fact that the k_ℓ are subject to an irreducible noise due to the internal variability of the mean concentration (Fig. 3). Note that we assume that the prior distribution of the GP is zero on average since POD reduced coefficients are centered on average.

Given the property that any finite subset of realizations of a GP follows a multivariate Gaussian distribution, the posterior probability distribution of the quantity of interest $k_\ell^*(\theta^*)$ knowing the training set $\{\theta^{train}, \mathbf{K}_\ell^{train}\}$ is

$$k_\ell^*(\theta^*) \Big|_{\{\theta^{train}, \mathbf{K}_\ell^{train}\}} \sim \mathcal{N}(\mu_\ell, \sigma_{GP}^2(k_\ell^*)), \quad (9)$$

with:

$$\begin{cases} \mu_\ell = r_\ell(\theta^*, \theta^{train}) [r_\ell(\theta^{train}, \theta^{train}) + s_\ell^2 \mathbf{I}]^{-1} \mathbf{K}_\ell^{train}, & (10a) \\ \sigma_{GP}^2(k_\ell^*) = r_\ell(\theta^*, \theta^*) + s_\ell^2 \\ \quad - r_\ell(\theta^*, \theta^{train}) [r_\ell(\theta^{train}, \theta^{train}) + s_\ell^2 \mathbf{I}]^{-1} r_\ell(\theta^{train}, \theta^*). & (10b) \end{cases}$$

In the regression context, these equations give the mean GPR prediction (Eq. (10a)) and the associated variance (Eq. (10b)), which quantifies two forms of uncertainty: (i) the uncertainty linked to the noise in the training data and related to the term $s_\ell^2 \mathbf{I}$, and (ii) the regression uncertainty that depends on the distance between the new input parameters θ^* and the training parameters θ^{train} . Both equations involve the kernel function r_ℓ to measure these distances. In this study, we use a standard Matérn kernel with the hyperparameter $\nu = 5/2$ [88].

In the end, each GP has four hyperparameters: the noise variance s_ℓ^2 , and three parameters involved in the Matérn kernel [88]: the maximum allowable covariance, and the two length scales associated with each of the two input parameters α_{inlet} and u_* . These hyperparameters are determined by maximum log-likelihood estimation [89] during GP optimization (Fig. 4a). Note that the input parameters θ are rescaled to $[0, 1]^2$ by min-max normalization to facilitate the GP optimization.

3.3. Uncertainty estimation of POD-GPR predictions

Below we explain how the GPR estimated uncertainty (Eq. (10b)) is propagated from latent space to physical space through the POD inverse projection. This is useful to quantify the uncertainty of POD-GPR field predictions.

POD-GPR predictions are defined as linear combinations of the POD reduced coefficients $k_\ell(\theta)$ (Eq. (6)), which are uncorrelated (by

POD modeling assumption) and normally distributed (Eq. (9)). Consequently, at each grid node \mathbf{x}_k , the POD–GPR prediction $\mathcal{T}(\mathbf{y}(\theta, \mathbf{x}_k))$ also follows a normal distribution of mean (Eq. (6)) and variance (Eq. (11)):

$$\sigma_{\text{POD-GPR}}^2(\mathcal{T}(\mathbf{y}(\theta, \mathbf{x}_k))) = \sum_{\ell=1}^L \Lambda_{\ell} \sigma_{\text{GP}}^2(k_{\ell}(\theta)) \boldsymbol{\psi}_{\ell}(\mathbf{x}_k)^2, \quad (11)$$

with $\sigma_{\text{GP}}^2(k_{\ell}(\theta))$ the ℓ th GP variance (Eq. (10b)).

Using logarithmic preprocessing (Eq. (7)), we deduce that the POD–GPR prediction of mean concentration follows a log-normal distribution of mean (Eq. (12a)) and variance (Eq. (12b)):

$$\left\{ \begin{aligned} \mathbf{y}_{\text{POD-GPR}}(\theta, \mathbf{x}_k) &= \sqrt{\frac{\Omega}{\omega(\mathbf{x}_k)}} \sum_{\ell=1}^L \sqrt{\Lambda_{\ell}} k_{\ell} \boldsymbol{\psi}_{\ell}(\mathbf{x}_k) + \langle \ln(\mathbf{y}_{\text{LES}} + y_t) \rangle_{\text{train}}, \\ \sigma_{\text{POD-GPR}}^2(\mathbf{y}(\theta, \mathbf{x}_k)) &= \left[\exp(s^2(\theta, \mathbf{x}_k)) - 1 \right] \times \exp(2\mathbf{y}(\theta, \mathbf{x}_k) + s^2(\theta, \mathbf{x}_k)), \end{aligned} \right. \quad (12a)$$

$$\sigma_{\text{POD-GPR}}^2(\mathbf{y}(\theta, \mathbf{x}_k)) = \left[\exp(s^2(\theta, \mathbf{x}_k)) - 1 \right] \times \exp(2\mathbf{y}(\theta, \mathbf{x}_k) + s^2(\theta, \mathbf{x}_k)), \quad (12b)$$

$$\text{with } s(\theta, \mathbf{x}_k)^2 = \left(\frac{\Omega}{\omega(\mathbf{x}_k)} \right) \sum_{\ell=1}^L \Lambda_{\ell} \sigma_{\text{GP}}^2(k_{\ell}(\theta)) \boldsymbol{\psi}_{\ell}(\mathbf{x}_k)^2.$$

Equation (12b) provides an estimate of the uncertainty around the POD–GPR mean prediction (Eq. (12a)). This uncertainty is the sum of the GP variances $\sigma_{\text{GP}}^2(k_{\ell}(\theta))$, which quantify the noise error in the training data and the regression error for each mode. In the present context, these two forms of error correspond to the uncertainty associated with the internal variability of the ABL (Section 2.5) and to the structural model error associated with model reduction. It is worth noting that this estimate does not include the error associated with the projection into the POD latent space.

3.4. A priori choice of latent space dimension

The choice of the POD latent space dimension is case-dependent and has a critical effect on the accuracy of the surrogate model. On the one hand, the higher the number of POD modes, the more variance of the original ensemble is captured in the POD reduced basis. On the other hand, high-order modes are likely to encode noise in the training set [90], and are therefore best set aside to prevent GPs from overfitting noise during learning. In this section, we present an innovative method to select L as a trade-off between the total variance embedded in the POD reduced basis and the amount of noise carried by the POD modes.

POD projection error. First, we evaluate the POD projection error, i.e. the error obtained after reconstructing the fields projected onto the POD latent space through inverse POD transformation, for varying number of modes L following the approach adopted by Nony et al. [39]. Fig. 5a shows that the POD projection normalized mean square error (NMSE) quickly decreases with the number of modes, and that a small number of modes ($L \approx 5 - 10$) allows to obtain very fine NMSE scores. We verify that the eigenvalues Λ_{ℓ} are a good proxy for quantifying the amount of information retrieved by each POD mode [49] and can therefore be used to select L as done by Xiao et al. [51].

Internal variability in the POD latent space. To quantify how the noise of the physical fields (see Section 2.5) is captured by each POD mode, we project $B = 1000$ bootstrap replicates of the LES fields onto the POD basis according to the procedure shown in Fig. 6. From the set of replicates $\left\{ \mu_b \left(\left\{ k_{\ell}^{(i)} \right\}_{\ell=1}^L \right) \right\}_{b=1}^B$ thus obtained, we can estimate the internal variability of the POD reduced coefficients for each mode order ℓ , and do so for each sample (i) in the training set. Fig. 5b shows that the standard deviation of the k_{ℓ} averaged over the training set increases significantly as the mode order ℓ increases. In particular, for $\ell \leq 5$, the standard deviation of the k_{ℓ} bootstrap replicates remain small ($< 3\%$), implying that these modes correspond to systematic patterns associated with the plume structure and its dependence on the wind conditions.

The standard deviation of the k_{ℓ} bootstrap replicates then increases rapidly reaching about 15%. This implies that field features linked to internal variability are mainly captured by higher order modes, which is consistent with the literature [90]. This in turn implies that we need to limit the number of modes L to avoid introducing noise into the POD–GPR surrogate model.

A priori criterion to choose the POD latent space dimension. Based on these findings, we propose to measure the ratio between the internal variability noise and the fraction of the total ensemble variance represented by each mode defined as

$$\frac{\langle \sigma_{\text{bootstrap}}^2(k_{\ell}) \rangle_{\text{train}}}{\Lambda_{\ell}}, \quad (13)$$

where $\langle \sigma_{\text{bootstrap}}^2(k_{\ell}) \rangle_{\text{train}}$ is the variance of the POD reduced coefficients replicates averaged over the training set, and where Λ_{ℓ} is the ℓ th eigenvalue in the POD decomposition.

The ratio in Eq. (13) is shown in Fig. 5c and provides a way to choose the latent space dimension L that minimizes both the noise and the POD projection error, and it has the advantage of being completely a priori as it does not require either the test set or the evaluation of the full POD–GPR model. Results show that this ratio is close to zero for the first six modes and then increases sharply with mode order. We therefore choose to truncate the POD decomposition before the inflection point using $L = 10$ modes to project the mean concentration fields. This approach for selecting the latent space dimension is evaluated a posteriori in Section 4.3.

3.5. Surrogate validation methodology

We present now the metrics used to quantify the surrogate model reduction error, before estimating the best values achievable for each metric given the internal variability.

3.5.1. Quantification of the surrogate error

The POD–GPR model accuracy is estimated on a set of independent test samples ($N_{\text{test}} = 40$, corresponding to 20% of the full LES dataset, see Fig. 1). This is essential to assess the ability of the model to generalize information from the training set to new meteorological forcing parameters ($\alpha_{\text{inlet}}, u_*$).

To assess the surrogate error, we use standard air quality metrics from Chang and Hanna [91] to compare the mean concentration field predicted by the surrogate model $\mathbf{c}_{\text{surrogate}}$ with the LES counterpart \mathbf{c}_{LES} . When validating POD–GPR predictions, which are probabilistic (see Section 3.3), we define $\mathbf{c}_{\text{surrogate}}$ as the mean of the probability distribution predicted by the POD–GPR (Eq. (12a)). The metrics used in this study are: the normalized mean square error (NMSE), the fraction of predictions within a factor of two of observations (FAC2), the geometric variance (VG), and the figure of merit in space (FMS):

$$\text{NMSE} = \frac{\langle (\mathbf{c}_{\text{LES}} - \mathbf{c}_{\text{surrogate}})^2 \rangle}{\langle \mathbf{c}_{\text{LES}} \rangle \langle \mathbf{c}_{\text{surrogate}} \rangle}, \quad (14)$$

$$\text{FAC2} = \langle \xi \rangle \text{ with } \xi(\mathbf{x}_k) = \begin{cases} 1 & \text{if } 0.5 \leq \mathbf{c}_{\text{surrogate}}(\mathbf{x}_k) / \mathbf{c}_{\text{LES}}(\mathbf{x}_k) \leq 2, \\ 1 & \text{if } \mathbf{c}_{\text{surrogate}}(\mathbf{x}_k) \leq c_t \text{ and } \mathbf{c}_{\text{LES}}(\mathbf{x}_k) \leq c_t, \\ 0 & \text{else,} \end{cases} \quad (15)$$

$$\text{VG} = \exp \left(\langle (\ln \tilde{\mathbf{c}}_{\text{LES}} - \ln \tilde{\mathbf{c}}_{\text{surrogate}})^2 \rangle \right), \quad (16)$$

$$\text{FMS}(c_{\ell}) = \frac{\Omega_{\cap}(c_{\ell})}{\Omega_{\cup}(c_{\ell})}, \quad (17)$$

where $\langle \cdot \rangle$ denotes spatial averaging weighted by the dual volume of the node \mathbf{x}_k , c_t is a concentration threshold defining $\tilde{\mathbf{c}} = \max(\mathbf{c}, c_t)$, as suggested by Chang and Hanna [91] and Schatzmann et al. [77]

