

Partitioning, ordering, and load balancing in a hierarchically parallel hybrid linear solver

Ichitaro Yamazaki

University of Tennessee, Knoxville

<u>Xiaoye Li</u>

Lawrence Berkeley National Laboratory

Francois-Henry Rouet

ENSEEIHT-IRIT and LIP

Bora Ucar

ENS-Lyon

Sparse Days, June 25th-26th, 2012, CERFACS, Toulouse

OUTLINE OF THE TALK



Overview of DD Schur complement method

- Need for hierarchical parallelism
- PDSLin (Parallel Domain decomposition Schur complement Linear solver)

Download: http://crd-legacy.lbl.gov/FASTMath-LBNL/Software/

• Combinatorial problems in hybrid sparse solvers

- K-way graph partitioning with multiple constraints
- Sparse triangular solution with sparse right-hand sides



Schur complement method

- a.k.a iterative substructuring method
 - or, non-overlapping domain decomposition
- Divide-and-conquer paradigm . . .
 - Divide entire problem (domain, graph) into subproblems (subdomains, subgraphs)
 - Solve the subproblems
 - Solve the interface problem (Schur complement)
- Variety of ways to solve subdomain problems and Schur complement ... lead to powerful poly-algorithm or hybrid solver framework

Algebraic view



1. Reorder into 2x2 block system, A₁₁ is block diagonal

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

2. Schur complement

$$S = A_{22} - A_{21} A_{11}^{-1} A_{12} = A_{22} - (U_{11}^{-T} A_{21}^{T})^{T} (L_{11}^{-1} A_{12}) = A_{22} - W \cdot G$$

where $A_{11} = L_{11}U_{11}$

S corresponds to interface (separator) variables no need to form explicitly

3. Compute the solution

(1)
$$x_2 = S^{-1}(b_2 - A_{21} A_{11}^{-1} b_1) \quad \leftarrow \text{ iterative solver}$$

(2) $x_1 = A_{11}^{-1}(b_1 - A_{12} x_2) \quad \leftarrow \text{ direct solver}$

Parallelism – extraction of many subdomains



- Partition adjacency graph of |A|+|A^T|
 Multiple goals: reduce size of separator, balance size of subdomains
 - nested dissection (e.g., PT-Scotch, ParMetis)
 - k-way partition (preferred)



Hierarchical parallelism



Multiple processors per subdomain

 E.g. one subdomain with 2x3 procs (e.g. SuperLU_DIST, MUMPS)



Load balance, communication

- Distinction between intra-group and inter-group
- Intra-group relies on SuperLU_DIST, inter-group new design

PDSLin software



- Download: http://crd-legacy.lbl.gov/FASTMath-LBNL/Software/
 - C and MPI, with Fortran interface.
 - Unsymmetric and symmetric, real and complex, multiple RHSs.
- Requires following external packages:
 - parallel graph partitioning:
 - PT-Scotch (http://labri.fr/perso/pelegrin/scotch), or
 - ParMetis (<u>http://glaros.dtc.umn.edu/gkhome/views/metis</u>)
 - subdomain solver options:
 - SuperLU (http://crd.lbl.gov/~xiaoye/SuperLU)
 - SuperLU_DIST or SuperLU_MT: two levels of parallelization
 - MUMPS
 - PDSLin
 - Inexact subdomain solution: ILU
 - Schur complement solver options:
 - PETSc (<u>http://mcs.anl.gov/petsc/petsc-as</u>)
 - SuperLU_DIST

PDSLin encompass Hybrid, Iterative, Direct





Combinatorial problems . . .



• K-way graph partitioning with multiple objectives

- Small separator
- Similar subdomains
- Similar connectivity

Sparse triangular solution with many sparse RHS

$$S = A_{22} - \sum_{l} (U_{l}^{-T} F_{l}^{T})^{T} (L_{l}^{-1} E_{l}) = A_{22} - \sum_{l} G \cdot W, \text{ where } D_{l} = L_{l} U_{l}$$

• Sparse matrix-matrix multiplication $\widetilde{G} \leftarrow \text{sparsify}(G, \sigma_1); \quad \widetilde{W} \leftarrow \text{sparsify}(W, \sigma_1)$ $T^{(p)} \leftarrow \widetilde{W}^{(p)} \times \widetilde{G}^{(p)}$ $\hat{S}^{(p)} \leftarrow A_{22}^{(p)} - \sum_q T^{(q)}(p); \quad \widetilde{S} \leftarrow \text{sparsify}(\hat{S}, \sigma_2)$

Two graph models



- Standard graph : G=(V, E)
 - GPVS: graph partitioning with vertex separator
 - GPES: graph partitioning with edge separator
- Hypergraph : H = (V, N), net = "edge", may include more than two vertices
 - Column-net hypergraph: H = (R, C) rows = vertices, columns = nets n₃ = {1, 3, 4}
 - Row-net hypergraph: H = (C, R) columns = vertices, rows = nets



- Partition problem: $\pi(V) = \{V_1, V_2, \dots, V_k\}$, disjoint parts
 - Graph: a cut edge connects two parts
 - Hypergraph: a cut net connects multiple parts
 - ➔ Want to minimize the cutsize in some metric (Objective), and keep equal weights among the parts (Constraints).

