

Dissection sparse direct solver and parallel task management

Atsushi Suzuki Cybermedia Center, Osaka University

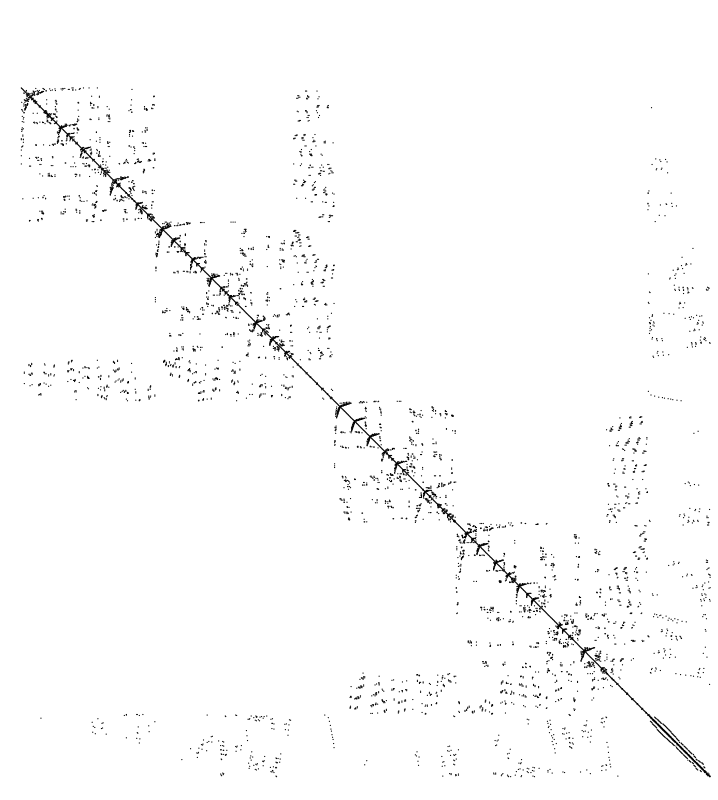
For numerical simulation of partial differential equations (PDE), we need to solve linear systems with symmetric or unsymmetric large sparse matrices obtained by discretization by finite difference, finite volume and finite element methods. Number of unknowns is more than one million and condition number of the large sparse matrix is more than 10^6 , due to large variety of physical coefficients and/or coupling of different physics. Sometimes direct solver could be only a possible tool to find solution of such difficult linear system. There are several sparse direct solvers for parallel computational environments, e.g., SuperLU_MT, Pardiso, SuperLU_DIST, and MUMPS. The first two codes run on shared memory systems and others run on distributed memory system. Among them, a sparse solver Dissection, we have first developed it for symmetric matrix in DOI:10.1002/nme.4729 on shared memory architecture, with keeping numerical stability by employing robust pivoting technique. Now the solver can factorize structurally symmetric matrix, which means nonzero pattern of the matrix is symmetric, in both double and quadruple precision arithmetic thanks to modern C++ implementation.

For parallel computation on shared memory system, tasks of local LDU-factorization for sparse and dense sub-matrices that are generated by nested-dissection ordering are assigned to available cores asynchronously. In precise, dependency of tasks for block factorization is analyzed and expressed as a set of directed acyclic graph (DAG), and quasi static assignment of tasks by considering complexity depending on various problem sizes due to sparsity of the matrix, with additional dynamic assignment to recover from rough estimation of complexity and execution noise.