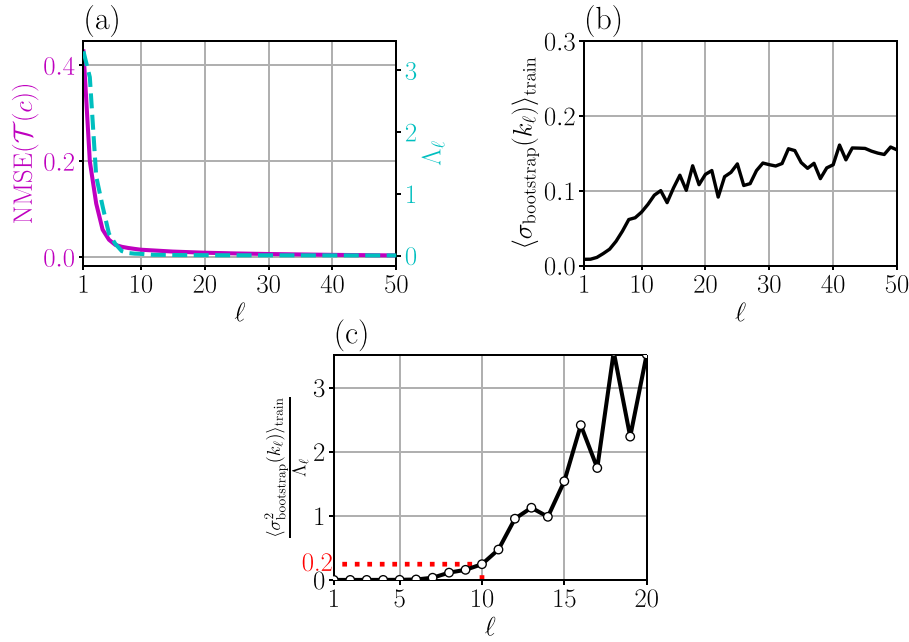


Fig. 5. (a) POD projection error evaluated over the training set with NMSE (Eq. (14)) as a function of the number of modes retained, and the POD eigenvalues Λ_ℓ associated with each mode ℓ . (b) Internal variability standard deviation of the POD reduced coefficients $\sigma_{\text{bootstrap}}(k_\ell)$ estimated from $B = 1000$ bootstrap replicates obtained using the procedure shown in Fig. 6, and averaged over the training set. (c) Ratio between the internal variability noise variance $\sigma_{\text{bootstrap}}^2(k_\ell)$ averaged over the training set and the POD eigenvalues Λ_ℓ associated with each mode ℓ . The red dotted line indicates the number of modes selected for this study.

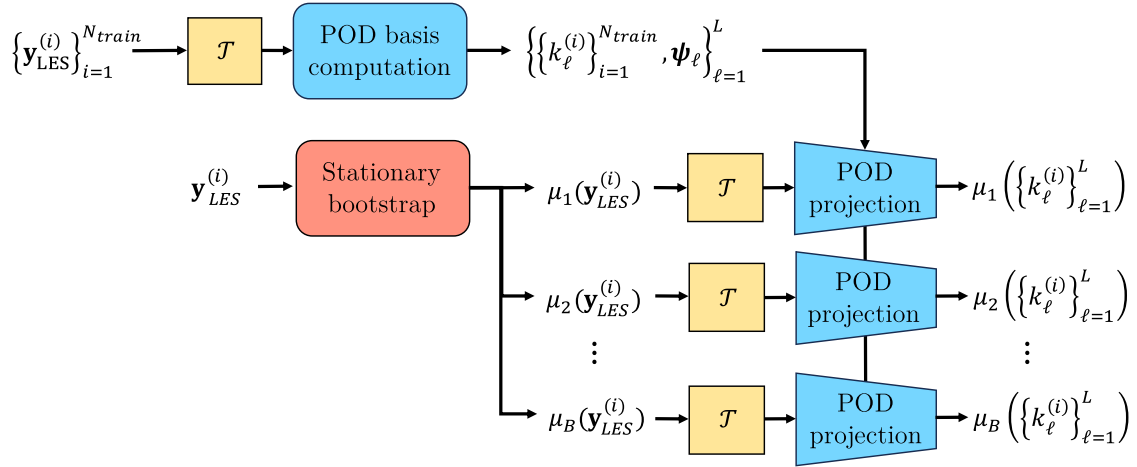


Fig. 6. Schematic of the propagation of the internal variability from the physical space to the POD latent coefficients. First, the POD basis is computed from the training set of LES fields $\{\mathbf{y}_{\text{LES}}^{(i)}\}_{i=1}^{N_{\text{train}}}$. Then, to propagate the internal variability of a given field $\mathbf{y}_{\text{LES}}^{(i)}$ into the latent space, B bootstrap replicates $\mu_b(\mathbf{y}_{\text{LES}}^{(i)})$ are computed using the bootstrap approach from Lumet et al. [29], and then projected onto the POD basis to get a set of B realizations $\mu_b(\{k_\ell^{(i)}\}_{\ell=1}^L)$ of the POD reduced coefficients of the field $\mathbf{y}_{\text{LES}}^{(i)}$.

to avoid issues with values close to zero in FAC2 and VG metrics. In this study, we use a threshold of $c_t = 10^{-4}$ ppm, considering that errors on lower concentrations are mainly due to numerical noise. Finally, $\Omega_\cap(c_\ell)$ denotes the volume, in m^3 , of the domain in which both $\mathbf{c}_{\text{surrogate}}$ and \mathbf{c}_{LES} are over a user-specified tracer value c_ℓ . Conversely, $\Omega_\cup(c_\ell)$ denotes the volume where $\mathbf{c}_{\text{surrogate}} \geq c_\ell$ or $\mathbf{c}_{\text{LES}} \geq c_\ell$.

The use of different metrics than the loss used during training is important because of the multi-order nature of the concentration field. NMSE is more sensitive to errors at high concentrations, while VG assesses prediction accuracy at low concentrations. FMS quantifies how close the two plume shapes are relative to a given concentration level. The scores that a perfect model would obtain are reported in Table 1.

3.5.2. Estimation of the internal variability

LES data are noisy due to internal variability (Section 2.5). It would therefore be pointless to try to build a surrogate model whose accuracy exceeds this uncertainty. To quantify the error due to internal variability alone, we use the bootstrap approach proposed in Lumet et al. [29] to generate two independent sets of bootstrap replicates of the same LES field. We then compute the average difference between each pair of replicates using the metrics introduced in Section 3.5.1. For each metric, we obtain the amount of error due to internal variability only, which is the expected error when comparing two independent realizations of the mean concentration fields for the same input parameters.

Table 1

Prediction accuracy of the POD–GPR surrogate model evaluated using the metrics defined in Section 3.5.1 and averaged over the test set. The standard deviations of the scores over the test set are also given, as well as the individual scores for test samples #81 and #187, which represent the lowest and highest FAC2 scores achieved by the POD–GPR, respectively. For comparison, the perfect scores for the metrics, the mean error due to internal variability only (Section 3.5.2) and the mean error due to standalone reduction dimension are given.

| | FAC2 | NMSE | VG | FMS (1 ppm) | FMS (0.01 ppm) |
|--------------------------|-------------|-------------|-------------|----------------|-------------------|
| Perfect score | 1 | 0 | 0 | 1 | 1 |
| Internal variability | 0.95 | 1.80 | 1.39 | 0.83 | 0.93 |
| POD projection error | 0.91 | 20.4 | 1.33 | 0.75 | 0.93 |
| POD–GPR prediction error | 0.91 | 20.6 | 1.39 | 0.75 | 0.92 |
| Standard deviation | 0.04 | 43.2 | 0.68 | 0.11 | 0.03 |
| Test sample #81 | 0.74 | 23.4 | 5.25 | 0.79 | 0.85 |
| Test sample #187 | 0.96 | 8.08 | 1.07 | 0.86 | 0.94 |

This is done for every LES sample in the dataset, and the ensemble-averaged internal variability errors give an upper bound estimate of the best overall accuracy achievable for each metric when validating the POD–GPR surrogate model.

4. Surrogate model validation

In this section, we present a thorough evaluation of the POD–GPR surrogate model. We first assess its accuracy over the test set and its efficiency (Section 4.1). We then validate the innovative aspects of our approach: the POD–GPR uncertainty estimation (Section 4.2), and the selection of the number of POD modes (Section 4.3). Finally, we study how the POD–GPR model behaves when reducing the training set (Section 4.4). All results are given for the mean concentration field, but the POD–GPR approach can be applied to emulate other fields, such as the wind velocity and turbulent kinetic energy or the concentration fluctuations and peaks, as shown in Lumet [46], Appendix B.1.

4.1. Evaluation of the surrogate model field predictions

We evaluate here the POD–GPR predictions of mean concentration following the methodology introduced in Section 3.5, using the mean internal variability error as the reference for validation. We use a latent space dimension of $L = 10$ in accordance with the informed choice made in Section 3.4.

Prediction accuracy. The overall performance of the surrogate model is quantified using standard air quality metrics (Section 3.5.1). Table 1 shows the obtained scores averaged over the test set. Overall, the POD–GPR model yields very satisfactory results, with most scores close to the error due to internal variability only, which is the best achievable accuracy. However, the scores for FMS(1 ppm) and especially NMSE remain relatively far from the internal variability error, indicating that POD–GPR is less good at predicting high concentration values. These large errors near the source can be particularly detrimental for assessing acute exposure to pollutant species. For this type of application, we recommend using a different preprocessing of the data to give more weight to high concentration areas during training.

Table 1 also shows that the POD–GPR prediction errors are almost identical to the standalone POD projection errors (i.e. errors obtained by simply reconstructing the test fields after projection onto the POD basis by inverse POD transformation). This implies that the accuracy of the POD–GPR model is mostly limited by the accuracy of the POD and not by the GPR. The poor prediction performance for high concentrations is thus explained by the fact that the POD is not well adapted to the multiscale and nonlinear nature of the concentration fields. In particular, the use of a logarithmic preprocessing before the POD degrades the reconstruction of high concentrations in the vicinity of the emission source, but has the advantage of preserving the other metrics and in particular the shape of the plume compared to linear processing (Lumet [46], Chapter IV).

There is quite a large spread of POD–GPR errors across the test samples, especially for the quadratic metrics NMSE and VG, indicating the presence of test sample outliers. This variability over the input parameter space is mainly explained by the fact that as the friction velocity decreases, the internal variability increases (Fig. 3), which makes the mean concentration noisier and therefore more difficult to predict. In addition, FMS(1 ppm), and to a lesser extent FMS(0.01 ppm) and FAC2, are subject to a zoning effect as they depend on the size of the plume within the domain of interest (Eqs. (15) and (17)). For example, these scores are improved when the wind direction carries the plume outside the container array (i.e. for $\alpha_{inlet} \approx 30^\circ$ or $\alpha_{inlet} \approx -90^\circ$).

Field prediction examples. For a more detailed assessment of the POD–GPR model accuracy, we also examine its predictions in the physical space. Figs. 7a, b, c, and d compare 2-D cuts of the mean concentration at $z = 1.6$ m predicted by LES and POD–GPR. Results are given for the test sample #187 ($\alpha_{inlet}^{(187)}, u_*^{(187)} = (21.8^\circ, 0.59 \text{ m s}^{-1})$) for which POD–GPR obtains the best FAC2 score over the test set, and for the test sample #81 ($\alpha_{inlet}^{(81)}, u_*^{(81)} = (-27.7^\circ, 0.08 \text{ m s}^{-1})$) associated with the worst FAC2 score. The global scores obtained for these two particular snapshots are summarized in Table 1.

In both cases, the POD–GPR model reproduces well the main features of the LES concentration field, in particular the shape and orientation of the plume. The spatial distribution of the different concentration levels is also well reproduced, which is confirmed by the near superposition of the 0.01 ppm and 10 ppm concentration contour lines between LES and POD–GPR (Fig. 7d, h).

However, for the sample with the worst FAC2 (#81), the POD–GPR underestimates the spanwise spread of the plume and significantly overestimates the mean concentration near the emission source (Fig. 7g). This is consistent with the poor NMSE obtained (Table 1) and this is due to the poor reproduction of high concentrations by the POD with logarithmic preprocessing. For this sample (corresponding to a low friction velocity and therefore subject to substantial internal variability), the POD–GPR tends to smooth the irregularities observed at the edges of the plume, thus poorly predicting the local abrupt decrease in concentration.