1. K-way subdomain extraction

Problem with ND:

Imbalance in separator size at different branches → Imbalance in subdomain size

- Alternative: directly partition into K parts, meet multiple constraints:
 - 1. Compute k-way partitioning → balanced subdomains
 - 2. Extract separator → balanced connectivity







k-way partition: objectives, constraints





Initial partition to extract vertex-separator impacts load balance:

Objectives to minimize:

- number of interface vertices
 → separator size
- number of interface edges
 → nonzeros in interface

Balance constraints:

- number of interior vertices and edges → LU of D_i
- number of interface vertices and edges → local update matrix

 $(U_l^{-T}F_l^T)^T (L_l^{-1}E_l)$

K-way edge partition



• Extract vertex-separator from k-way edge partition

- Compute k-way edge partition satisfying balanced subdomains (e.g., PT-Scotch, Metis)
- Extract vertex-separator from edge-separators to minimize and balance interfaces (i.e. minimum vertex cover)
- Heuristics to pick the next vertex: pick a vertex from largest subdomain to maintain balance among subdomains
 - pick the vertex with largest degree to minimize separator size
 - pick the vertex to obtain best balance of interfaces (e.g., nnz).



Balance & Solution time with edge part. + VC



- tdr190k from accelerator cavity design: N = 1.1M, k = 32
- Compared to ND of SCOTCH
- balance = max value_e / min value_e, for l = 1, 2, ..., k
- Results
 - Improved balance of subdomains, but not of interfaces due to larger separator → total time not improved
 - Post-processing already-computed partition is not effective to balance multiple constraints



Recursive Hypergraph Bisection



 Column-net HG partition to permute an m-by-n matrix M to a singly-bordered form (e.g., Patoh):

$$P_{r}MP_{c}^{T} = \begin{pmatrix} M_{1} & & C_{1} \\ M_{2} & & C_{2} \\ & \ddots & & \vdots \\ & & M_{k} & C_{k} \end{pmatrix} \rightarrow P_{c}(M^{T}M)P_{c}^{T} = \begin{pmatrix} M_{1}^{T}M_{1} & & M_{1}^{T}C_{1} \\ & M_{2}^{T}M_{2} & & M_{2}^{T}C_{2} \\ & & \ddots & & \vdots \\ & & & M_{k}^{T}M_{k} & M_{k}^{T}C_{k} \\ & & & C_{1}^{T}M_{1} & C_{2}^{T}M_{2} & \cdots & C_{k}^{T}M_{k} & \sum C_{l}^{T}C_{l} \end{pmatrix}$$

• Recursive hypergraph bisection (RHB)

- Compute structural decomposition str(A) = str(M^TM)
 - e.g., using edge clique cover of G(A) [Catalyurek '09]
- Compute 2-way partition of M (with multiple constraints) into a singly-bordered form based on recursive bisection
- Permute A as:

as.

$$\operatorname{str}(\mathbf{P}_{\mathbf{c}}\mathbf{A}\mathbf{P}_{\mathbf{c}}^{\mathrm{T}}) = \begin{pmatrix} D_{1} & E_{1} \\ D_{2} & E_{2} \\ F_{1} & F_{2} & C \end{pmatrix} = \begin{pmatrix} M_{1}^{T} & \\ & M_{2}^{T} \\ C & C_{2}^{T} \end{pmatrix} \bullet \begin{pmatrix} M_{1} & C_{1} \\ & M_{2} & C_{2} \end{pmatrix}$$

$$15$$

RHB: objectives, constraints



(interface nz columns)

<u>Objective</u>: minimize cutsize cutsize metrics:

- Connectivity 1, $\sum_{n \in M} (\lambda(j) 1)$
- Cut-net, $\sum_{n \in N, \lambda(j) > 1} 1$ (separator size)
- Sum-of-external degree (soed), $\sum_{j \in N, \lambda(j)=1} \lambda(j)$ (sum of above)

Where, j-th net n_j is in the cut, and λ_j is the number of parts n_j is connected to.