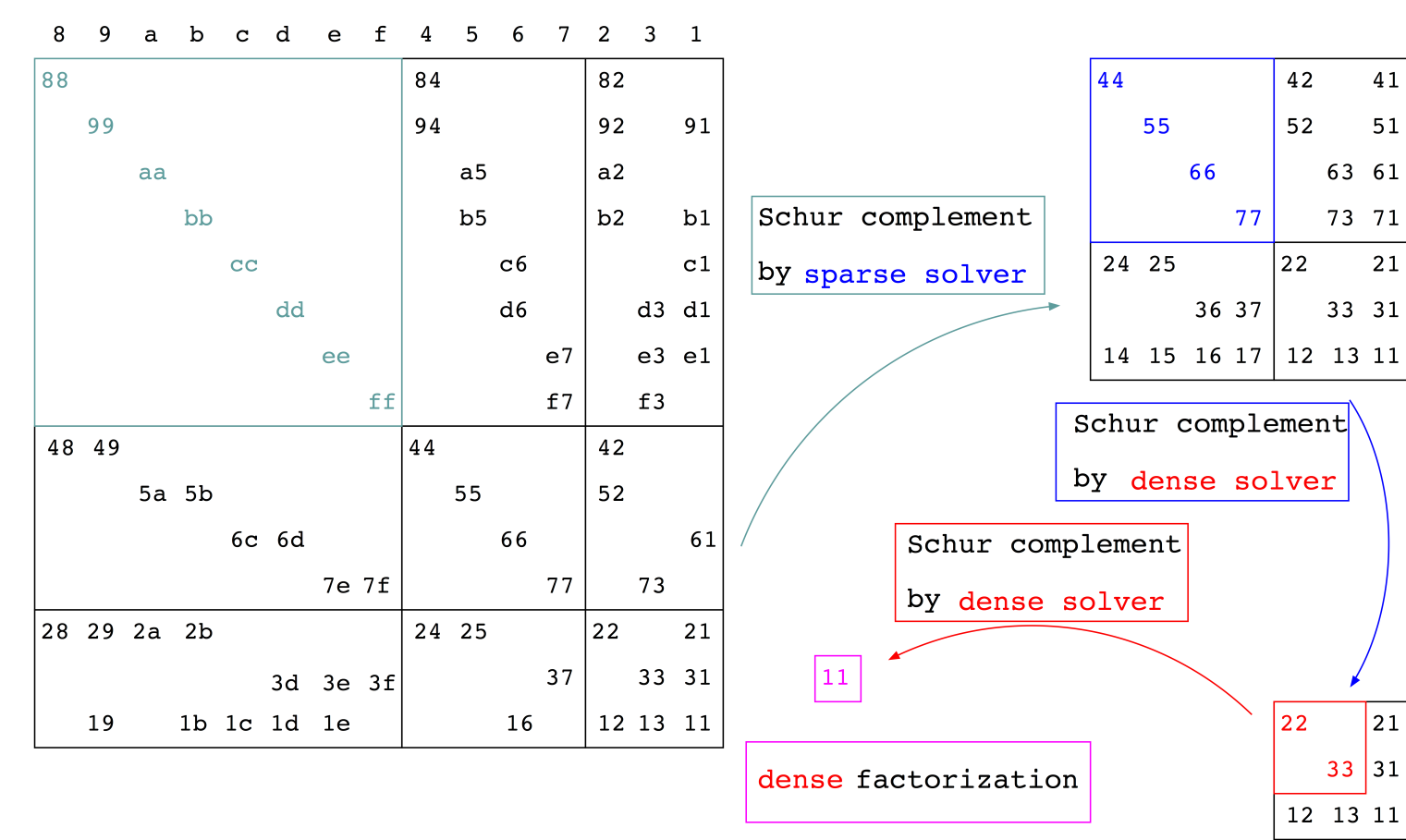
For numerical stability of LDU-factorization, threshold pivoting technique is employed, where additional Schur complement concerning on postponed pivots is build and factorized at the end of elimination procedure. This approach works well numerically for sparse matrix from PDE thanks to nested-dissection ordering which is originated by domain decomposition strategy and archives high parallel efficiency thanks to static structure of DAG with only one task added at the end of the elimination tree.

Dissection software is developed with F.-X. Roux @ ONERA/LJLL UPMC, France and used in FreeFEM @ LJLL UPMC, France, hpddm @ ENSEEIHT, France under GPL, and licensed to KIOXIA, Japan under CeCILL-C.

nested-dissection ordering of sparse matrix

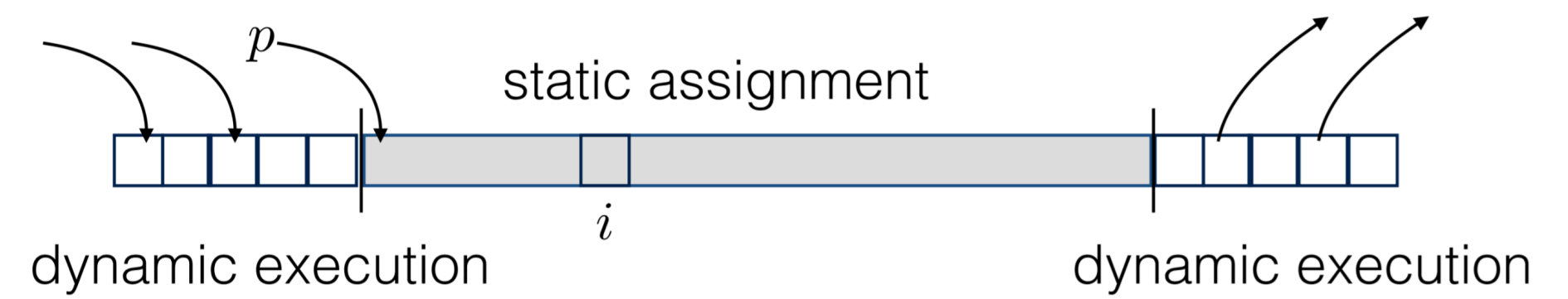


7 stencil of Poisson equation in 3D,
11³ nodes : nested dissection ordering
by SCOTCH