Efficiency. In terms of computational cost, it takes approximately 30 s to train the POD–GPR model using a single core of an Intel Ice Lake CPU. This includes field preprocessing, POD basis decomposition and GPR optimization. This training cost is insignificant compared to the cost of building the training dataset (Section 2.4). Once trained, the model provides a prediction of the full 3-D concentration field in about 0.03 s. This approach is therefore compatible with applications requiring a large ensemble of predictions and/or real-time predictions.

4.2. Assessment of the surrogate model uncertainty estimation

We evaluate here the ability of the POD–GPR model to provide realistic uncertainty estimates by comparing them to the actual surrogate error over the test set and to the internal variability present in the LES dataset.

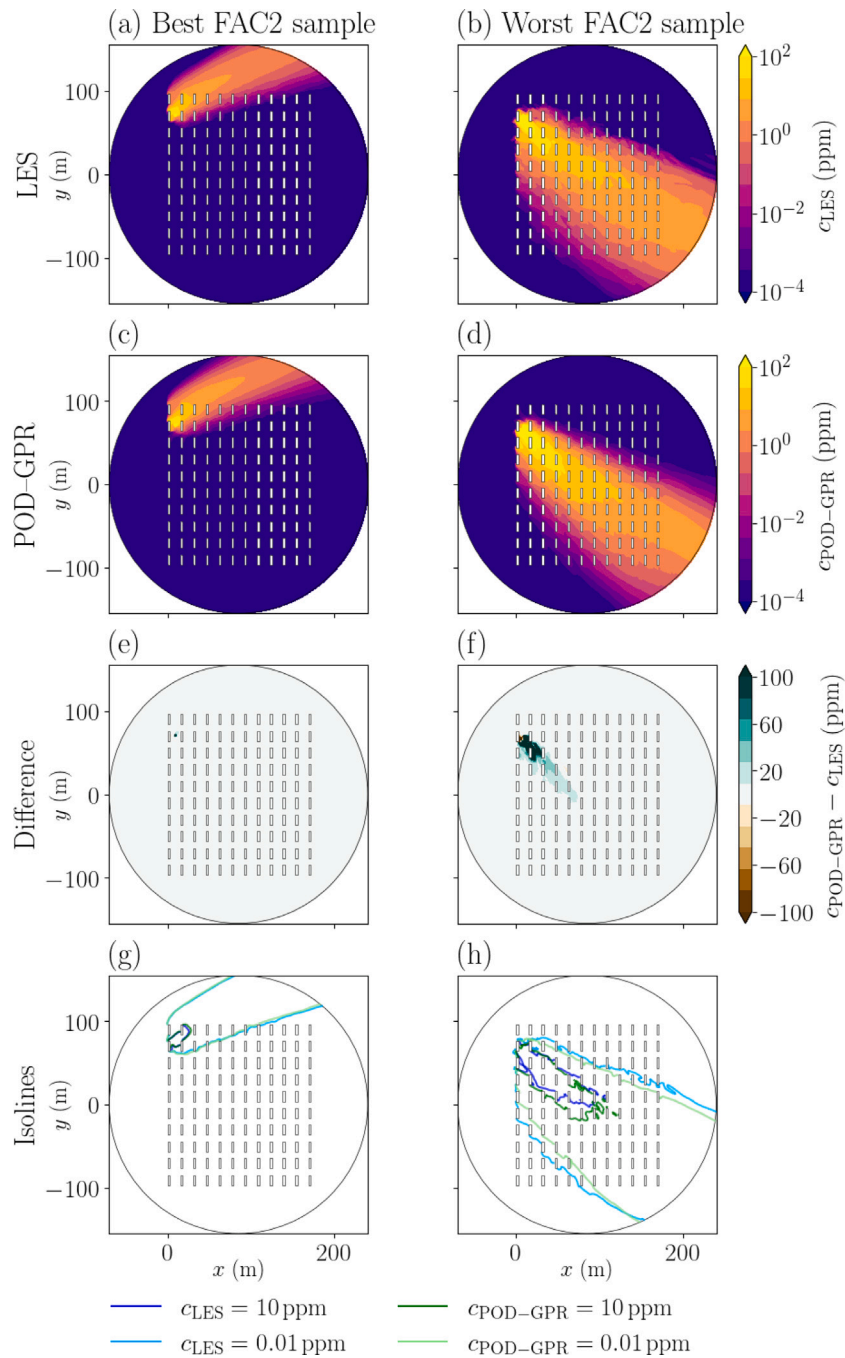


Fig. 7. Horizontal cuts at $z = 1.6$ m of two test mean concentration fields estimated by LES (a, b) and POD-GPR (c, d), and the absolute difference between the two (e, f). The left column corresponds to the test sample #187 for which POD-GPR achieves the best FAC2 (Eq. (15)) score over the test set, and the right column corresponds to the test sample #81 which is associated with the worst FAC2 score. The LES and POD-GPR predictions of 0.01 ppm and 10 ppm iso concentration levels are shown in (g, h).

Uncertainty reliability. Fig. 8a shows the uncertainty reliability diagram comparing actual surrogate error (y-axis) and surrogate model uncertainty estimates (x-axis). The POD-GPR uncertainty is underestimated compared to the actual POD-GPR error for most domain nodes, especially for the lowest concentration values. Nevertheless, the estimated trend is consistent, i.e. the larger the actual error, the larger the prior estimate. Furthermore, the overall level of precision is satisfactory as the estimated uncertainty is in the right order of magnitude (within the green dashed lines) for 98% of the domain nodes. This is confirmed by the response surface of the POD-GPR (Fig. 13a, b), as the predicted envelopes appear to cover the test samples well. We can therefore be confident in the uncertainty predicted by the POD-GPR surrogate model despite a tendency to underestimate.

To further investigate the cause of this underestimation, the uncertainty reliability is examined directly in the latent space in Fig. 8b. We find that for the estimation of the reduced POD coefficients by the GPs, the uncertainty estimate is very close to the error made on average, except for the high-order modes 8 and 10. This increase in error for higher-order modes is consistent with the fact that they are more affected by internal variability (Fig. 5b). The following conclusions can be drawn: (i) the variance of the GP posterior distribution (Eq. (10b)) is realistic, and (ii) the underestimation observed in the physical space in Fig. 8a comes from the inverse POD projection. This is consistent with the fact that the POD projection error is not taken into account when estimating the total POD-GPR uncertainty (Section 3.3), yet the

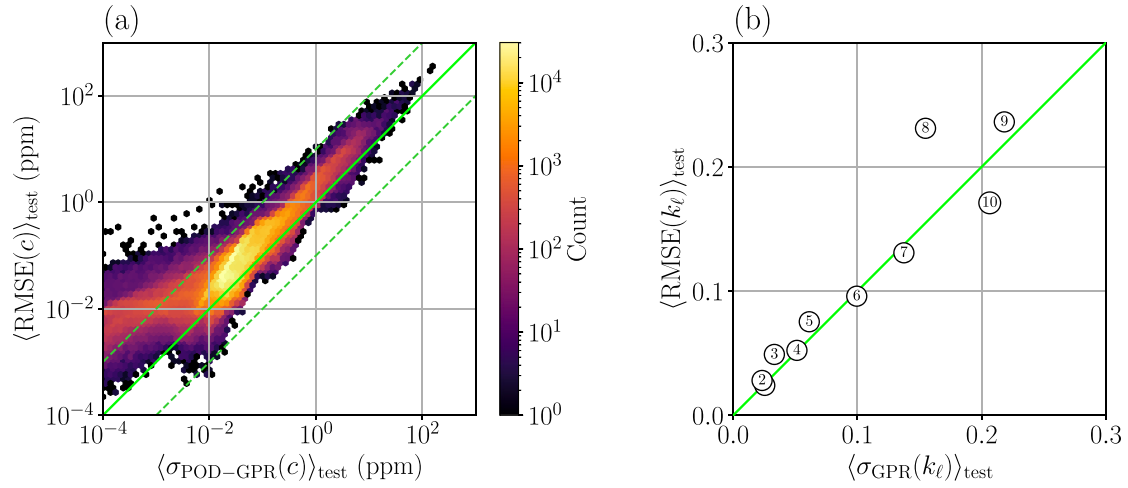


Fig. 8. Reliability diagrams in the physical space and in the latent space: (a) Root mean square error (RMSE) of the POD-GPR concentration prediction over the test set versus the POD-GPR estimated uncertainty at each node where the concentration is larger than the tolerance $c_t = 10^{-4}$ ppm. Each hexagon is colored according to the number of node points in the hexagon. (b) RMSE of the GP prediction of the POD reduced coefficients k_ℓ over the test set versus the GP estimated uncertainty, each mode ℓ is represented by a numbered circle (the POD latent space dimension is $L = 10$). The green solid lines correspond to the identity function, and the dashed lines in (a) show the range of plus or minus one order of magnitude.

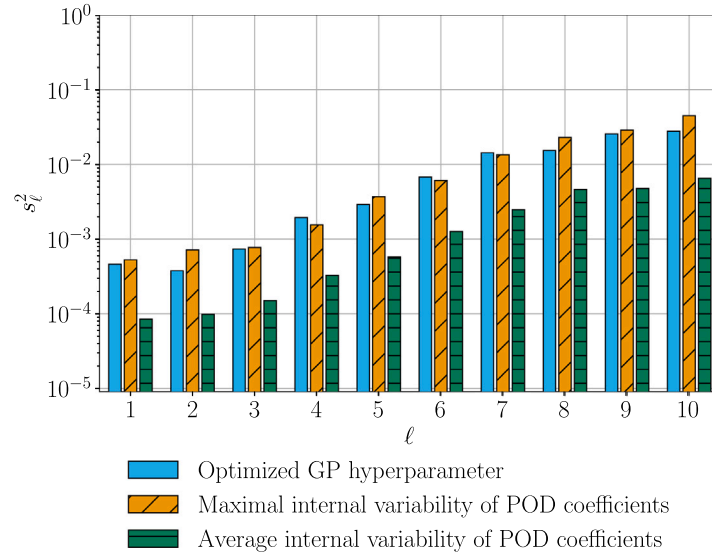


Fig. 9. GP noise variance s_ℓ^2 hyperparameter obtained by log-likelihood maximization for each mode ℓ as blue bars, and maximal (resp. average) noise on the POD reduced coefficients over the training set as orange (resp. green) bars.

total POD-GPR error is essentially due to the POD projection error as indicated in Table 1.

Ability to estimate internal variability a posteriori. We now examine the nature of the estimated uncertainty in more detail, and assess the proportion due to internal variability. The first point is to study how the noise of the LES fields projected onto the POD latent space is captured by GPR. Fig. 9 shows that the values of the GP variance hyperparameters s_ℓ^2 obtained by maximum likelihood estimation are very close to the maximum level of internal variability of the POD reduced coefficients over the training set estimated by bootstrap. This is a strong result because the bootstrap estimates of the internal variability are not used to train the GPs.

The fact that the GP noise variance parameter matches the maximum level of internal variability (Fig. 9) implies that GPs overestimate the variance of the POD reduced coefficients for most samples where the internal variability is low. This is a structural limitation due to the fact that the GP additive noise does not depend on the input parameter space (Eq. (8)), while the variance due to internal variability does

(Fig. 3). As a result, in the physical space, the POD-GPR uncertainty predictions tend to be underestimated compared to the actual internal variability for samples where the internal variability is high, while they are overestimated for samples with low internal variability. This could be partially overcome in the future by implementing input-dependent noise variance hyperparameters, as suggested by Miyagusuku et al. [92].