- <u>Constraints</u>: equal weights of different parts
 i-th vertex weights (based on previous partition):
 - unit weight (subdo
 - Nnz(M_k(i, :))

- (subdomain dimension)
- (subdomain nnz)
- $Nnz(M_k(i, :)) + nnz(C_k(i, :))$ (interface nnz)

Balance & Time results with RHB



- tdr190k from accelerator cavity design: N = 1.1M, k = 32•
- Compared to ND of SCOTCH
- Results
 - Single-constraint improves balance without much increase of separator size \rightarrow 1.7x faster than ND
 - Multi-constraints improves balance, but larger separator



tdr190k with 32 subdomains, single-constraint with nnz weights

2. Sparse triangular solution with sparse RHS (intra-group within a subdomain)



RHS vectors E_e and F_e are sparse (e.g., about 20 nnz per column); There are many RHS vectors (e.g., O(10⁴) columns)



Blocking RHS vectors

- Reduce number of calls to the symbolic routine and number of messages, and improve read reuse of the LU factors
- Achieved over 5x speedup
- **x** zeros must be padded to fill the block \rightarrow memory cost !

Sparse triangular solution with sparse RHSs



- Objective: Reorder columns of E_e to maximize structural similarity among the adjacent columns.
- Where are the fill-ins?

<u>Path Theorem [Gilbert'94]</u> Given the elimination tree of D_{I_i} fills generated in G_I at the positions associated with nodes on the path from nodes of the nonzeros in E_I to the root



Sparse rhs ... postordering



• Postorder-conforming ordering of the RHS vectors

- Postorder elimination tree of D_I
- Permute columns of E₁ s.t. row indices of the first nonzeros are in ascending order
- Increased overlap of the paths to the root, fewer padded zeros
- 20–40% reduction in solution time over ND



Sparse triangular solution ... Hypergraph



- Partition/group columns using row-net HG
- Define a cost function ≈ padded zeros

$$\cos t(\Pi) = \sum_{\text{row } i=1}^{n} \left(\lambda_{i} B - nnz(G_{l}(i,:)) \right) \quad \text{row i}$$

$$= \sum_{i=1}^{n} \left(\lambda_{i} - 1 \right) B \quad + \quad \sum_{i=1}^{n} B + nnz(G)$$

R

"connectivity-1" metric con

constant

- Minimize cost(π) using Patoh
- Additional 10% reduction in time

Sparse RHS: memory saving



- tdr190k from accelerator cavity design: N = 1.1M, k = 8
- Fraction of padded zeros, with different block size



Summary



Graph partitioning

- Direct edge partition + min. vertex cover not effective
- Recursive Hypergraph Bisection: 1.7x faster
- Reordering sparse r.h.s. in sparse triangular solution
 - postordering: 20-40% faster; hypergraph: additional 10% faster

• Future work

Parallel hypergraph partitioning



Extra Slides

Parallelization with serial subdomain



- No. of subdomains increases with increasing core count.
 Schur complement size and iteration count increase
- HIPS (serial subdomain) vs. PDSLin (parallel subdomain)
 - M3D-C¹, Extended MHD to model fusion reactor tokamak, 2D slice of 3D torus
 - Dimension 801k, 70 nonzeros per row, real unsymmetric

Р	N _s	HIPS 1.0 sec (iter)	PDSLin sec (iter)
8	13k	284.6 (26)	79.9 (15)
32	29k	55.4 (64)	25.3 (16)
128	62k		17.1 (16)
512	124k		21.9 (17)

PDSLin for largest system



- Matrix properties:
 - 3D cavity design in Omega3P, 3rd order basis function for each matrix element
 - dimension = 52.7 M, nonzeros = 4.3 B (~80 nonzeros per row), real, symmetric, highly indefinite
- Experimental setup:
 - PT-Scotch extracts 128 subdomains of size ~410k
 - SuperLU DIST factors each subdomain, computes preconditioner $LU(\tilde{S})$ of size 247k (32 cores)
 - BiCGStab of PETSc to solve Sy = d
- Performance:
 - Fill-ratio (nnz(Precond.)/nnz(A)): ~ 250
 - Using 2048 cores:
 - preconditioner construction: 493.1 sec.
 - solution: 108.1 second (32 iterations)

Application 2: Accelerator cavity design



- DOE SciDAC: Community Petascale Project for Accelerator Science and Simulation (ComPASS), PI: P. Spentzouris, Fermilab
- Development of a comprehensive computational infrastructure for accelerator modeling and optimization
- RF cavity: Maxwell equations in electromagnetic field
- FEM in frequency domain leads to large sparse eigenvalue problem; needs to solve shifted linear systems