Fig. 10 shows that the uncertainty estimated by the POD-GPR is overall consistent with the LES internal variability over the training set, as the level of variability is in the right order of magnitude for 99% of the domain nodes. For most of the domain, the POD-GPR tends to overestimate the internal variability (hexagonal cells of high density beyond the green line), which is consistent with the GP noise matching the maximum level of internal variability in the latent space (Fig. 9). Note that this analysis is performed over the training set since for these samples the GPR regression covariance is zero, and thus the POD-GPR uncertainty estimate only corresponds to the estimated internal variability. Finally, we note that the estimated uncertainty envelopes

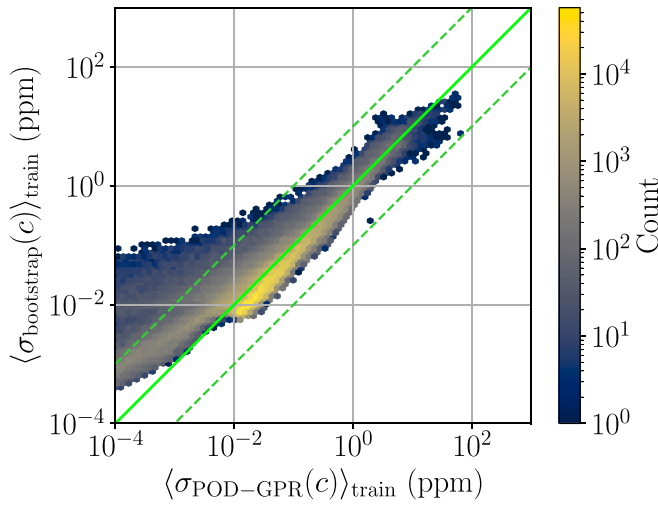


Fig. 10. Internal variability of the mean concentration estimated by bootstrap and averaged over the training set versus the POD-GPR estimated uncertainty at each node where the concentration is larger than the tolerance $c_t = 10^{-4}$ ppm. Each hexagon is colored according to the number of node points in the hexagon. The green solid lines correspond to the identity function and the dashed lines show the range of plus or minus one order of magnitude.

are consistent with the LES internal variability when looking at the POD-GPR response surfaces (Fig. 13a, b).

In this internal variability analysis, the second point is to evaluate the spatial consistency of the POD-GPR uncertainty estimates with respect to the spatial distribution of internal variability to verify that the uncertainty is properly propagated from the POD latent space to the physical space (Section 3.3). We find that the variance predicted by the POD-GPR is consistent with the internal variability of the LES in terms of magnitude and structure, as shown in Fig. 11 for the training sample #016 ($\alpha_{inlet}^{(016)}, u_*^{(016)} = (-79.5^\circ, 0.14 \text{ m s}^{-1})$) for which the POD-GPR uncertainty estimate is the closest to the internal variability estimated by bootstrap and for the training sample #180 ($\alpha_{inlet}^{(180)}, u_*^{(180)} = (-58.1^\circ, 0.56 \text{ m s}^{-1})$), where POD-GPR overestimates the internal variability the most. Despite the overall agreement, the POD-GPR variability estimates appear to be overestimated within the plume and significantly underestimated near the plume edges (Fig. 11e, f), which is consistent with the overall tendency to underestimate low internal variability levels (Fig. 10). This is explained by the fact that there are high concentration gradients near the plume edges and thus high internal variability levels, a feature not well represented by the POD projection, which is based solely on mean concentration and not on its variability.

In summary, the POD-GPR uncertainty estimates derived in Section 3.3 (i) represent, in a spatially coherent manner, the inherent internal variability of the mean concentration field thanks to the ability of the GPs to accurately infer the level of noise in the training set, and (ii) properly explain the actual surrogate errors at predicting the mean concentration. This particularly reinforces the robustness of the POD-GPR and its relevance to uncertainty quantification applications.

4.3. A posteriori validation of the latent space dimension

We revisit our choice of the number of POD modes ($L = 10$) obtained by following the a priori statistical approach we propose in Section 3.4. For this purpose, we evaluate the effect of the number of modes L on the performance of the full POD-GPR model on the test set (i.e. by varying L from 5 to 50 in the construction of the POD-GPR model).

Validation metrics. Fig. 12 shows how the metrics defined in Section 3.5.1 change when modifying the POD latent space dimension L . The POD-GPR prediction accuracy over the test set increases with the number of modes and reaches a plateau for a larger number of modes ($L \approx 15-25$) than the NMSE on the training set used in our mode choice approach (Fig. 5). This may indicate that integrating a larger number of modes into the POD-GPR model could lead to improved surrogate model accuracy.

Response surfaces. As an additional diagnostic, Fig. 13 shows that using a larger number of modes significantly deteriorates the POD-GPR response surfaces, making them very noisy and implausible as, with $L = 50$ modes (Fig. 13e, f), the POD-GPR model is no longer able to retrieve the inversely proportional dependence of concentration on friction velocity expected from theory and retrieved for the configuration with $L = 10$ modes (Fig. 13a, b). This degradation is due to the fact that high-order modes mostly account for noisy structures due to internal variability (Fig. 5b), and are therefore not informative on systematic structures related to the wind conditions. As a result, when including high-order modes, the GPs attempt to learn unphysical dependence on the input parameters, resulting in the shortwave noise observed in Fig. 13. Still, the increase in uncertainty with the response surface deterioration suggests that the POD-GPR uncertainty estimate is robust. However, the fact that the degradation of the POD-GPR response surface is not seen by the global metrics, which continue to improve as the number of modes increases (Fig. 12), shows that one should not draw conclusions based on scalar metrics alone.

In the light of these tests, our prior selection method for the latent space dimension is convincing. The resulting trade-off of $L = 10$ yields good validation scores, while avoiding the problem of response surface noise. However, we acknowledge that using a slightly larger number of modes ($L \approx 15-20$) would also be appropriate and even slightly improve the surrogate model accuracy. Defining an optimal criterion for latent space dimension selection based on the noise/signal ratio defined in Eq. (13) is therefore an interesting prospect, but requires more validation cases.

4.4. Robustness to training set reduction

In order to assess the potential of the POD-GPR approach for future applications, we examine how the POD-GPR accuracy evolves as the size of the training set decreases (without changing the test set). This is particularly important to investigate the possible trade-offs between the ability of the model to generalize from training data and the substantial cost of building the LES training dataset.

The surrogate model is trained for decreasing training set sizes from $N_{train} = 160$ to $N_{train} = 40$ by keeping only the first samples in Halton's sequence. To make the comparison fair, we systematically evaluate the averaged prediction errors over the same test set of $N_{test} = 40$ samples. Results are shown in Fig. 14a, b, c and d in terms of FAC2, VG, FMS(1 ppm), FMS(0.01 ppm). The decrease in accuracy is fairly constrained and evolves linearly with the training set size, with a loss of 0.08 in FAC2 and 0.12 in VG for every 10 training samples removed. More importantly, the accuracy decreases less rapidly than that of the nearest neighbor model (1-NN), which trivially predicts the mean concentration field as equal to the closest training field in the parameter space (see Appendix). This is especially true for the low concentration values, as the VG score is significantly higher with the 1-NN model than with the POD-GPR model for small training set sizes (Fig. 14b).

Regarding the NMSE metric (Fig. 14e), the evolution with N_{train} is quite chaotic for the POD-GPR and worse than for the 1-NN approach. As previously mentioned, this is related to the high POD projection error near the source when using the logarithmic transformation, and we can consider that the POD-GPR approach with the present preprocessing is not designed to make predictions near the source, regardless of the training set size.

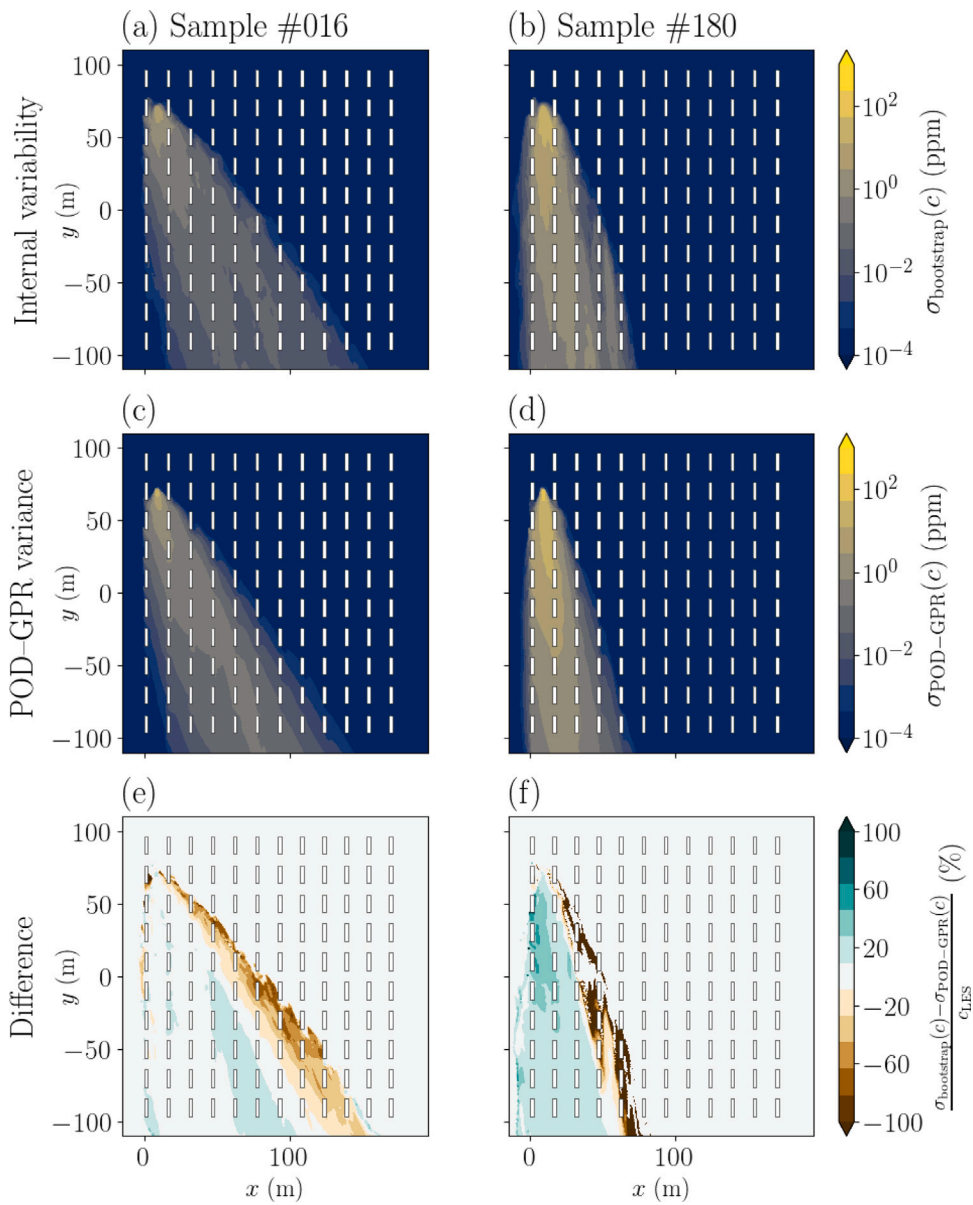


Fig. 11. Horizontal cuts at $z = 1.6$ m of the standard deviation of the mean concentration induced by internal variability estimated using bootstrap (a, b), predicted by POD-GPR (c, d), and the relative difference between the two (e, f). The left column corresponds to the training sample #016 and the right column corresponds to the training sample #180.

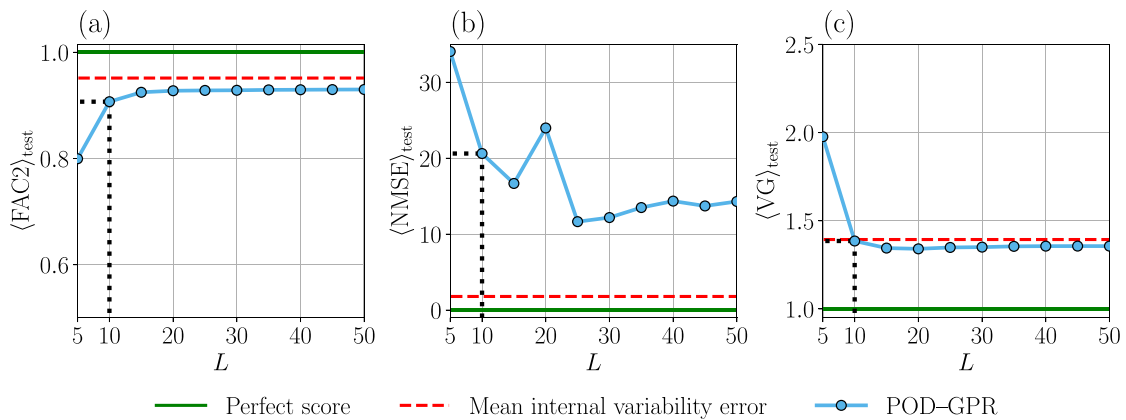


Fig. 12. POD-GPR prediction error as a function of the number of the modes L and evaluated with FAC2 (a), NMSE (b), VG (c) averaged over the test set. Green lines correspond to perfect scores; and red dashed lines correspond to the mean level of error due to internal variability only. Error levels corresponding to the selected number of modes ($L = 10$) are shown as black dotted lines.

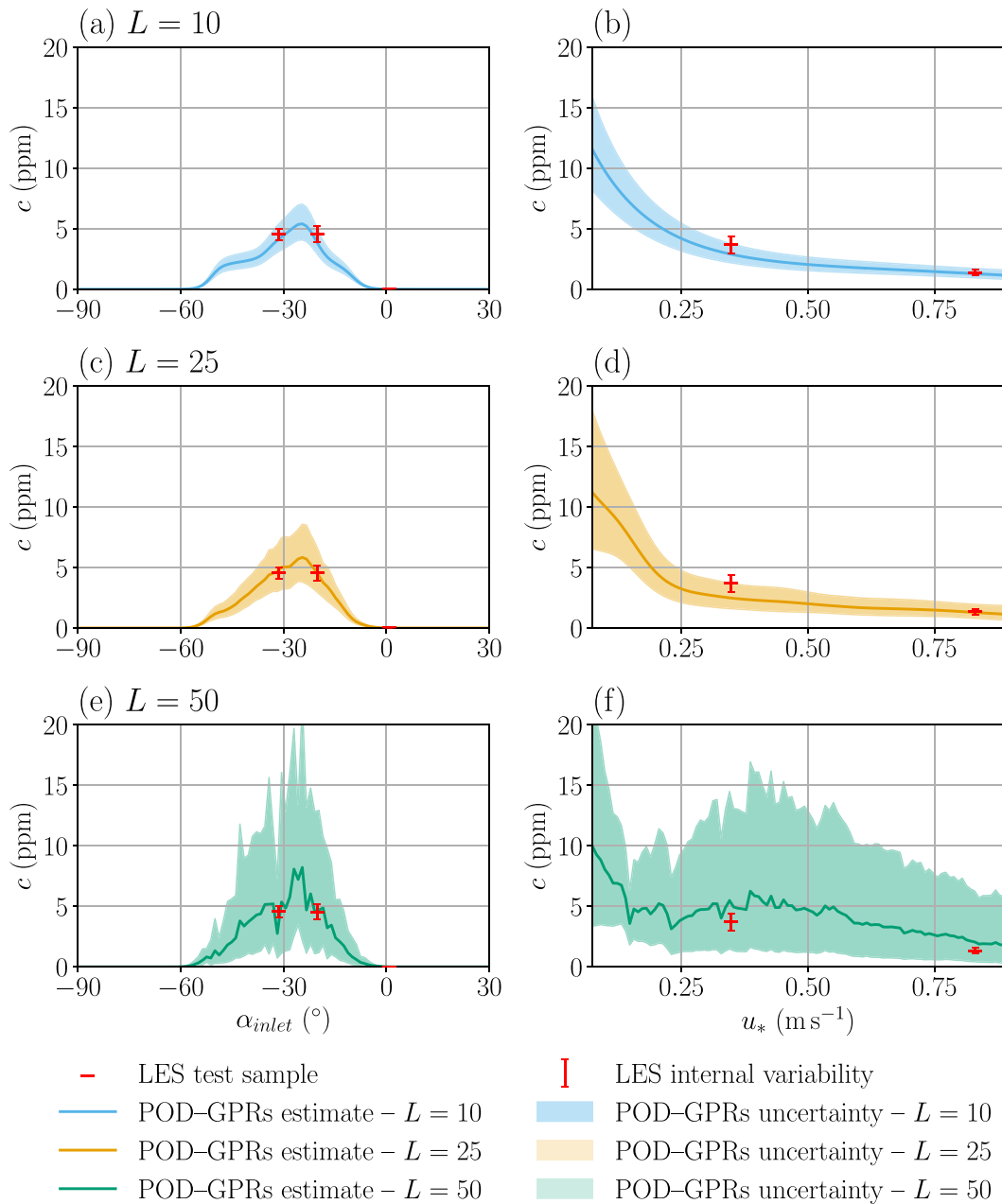


Fig. 13. POD-GPR prediction of the mean concentration at tower B at $z = 2\text{m}$ (see tower location in Fig. 2) as a function of the inlet wind direction α_{inlet} (a, c, e), and of the friction velocity u_* (b, d, f). Shaded areas correspond to the 95% confidence intervals estimated by the POD-GPR according to the procedure detailed in Section 3.3. Each row corresponds to the results obtained with a different latent space dimension $L \in \{10, 25, 50\}$. When varying one parameter, the other is set constant to either $u_*^{plot} = 0.45\text{ m s}^{-1}$ (a, c, e), or $\alpha_{inlet}^{plot} = -43^\circ$ (b, d, f), and the test samples closest to the two segments of parameter space thus scanned (see Fig. 1) are represented by horizontal red bars. The uncertainty on LES test samples induced by internal variability is depicted as red vertical error bars.

Fig. 15 shows that the POD-GPR uncertainty predictions are very robust to training set size reduction. We find that, on average, the POD-GPR uncertainty predictions explain overall well its actual error over the test set even with only 40 training samples (Fig. 15a, b, c). Similarly, the ability of the POD-GPR to represent the internal variability of the mean concentration is well preserved (Fig. 15d, e, f), although we note a tendency to underestimate it when the training set size is reduced, as there are fewer close neighboring points for the GPs to estimate the noise in this case.

In summary, the ability of the POD-GPR model to generalize from a training set of limited size is better than for the 1-NN baseline approach, justifying the use of such a more sophisticated surrogate model. We find that for this problem, 40 LES training samples are sufficient to achieve good levels of accuracy for most metrics. Furthermore, the

uncertainty estimates provided by POD-GPR remain consistent as the training set size decreases, despite a tendency to overestimate.

5. Conclusion

In this study, a data-driven surrogate dispersion model based on the two-stage POD-GPR approach was trained and rigorously evaluated using a large dataset of 200 LES simulations reproducing microscale dispersion scenarios of the field-scale MUST experiment for varying meteorological forcing. The resulting surrogate model is able to capture well the general plume shape within the canopy, approaching the best achievable accuracy given the internal variability in the LES data, while being very computationally efficient.

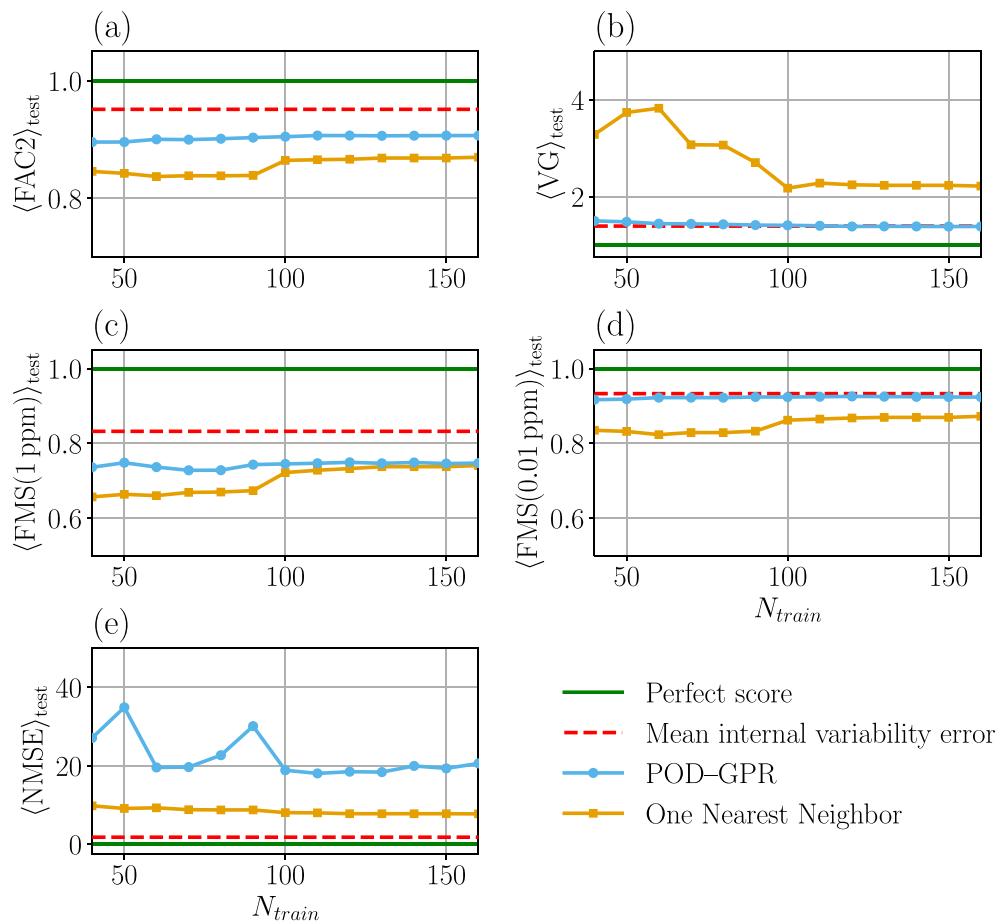


Fig. 14. Surrogate modeling errors for decreasing training set sizes. The mean concentration prediction error is assessed using the metrics defined in Section 3.5.1: namely FAC2 (a), VG (b), FMS (c, d), and NMSE (e). Results are given for the POD-GPR as blue circles and the 1-NN model as orange squares. Perfect scores are represented as green lines; and red dashed lines correspond to the mean level of error due to internal variability only.

The main novelty of this study is the in-depth analysis of the POD-GPR surrogate model uncertainty and of the weight of internal variability, thus meeting the need expressed by Tominaga et al. [7], Dauxois et al. [10] and Wu and Quan [30]. Future developments are required to account for the POD projection error in the POD-GPR uncertainty estimate to avoid error underestimation. But the present uncertainty estimates already explain the differences between the POD-GPR predictions and the LES references quite well, being in the right order of magnitude in 97% of cases. This work thus represents an important methodological step towards the representation of total uncertainty in microscale urban pollutant dispersion, as aleatory and modeling uncertainties have not been considered in most uncertainty quantification [11,42] and data assimilation [25,43–45,93,94] studies to date.

A second important contribution of this study is the method for selecting a priori the POD latent space dimension, which is based on a trade-off between the accuracy of the POD reconstruction and the noise captured by the POD modes estimated by bootstrap as in Lumet et al. [29]. The threshold used here to make this trade-off need to be consolidated and made more objective in future studies by considering a wide range of cases. For this study, the retained dimension ($L = 10$) is smaller than the dimension chosen based on the standalone reconstruction error [39,51], but this choice is justified by the fact that using more modes ($L > 25$) significantly noises and degrades the POD-GPR response surface despite slightly better global metrics such as FAC2 and NMSE. This highlights that a surrogate model validation process learning from LES data, especially for the concentration variable, should not be based solely on global metrics but requires more local and structural analyses.

In this study, the main shortcoming of the POD-GPR approach is its lack of accuracy in areas of high concentration, i.e. close to the source. This is mainly due to POD, as a linear transformation is not well suited to the wide disparity in concentration scales and introduces projection errors. A promising way to overcome this issue is the mixture-of-experts approach, inspired by the work of El Garroussi et al. [95], whose key idea is to train several POD-GPR models, each corresponding to a different preprocessing, to capture the different concentration scales ([46], Appendix B.3). Another promising perspective is the use of nonlinear dimension reduction techniques such as neural network autoencoders [38,53,87,96]. However, a difficulty lies in the interpretation of the nonlinear latent space and in the identification of the internal variability noise.

We emphasize that the drastic reduction in prediction time offered by the POD-GPR approach comes at the expense of a very high computational cost for building the LES learning database (on the order of a million core hours), which may hinder the use of this approach in practical engineering applications. Therefore, defining the minimum number of LES samples required for training is a key issue in LES surrogate modeling. In this study, we show that the POD-GPR approach copes very well with a reduction of the training set down to 40 samples for two input parameters. The number of training samples could be further reduced by applying adaptive sampling methods to target learning zones [97,98].

Finally, it should be recognized that the POD-GPR surrogate model presented in this study is limited in its generalization ability as it was only trained for the MUST building layout and source location. In the future, learning the dependence of pollutant dispersion on

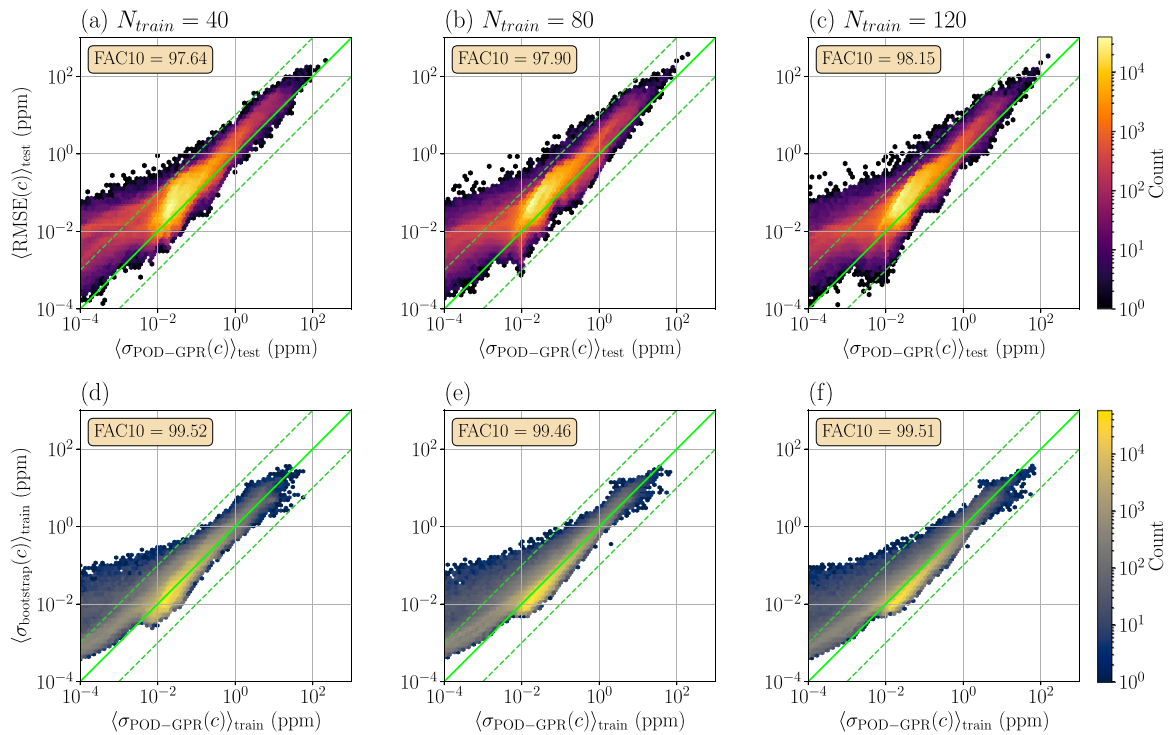


Fig. 15. Reliability diagrams of the POD-GPR uncertainty estimates for varying training set size $N_{train} \in \{40, 80, 120\}$. (a, b, c) Root mean square error (RMSE) of the POD-GPR concentration prediction over the test set and (d, e, f) internal variability of the mean concentration estimated by bootstrap and averaged over the training set, both versus the POD-GPR estimated uncertainty at each node where the concentration is larger than the tolerance $c_t = 10^{-4}$ ppm. Each hexagon is colored according to the number of node points in the hexagon. The green solid lines correspond to the identity function and the dashed lines show the range of plus or minus one order of magnitude. The FAC10 scores give the percentage of points between the two dashed lines (similarly as in Eq. (15)).

urban geometry and source location will require significantly larger training datasets, which may not be feasible due to the computational cost of LES. Multi-fidelity approaches that combine the high-fidelity information provided by LES with a less expensive model such as a RANS model [74,99,100], are a promising way to enrich the learning dataset while minimizing computational cost, thus paving the way for the uncertainty-aware POD-GPR surrogate model to be used for more general and complex urban pollutant dispersion studies.

CRedit authorship contribution statement

Elliott Lumet: Writing – review & editing, Writing – original draft, Visualization, Investigation. **Mélanie C. Rochoux:** Writing – review & editing, Writing – original draft, Supervision. **Thomas Jaravel:** Writing – review & editing, Supervision. **Simon Lacroix:** Writing – review & editing, Supervision.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

Elliott Lumet's Ph.D. thesis was funded by the Université Fédérale Toulouse Midi-Pyrénées, France together with the Région Occitanie (ADI20, AtmoDrones project, 2020–2023). This work was granted access to the HPC resources of GENCI-TGCC/CINES (A0062A10822 project, 2020–2022). The authors acknowledge Bastien Nony for preliminary code development and helpful discussion.

Appendix. The nearest neighbor control surrogate model

We use a Nearest Neighbor model (1-NN) as a simple baseline model against which we compare the POD-GPR accuracy. It is an appropriate control model because it represents the generalization error obtained by simply querying the available simulation dataset, and thus represents the minimum level of error that the POD-GPR must exceed to be worth using. The 1-NN is a classical k -Nearest Neighbor (k -NN) model [89] with only one neighbor ($k = 1$). The 1-NN prediction is simply defined as the nearest LES field in the training set:

$$y_{\text{surrogate}}(\theta) = y_{\text{LES}}^{\text{train}}(\theta^*), \quad \text{with } \theta^* = \min_{1 \leq i \leq N_{\text{train}}} d(\theta_i^{\text{train}}, \theta), \quad (\text{A.1})$$

where d is the Euclidean distance in a rescaled input space:

$$d(\theta^{(1)}, \theta^{(2)}) = \sqrt{\left(\frac{\alpha_{\text{inlet}}^{(2)} - \alpha_{\text{inlet}}^{(1)}}{\alpha_{\text{inlet}}^{\text{max}} - \alpha_{\text{inlet}}^{\text{min}}} \right)^2 + \zeta^2 \left(\frac{u_*^{(2)} - u_*^{(1)}}{u_*^{\text{max}} - u_*^{\text{min}}} \right)^2} \quad (\text{A.2})$$

where $\alpha_{\text{inlet}}^{\text{min}}$, $\alpha_{\text{inlet}}^{\text{max}}$, u_*^{min} , and u_*^{max} are the input space boundaries, and ζ is a rescaling factor that distorts the distances in the parameter space.

The hyperparameter ζ gives more or less weight to the friction velocity when searching for the closest LES field in the dataset (Eq. (A.1)). It is optimized during training by cross-validation [89] with 8-fold resampling of the training set. The best compromise between RMSE, VG and FMS(1 ppm) scores is obtained for $\zeta = 0.275$, which reduces the distances along the friction velocity axis and therefore gives more weight to the inlet wind direction parameter.

Data availability

The dataset used in this paper is openly available on a public repository (PPMLES, [76]). A notebook describing the construction and validation of the POD-GPR surrogate model is openly available at https://github.com/elliott-lumet/pod_gpr_ppmles. Other analysis codes

developed for this study are available from the corresponding author upon reasonable request.

References

- [1] EEA, Air quality in Europe, report, European Environment Agency, 2020, URL <https://www.eea.europa.eu/publications/air-quality-in-europe-2020>.
- [2] I. Manisalidis, E. Stavropoulou, A. Stavropoulos, E. Bezirtzoglou, Environmental and health impacts of air pollution: A review, *Front. Public Health* 8 (2020) <http://dx.doi.org/10.3389/fpubh.2020.00014>.
- [3] H.J.S. Fernando, S.M. Lee, J. Anderson, M. Princevac, E. Pardyjak, S. Grossman-Clarke, Urban fluid mechanics: Air circulation and contaminant dispersion in cities, *Environ. Fluid Mech.* 1 (1) (2001) 107–164, <http://dx.doi.org/10.1023/A:1011504001479>.
- [4] P. Klein, B. Leitl, M. Schatzmann, Driving physical mechanisms of flow and dispersion in urban canopies, *Int. J. Climatol.* 27 (14) (2007) 1887–1907, <http://dx.doi.org/10.1002/joc.1581>.
- [5] M. Pasquier, S. Jay, J. Jacob, P. Sagaut, A lattice-Boltzmann-based modelling chain for traffic-related atmospheric pollutant dispersion at the local urban scale, *Build. Environ.* 242 (2023) 110562, <http://dx.doi.org/10.1016/j.buildenv.2023.110562>.
- [6] B. Blocken, Computational fluid dynamics for urban physics: Importance, scales, possibilities, limitations and ten tips and tricks towards accurate and reliable simulations, *Build. Environ.* 91 (2015) 219–245, <http://dx.doi.org/10.1016/j.buildenv.2015.02.015>, Fifty Year Anniversary for Building and Environment.
- [7] Y. Tominaga, L.L. Wang, Z.J. Zhai, T. Stathopoulos, Accuracy of CFD simulations in urban aerodynamics and microclimate: Progress and challenges, *Build. Environ.* 243 (2023) 110723, <http://dx.doi.org/10.1016/j.buildenv.2023.110723>.
- [8] M. Schatzmann, B. Leitl, Issues with validation of urban flow and dispersion CFD models, *J. Wind Eng. Ind. Aerodyn.* 99 (4) (2011) 169–186, <http://dx.doi.org/10.1016/j.jweia.2011.01.005>, The Fifth International Symposium on Computational Wind Engineering.
- [9] B. Blocken, 50 years of computational wind engineering: Past, present and future, *J. Wind Eng. Ind. Aerodyn.* 129 (2014) 69–102, <http://dx.doi.org/10.1016/j.jweia.2014.03.008>.
- [10] T. Dauxois, T. Peacock, P. Bauer, C.P. Caulfield, C. Cenedese, C. Gorié, G. Haller, G.N. Ivey, P.F. Linden, E. Meiburg, N. Pardini, N.M. Vriend, A.W. Woods, Confronting grand challenges in environmental fluid mechanics, *Phys. Rev. Fluids* 6 (2021) 020501, <http://dx.doi.org/10.1103/PhysRevFluids.6.020501>.
- [11] C. García-Sánchez, D. Phillips, C. Gorié, Quantifying inflow uncertainties for CFD simulations of the flow in downtown Oklahoma City, *Build. Environ.* 78 (2014) 118–129, <http://dx.doi.org/10.1016/j.buildenv.2014.04.013>.
- [12] D.D. Lucas, A. Gowardhan, P. Cameron-Smith, R.L. Baskett, Impact of meteorological inflow uncertainty on tracer transport and source estimation in urban atmospheres, *Atmos. Environ.* 143 (2016) 120–132, <http://dx.doi.org/10.1016/j.atmosenv.2016.08.019>.
- [13] D. Wise, V. Boppana, K. Li, H. Poh, Effects of minor changes in the mean inlet wind direction on urban flow simulations, *Sustain. Cities Soc.* 37 (2018) 492–500, <http://dx.doi.org/10.1016/j.scs.2017.11.041>.
- [14] J.L. Santiago, A. Dejoan, A. Martilli, F. Martin, A. Pinelli, Comparison between large-eddy simulation and Reynolds-averaged Navier–Stokes computations for the MUST field experiment. Part I: Study of the flow for an incident wind directed perpendicularly to the front array of containers, *Bound.-Layer Meteorol.* 135 (1) (2010) 109–132, <http://dx.doi.org/10.1007/s10546-010-9466-3>.
- [15] H. Montazeri, B. Blocken, CFD simulation of wind-induced pressure coefficients on buildings with and without balconies: Validation and sensitivity analysis, *Build. Environ.* 60 (2013) 137–149, <http://dx.doi.org/10.1016/j.buildenv.2012.11.012>.
- [16] C. Gromke, N. Jarmarkattel, B. Ruck, Influence of roadside hedgerows on air quality in urban street canyons, *Atmos. Environ.* 139 (2016) 75–86, <http://dx.doi.org/10.1016/j.atmosenv.2016.05.014>.
- [17] V. Winiarek, M. Bocquet, O. Saunier, A. Mathieu, Estimation of errors in the inverse modeling of accidental release of atmospheric pollutant: Application to the reconstruction of the cesium-137 and iodine-131 source terms from the Fukushima Daiichi power plant, *J. Geophys. Res.: Atmos.* 117 (D5) (2012) <http://dx.doi.org/10.1029/2011JD016932>.
- [18] T.O. Spicer, G. Tickle, Simplified source description for atmospheric dispersion model comparison of the Jack Rabbit II chlorine field experiments, *Atmos. Environ.* 244 (2021) 117866, <http://dx.doi.org/10.1016/j.atmosenv.2020.117866>.
- [19] Y. Tominaga, T. Stathopoulos, Turbulent Schmidt numbers for CFD analysis with various types of flowfield, *Atmos. Environ.* 41 (37) (2007) 8091–8099, <http://dx.doi.org/10.1016/j.atmosenv.2007.06.054>.
- [20] B. Blocken, T. Stathopoulos, P. Saathoff, X. Wang, Numerical evaluation of pollutant dispersion in the built environment: Comparisons between models and experiments, *J. Wind Eng. Ind. Aerodyn.* 96 (10) (2008) 1817–1831, <http://dx.doi.org/10.1016/j.jweia.2008.02.049>, 4th International Symposium on Computational Wind Engineering (CWE2006).
- [21] G.-W.H. Yue Yang, L.-P. Wang, Effects of subgrid-scale modeling on Lagrangian statistics in large-eddy simulation, *J. Turbul.* 9 (2008) N8, <http://dx.doi.org/10.1080/14685240801905360>.
- [22] Y. Tominaga, T. Stathopoulos, Numerical simulation of dispersion around an isolated cubic building: Comparison of various types of $k-\epsilon$ models, *Atmos. Environ.* 43 (20) (2009) 3200–3210, <http://dx.doi.org/10.1016/j.atmosenv.2009.03.038>.
- [23] C. Gorié, G. Iaccarino, A framework for epistemic uncertainty quantification of turbulent scalar flux models for Reynolds-averaged Navier–Stokes simulations, *Phys. Fluids* 25 (5) (2013) 055105, <http://dx.doi.org/10.1063/1.4807067>.
- [24] C. Gorié, C. Garcia-Sanchez, G. Iaccarino, Quantifying inflow and RANS turbulence model form uncertainties for wind engineering flows, *J. Wind Eng. Ind. Aerodyn.* 144 (2015) 202–212, <http://dx.doi.org/10.1016/j.jweia.2015.03.025>.
- [25] H. Xiao, J.-L. Wu, J.-X. Wang, R. Sun, C. Roy, Quantifying and reducing model-form uncertainties in Reynolds-averaged Navier–Stokes simulations: A data-driven, physics-informed Bayesian approach, *J. Comput. Phys.* 324 (2016) 115–136, <http://dx.doi.org/10.1016/j.jcp.2016.07.038>.
- [26] M. Neophytou, A. Gowardhan, M. Brown, An inter-comparison of three urban wind models using Oklahoma City Joint Urban 2003 wind field measurements, *J. Wind Eng. Ind. Aerodyn.* 99 (4) (2011) 357–368, <http://dx.doi.org/10.1016/j.jweia.2011.01.010>.
- [27] G. Antonioni, S. Burkhart, J. Burman, A. Dejoan, A. Fusco, R. Gaasbeek, T. Gjesdal, A. Jäppinen, K. Riikonen, P. Morra, O. Parmhed, J. Santiago, Comparison of CFD and operational dispersion models in an urban-like environment, *Atmos. Environ.* 47 (2012) 365–372, <http://dx.doi.org/10.1016/j.atmosenv.2011.10.053>.
- [28] C. García-Sánchez, J. van Beeck, C. Gorié, Predictive large eddy simulations for urban flows: Challenges and opportunities, *Build. Environ.* 139 (2018) 146–156, <http://dx.doi.org/10.1016/j.buildenv.2018.05.007>.
- [29] E. Lumet, T. Jaravel, M.C. Rochoux, O. Vermorel, S. Lacroix, Assessing the internal variability of large-Eddy simulations for microscale pollutant dispersion prediction in an idealized urban environment, *Bound.-Layer Meteorol.* 190 (2) (2024) 9, <http://dx.doi.org/10.1007/s10546-023-00853-7>.
- [30] Y. Wu, S.J. Quan, A review of surrogate-assisted design optimization for improving urban wind environment, *Build. Environ.* 253 (2024) 111157, <http://dx.doi.org/10.1016/j.buildenv.2023.111157>.
- [31] T. Lassila, A. Manzoni, A. Quarteroni, G. Rozza, Model order reduction in fluid dynamics: challenges and perspectives, *Reduced Order Methods for modeling and computational reduction* (2014) 235–273, http://dx.doi.org/10.1007/978-3-319-02090-7_9.
- [32] R. Vinuesa, S.L. Brunton, Enhancing computational fluid dynamics with machine learning, *Nat. Comput. Sci.* 2 (6) (2022) 358–366, <http://dx.doi.org/10.1038/s43588-022-00264-7>.
- [33] Y. Wu, Q. Zhan, S.J. Quan, Y. Fan, Y. Yang, A surrogate-assisted optimization framework for microclimate-sensitive urban design practice, *Build. Environ.* 195 (2021) 107661, <http://dx.doi.org/10.1016/j.buildenv.2021.107661>.
- [34] C. Huang, G. Zhang, J. Yao, X. Wang, J.K. Calautit, C. Zhao, N. An, X. Peng, Accelerated environmental performance-driven urban design with generative adversarial network, *Build. Environ.* 224 (2022) 109575, <http://dx.doi.org/10.1016/j.buildenv.2022.109575>.
- [35] M. Mendil, S. Leirens, P. Armand, C. Duchenne, Hazardous atmospheric dispersion in urban areas: A deep learning approach for emergency pollution forecast, *Environ. Model. Softw.* 152 (2022) 105387, <http://dx.doi.org/10.1016/j.envsoft.2022.105387>.
- [36] P. Kastner, T. Dogan, A GAN-based surrogate model for instantaneous urban wind flow prediction, *Build. Environ.* 242 (2023) 110384, <http://dx.doi.org/10.1016/j.buildenv.2023.110384>.
- [37] L. Margheri, P. Sagaut, A hybrid anchored-ANOVA – POD/Kriging method for uncertainty quantification in unsteady high-fidelity CFD simulations, *J. Comput. Phys.* 324 (2016) 137–173, <http://dx.doi.org/10.1016/j.jcp.2016.07.036>.
- [38] S. Xiang, X. Fu, J. Zhou, Y. Wang, Y. Zhang, X. Hu, J. Xu, H. Liu, J. Liu, J. Ma, S. Tao, Non-intrusive reduced order model of urban airflow with dynamic boundary conditions, *Build. Environ.* 187 (2021) 107397, <http://dx.doi.org/10.1016/j.buildenv.2020.107397>.
- [39] B.X. Nony, M.C. Rochoux, T. Jaravel, D. Lucor, Reduced-order modeling for parameterized large-eddy simulations of atmospheric pollutant dispersion, *Stoch. Environ. Res. Risk Assess.* 37 (6) (2023) 2117–2144, <http://dx.doi.org/10.1007/s00477-023-02383-7>.
- [40] K. Cheng, Z. Lu, C. Ling, S. Zhou, Surrogate-assisted global sensitivity analysis: an overview, *Struct. Multidiscip. Optim.* 61 (2020) 1187–1213, <http://dx.doi.org/10.1007/s00158-019-02413-5>.
- [41] N. Fellmann, M. Pasquier, C. Blanchet-Scalliet, C. Helbert, A. Spagnol, D. Sinoquet, Sensitivity analysis for sets : application to pollutant concentration maps, 2023.
- [42] C. García-Sánchez, G. Van Tendeloo, C. Gorié, Quantifying inflow uncertainties in RANS simulations of urban pollutant dispersion, *Atmos. Environ.* 161 (2017) 263–273, <http://dx.doi.org/10.1016/j.atmosenv.2017.04.019>.
- [43] V. Mons, L. Margheri, J.-C. Chassaing, P. Sagaut, Data assimilation-based reconstruction of urban pollutant release characteristics, *J. Wind Eng. Ind. Aerodyn.* 169 (2017) 232–250, <http://dx.doi.org/10.1016/j.jweia.2017.07.007>.

- [44] J. Sousa, C. García-Sánchez, C. Gorié, Improving urban flow predictions through data assimilation, *Build. Environ.* 132 (2018) 282–290, <http://dx.doi.org/10.1016/j.buildenv.2018.01.032>.
- [45] J. Sousa, C. Gorié, Computational urban flow predictions with Bayesian inference: Validation with field data, *Build. Environ.* 154 (2019) 13–22, <http://dx.doi.org/10.1016/j.buildenv.2019.02.028>.
- [46] E. Lumet, Assessing and reducing uncertainty in large-eddy simulation for microscale atmospheric dispersion (Ph.D. thesis), Université Toulouse III - Paul Sabatier, 2024, URL <https://theses.fr/2024TLSES003>, (Accessed 30 May 2024).
- [47] A. Marrel, N. Perot, C. Mottet, Development of a surrogate model and sensitivity analysis for spatio-temporal numerical simulators, *Stoch. Environ. Res. Risk Assess.* 29 (3) (2015) 959–974, <http://dx.doi.org/10.1007/s00477-014-0927-y>.
- [48] L. Sirovich, Turbulence and the dynamics of coherent structures. I. Coherent structures, *Quart. Appl. Math.* 45 (3) (1987) 561–571, <http://dx.doi.org/10.1090/qam/910462>.
- [49] G. Berkooz, P. Holmes, J.L. Lumley, The proper orthogonal decomposition in the analysis of turbulent flows, *Annu. Rev. Fluid Mech.* 25 (1) (1993) 539–575, <http://dx.doi.org/10.1146/annurev.fl.25.010193.002543>.
- [50] C.E. Rasmussen, C.K. Williams, et al., Gaussian processes for machine learning, vol. 1, Springer, 2006, <http://dx.doi.org/10.7551/mitpress/3206.001.0001>.
- [51] D. Xiao, C. Heaney, F. Fang, L. Mottet, R. Hu, D. Bistrrian, E. Aristodemou, I. Navon, C. Pain, A domain decomposition non-intrusive reduced order model for turbulent flows, *Comput. & Fluids* 182 (2019) 15–27, <http://dx.doi.org/10.1016/j.compfluid.2019.02.012>.
- [52] A. Weerasuriya, X. Zhang, B. Lu, K. Tse, C. Liu, A Gaussian process-based emulator for modeling pedestrian-level wind field, *Build. Environ.* 188 (2021) 107500, <http://dx.doi.org/10.1016/j.buildenv.2020.107500>.
- [53] S. Masoumi-Verki, F. Haghghat, U. Eicker, A review of advances towards efficient reduced-order models (ROM) for predicting urban airflow and pollutant dispersion, *Build. Environ.* 216 (2022) 108966, <http://dx.doi.org/10.1016/j.buildenv.2022.108966>.
- [54] E. Yee, C.A. Biltoft, Concentration fluctuation measurements in a plume dispersing through a regular array of obstacles, *Bound.-Layer Meteorol.* 111 (3) (2004) 363–415, <http://dx.doi.org/10.1023/B:BOUN.0000016496.83909.ee>.
- [55] P. Gousseau, B. Blocken, T. Stathopoulos, G. van Heijst, CFD simulation of near-field pollutant dispersion on a high-resolution grid: A case study by LES and RANS for a building group in downtown montreal, *Atmos. Environ.* 45 (2) (2011) 428–438, <http://dx.doi.org/10.1016/j.atmosenv.2010.09.065>.
- [56] C. Biltoft, Customer report for Mock Urban Setting Test. DPG Document No. WDT-FCR-01-121, West Desert Test Center, U.S. Army Dugway Proving Ground, Utah, USA, 2001.
- [57] J. Franke, A. Hellsten, H. Schlünzen, B. Carissimo, Best practice guideline for the CFD simulation of flows in the urban environment, Technical report, COST European Cooperation in Science and Technology, 2007, URL <https://hal.science/hal-04181390>, (Accessed 01 December 2023).
- [58] S.R. Hanna, O.R. Hansen, S. Dharmavaram, FLACS CFD air quality model performance evaluation with kit fox, MUST, prairie grass, and EMU observations, *Atmos. Environ.* 38 (28) (2004) 4675–4687, <http://dx.doi.org/10.1016/j.atmosenv.2004.05.041>.
- [59] K.-J. Hsieh, F.-S. Lien, E. Yee, Numerical modeling of passive scalar dispersion in an urban canopy layer, *J. Wind Eng. Ind. Aerodyn.* 95 (12) (2007) 1611–1636, <http://dx.doi.org/10.1016/j.jweia.2007.02.028>.
- [60] M. Milliez, B. Carissimo, Numerical simulations of pollutant dispersion in an idealized urban area, for different meteorological conditions, *Bound.-Layer Meteorol.* 122 (2) (2007) 321–342, <http://dx.doi.org/10.1007/s10546-006-9110-4>.
- [61] R. Donnelly, T. Lyons, T. Flassak, Evaluation of results of a numerical simulation of dispersion in an idealised urban area for emergency response modelling, *Atmos. Environ.* 43 (29) (2009) 4416–4423, <http://dx.doi.org/10.1016/j.atmosenv.2009.05.038>.
- [62] G.C. Efthimiou, J.G. Bartzis, N. Koutsourakis, Modelling concentration fluctuations and individual exposure in complex urban environments, *J. Wind Eng. Ind. Aerodyn.* 99 (4) (2011) 349–356, <http://dx.doi.org/10.1016/j.jweia.2010.12.007>, The Fifth International Symposium on Computational Wind Engineering.
- [63] P. Kumar, A.-A. Feiz, P. Ngai, S.K. Singh, J.-P. Issartel, CFD simulation of short-range plume dispersion from a point release in an urban like environment, *Atmos. Environ.* 122 (2015) 645–656, <http://dx.doi.org/10.1016/j.atmosenv.2015.10.027>.
- [64] M.L. Bahlali, E. Dupont, B. Carissimo, Atmospheric dispersion using a Lagrangian stochastic approach: Application to an idealized urban area under neutral and stable meteorological conditions, *J. Wind Eng. Ind. Aerodyn.* 193 (2019) 103976, <http://dx.doi.org/10.1016/j.jweia.2019.103976>.
- [65] F. Camelli, R. Lohner, S. Hanna, VLES study of MUST experiment, in: 43rd AIAA Aerospace Sciences Meeting and Exhibit, 2005, <http://dx.doi.org/10.2514/6.2005-1279>.
- [66] M. König, Large-eddy simulation modelling for urban scale (Ph.D. thesis), University of Leipzig, 2014, URL <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=baab9d7b41623099c1b6d840c11821b8e31fac9b>, (Accessed 05 August 2024).
- [67] T. Nagel, R. Schoetter, V. Masson, C. Lac, B. Carissimo, Numerical analysis of the atmospheric boundary-layer turbulence influence on microscale transport of pollutant in an idealized urban environment, *Bound.-Layer Meteorol.* 184 (1) (2022) 113–141, <http://dx.doi.org/10.1007/s10546-022-00697-7>.
- [68] T. Schönfeld, M. Rudgyard, Steady and unsteady flow simulations using the hybrid flow solver AVBP, *AIAA J.* 37 (11) (1999) 1378–1385, <http://dx.doi.org/10.2514/2.636>.
- [69] L.Y. Gicquel, N. Gourdain, J.-F. Boussuge, H. Deniau, G. Staffelbach, P. Wolf, T. Poinsot, High performance parallel computing of flows in complex geometries, *Comptes Rendus Mécanique* 339 (2) (2011) 104–124, <http://dx.doi.org/10.1016/j.crme.2010.11.006>, High Performance Computing.
- [70] J. Ramshaw, P. O'Rourke, A. Amsden, Acoustic damping for explicit calculations of fluid flow at low Mach number, Technical report no. LA-10641-MS, Los Alamos National Laboratories, USA, 1986, URL https://inis.iaea.org/collection/NCLCollectionStore/_Public/17/074/17074782.pdf, (Accessed 01 December 2023).
- [71] F. Nicoud, F. Ducros, Subgrid-scale stress modelling based on the square of the velocity gradient tensor, *Flow Turbul. Combust.* 62 (3) (1999) 183–200, <http://dx.doi.org/10.1023/A:1009995426001>.
- [72] A. Smirnov, S. Shi, I. Celik, Random flow generation technique for large eddy simulations and particle-dynamics modeling, *J. Fluids Eng.* 123 (2) (2001) 359–371, <http://dx.doi.org/10.1115/1.1369598>.
- [73] R. Vasaturo, I. Kalkman, B. Blocken, P. van Wesemael, Large eddy simulation of the neutral atmospheric boundary layer: Performance evaluation of three inflow methods for terrains with different roughness, *J. Wind Eng. Ind. Aerodyn.* 173 (2018) 241–261, <http://dx.doi.org/10.1016/j.jweia.2017.11.025>.
- [74] G. Lamberti, C. Gorié, A multi-fidelity machine learning framework to predict wind loads on buildings, *J. Wind Eng. Ind. Aerodyn.* 214 (2021) 104647, <http://dx.doi.org/10.1016/j.jweia.2021.104647>.
- [75] J.H. Halton, Algorithm 247: Radical-inverse quasi-random point sequence, *Commun. ACM* 7 (12) (1964) 701–702, <http://dx.doi.org/10.1145/355588.365104>.
- [76] E. Lumet, T. Jaravel, M.C. Rochoux, PPMLES – Perturbed-Parameter ensemble of MUST Large-Eddy Simulations, 2024, <http://dx.doi.org/10.5281/zenodo.11394347>.
- [77] M. Schatzmann, H. Olesen, J. Franke, COST 732 model evaluation case studies: approach and results, Technical report, University of Hamburg, Meteorological Institute, 2010, URL https://www.researchgate.net/profile/George-Efthimiou-3/post/Halo-fluent-been-compared-to-starcem/attachment/59d6585379197b80779ae4bd/AS%3A538043318628353%401505290931380/download/5th_Docu_May_10.pdf, (Accessed 01 December 2023).
- [78] D.N. Politis, J.P. Romano, The stationary bootstrap, *J. Amer. Statist. Assoc.* 89 (428) (1994) 1303–1313, <http://dx.doi.org/10.1080/01621459.1994.10476870>.
- [79] M. Guo, J.S. Hesthaven, Reduced order modeling for nonlinear structural analysis using Gaussian process regression, *Comput. Methods Appl. Mech. Engrg.* 341 (2018) 807–826, <http://dx.doi.org/10.1016/j.cma.2018.07.017>.
- [80] F. Chinesta, P. Ladeveze, E. Cueto, A short review on model order reduction based on proper generalized decomposition, *Arch. Comput. Methods Eng.* 18 (4) (2011) 395–404, <http://dx.doi.org/10.1007/s11831-011-9064-7>.
- [81] K. Taira, S.L. Brunton, S.T.M. Dawson, C.W. Rowley, T. Colonius, B.J. McKeon, O.T. Schmidt, S. Gordeyev, V. Theofilis, L.S. Ukeiley, Modal analysis of fluid flows: An overview, *AIAA J.* 55 (12) (2017) 4013–4041, <http://dx.doi.org/10.2514/1.J056060>.
- [82] A. Kessy, A. Lewin, K. Strimmer, Optimal whitening and decorrelation, *Amer. Statist.* 72 (4) (2018) 309–314, <http://dx.doi.org/10.1080/00031305.2016.1277159>.
- [83] L. Cordier, M. Bergmann, Réduction de dynamique par décomposition orthogonale aux valeurs propres (POD) (in French). Lecture notes, 7563, Ecole de printemps OCET, 2006, p. 107, URL <https://www.math.u-bordeaux.fr/~mbergman/PDF/OuvrageSynthese/OCET06.pdf>, (Accessed 01 December 2023).
- [84] O.T. Schmidt, T. Colonius, Guide to spectral proper orthogonal decomposition, *AIAA J.* 58 (3) (2020) 1023–1033, <http://dx.doi.org/10.2514/1.J058809>.
- [85] M. Cassiani, M.B. Bertagni, M. Marro, P. Salizzoni, Concentration fluctuations from localized atmospheric releases, *Bound.-Layer Meteorol.* 177 (2) (2020) 461–510, <http://dx.doi.org/10.1007/s10546-020-00547-4>.
- [86] S.L. Brunton, J.N. Kutz, *Data-Driven Science and Engineering: Machine Learning, Dynamical Systems, and Control*, Cambridge University Press, 2019, <http://dx.doi.org/10.1017/9781108380690>.
- [87] B.X. Nony, Reduced-order models under uncertainties for microscale atmospheric pollutant dispersion in urban areas: exploring learning algorithms for high-fidelity model emulation (Ph.D. thesis), Université de Toulouse, France, 2023, URL <https://theses.fr/2023TOU30156>, (Accessed 10 October 2024).
- [88] M.L. Stein, Interpolation of spatial data: some theory for kriging, in: *Springer Series in Statistics*, Springer Science & Business Media, 1999, <http://dx.doi.org/10.1007/978-1-4612-1494-6>.
- [89] T. Hastie, R. Tibshirani, J.H. Friedman, J.H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*, vol. 2, Springer, 2009, <http://dx.doi.org/10.1007/978-0-387-21606-5>.

- [90] J. Forkman, J. Josse, H.-P. Piepho, Hypothesis tests for principal component analysis when variables are standardized, *J. Agric. Biol. Environ. Stat.* 24 (2019) 289–308, <http://dx.doi.org/10.1007/s13253-019-00355-5>.
- [91] J. Chang, S. Hanna, Air quality model performance evaluation, *Meteorol. Atmos. Phys.* 87 (1) (2004) 167–196, <http://dx.doi.org/10.1007/s00703-003-0070-7>.
- [92] R. Miyagusuku, A. Yamashita, H. Asama, Gaussian processes with input-dependent noise variance for wireless signal strength-based localization, in: 2015 IEEE International Symposium on Safety, Security, and Rescue Robotics, SSR, 2015, pp. 1–6, <http://dx.doi.org/10.1109/SSRR.2015.7442993>.
- [93] C.L. Defforge, B. Carissimo, M. Bocquet, R. Bresson, P. Armand, Improving CFD atmospheric simulations at local scale for wind resource assessment using the iterative ensemble Kalman smoother, *J. Wind Eng. Ind. Aerodyn.* 189 (2019) 243–257, <http://dx.doi.org/10.1016/j.jweia.2019.03.030>.
- [94] C.L. Defforge, B. Carissimo, M. Bocquet, R. Bresson, P. Armand, Improving numerical dispersion modelling in built environments with data assimilation using the iterative ensemble Kalman smoother, *Bound.-Layer Meteorol.* 179 (2) (2021) 209–240, <http://dx.doi.org/10.1007/s10546-020-00588-9>.
- [95] S. El Garroussi, S. Ricci, M. De Lozzo, N. Goutal, D. Lucor, Assessing uncertainties in flood forecasts using a mixture of generalized polynomial chaos expansions, in: 2020 TELEMAC-MASCARET User Conference, 2020, URL <https://hal.science/hal-03444227/document>, (Accessed 01 December 2023).
- [96] T. Murata, K. Fukami, K. Fukagata, Nonlinear mode decomposition with convolutional neural networks for fluid dynamics, *J. Fluid Mech.* 882 (2020) A13, <http://dx.doi.org/10.1017/jfm.2019.822>.
- [97] V. Picheny, D. Ginsbourger, O. Roustant, R.T. Haftka, N.-H. Kim, Adaptive designs of experiments for accurate approximation of a target region, *J. Mech. Des.* 132 (7) (2010) 071008, <http://dx.doi.org/10.1115/1.4001873>.
- [98] T. Braconnier, M. Ferrier, J.-C. Jouhaud, M. Montagnac, P. Sagaut, Towards an adaptive POD/SVD surrogate model for aeronautic design, *Comput. & Fluids* 40 (1) (2011) 195–209, <http://dx.doi.org/10.1016/j.compfluid.2010.09.002>.
- [99] M. Shirzadi, Y. Tominaga, Multi-fidelity shape optimization methodology for pedestrian-level wind environment, *Build. Environ.* 204 (2021) 108076, <http://dx.doi.org/10.1016/j.buildenv.2021.108076>.
- [100] B.X. Nony, M.C. Rochoux, T. Jaravel, D. Lucor, Reduced-order model for microscale atmospheric dispersion combining multi-fidelity LES and RANS data, in: ECCOMAS Proceedings, 2023, <http://dx.doi.org/10.7712/120223.10337.19817>, ID 10337, pp. 265–283, presented at 5th International Conference on Uncertainty Quantification in Computational Science and Engineering, UNCECOMP Congress 2023, Athens (Greece).