multi-frontal LDU-factorization by recursive computation
of Schur complements

task execution by static + dynamic scheduling



- ▶ dependency of tasks is determined after symbolic analysis
- ▶ estimation of complexity of task based on block size
- ▶ all tasks are described though some tasks are not executed due to numerical pivoting
- ▶ for atomic operation to update value of the index for task, mutex is not necessary
- ▶ on-the-fly measurement of task could improve complexity estimation

details of task scheduling

- ▶ $s[i]$ $1 \leq i \leq N$ tasks in the critical path
 - ▶ $d[j]$ $1 \leq j \leq M$ other tasks, independent of $s[i]$
 - ▶ θ : ratio of static scheduling, $n = \theta N$
 - ▶ $1 \leq p \leq P$: processor id
- $i = 1, j = 1$ before arrival of processes.
while (all processes arrive and $i \leq N$) {
 while (parents of $s[i]$ are not finished) {
 verify parents of $d[j]$ are finished.
 if finished then increase index j and execute $d[j-1]$: **dynamic**
 otherwise sleep until receive a wake-up signal.
 }
 increase index i and execute $s[i-1]$: **dynamic**
}
if (p is the last arrived process) {
 divide tasks $s[i], \dots, s[n]$ into P groups $\{b_1, \dots, b_P\}$ with
 $i = b_1 < b_2 < \dots < b_P < n$, where $\sum_{b_q \leq k < b_{q+1}} [\text{complexity of } s[k]]$ are
 homogeneous for all $1 \leq q \leq P$.
 set $i = n$.
 execute $s[k]$ for $b_p \leq k < b_{p+1}$ without checking status of parents : **static**
 while ($i \leq N$) {
 increase index i and execute $s[i-1]$: **dynamic greedy**
 }
 mutex is used to increase global variable i , cond_wait/cond_broadcast

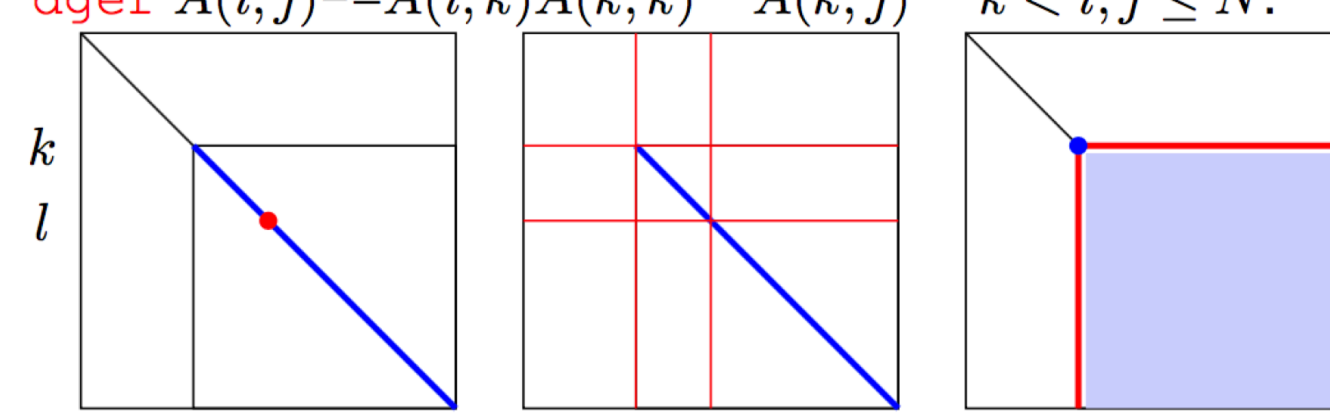
LDU-factorization and local symmetric pivoting

$$\begin{bmatrix} A_{11} & A_{21} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & 0 \\ A_{21} & S_{22} \end{bmatrix} \begin{bmatrix} I_1 & A_{11}^{-1}A_{12} \\ 0 & I_2 \end{bmatrix}$$

$$\begin{aligned} S_{22} &= A_{22} - A_{21}A_{11}^{-1}A_{12} \\ &= A_{22} - A_{21}(L_{11}D_{11}U_{11})^{-1}A_{12} \\ &= A_{22} - (U_{11}^{-T}A_{21}^T)(D_{11}^{-1}L_{11}^{-1}A_{12}) \end{aligned}$$

block LDU-factorization by forward substitution
of multiple-RHS and updating Schur complement

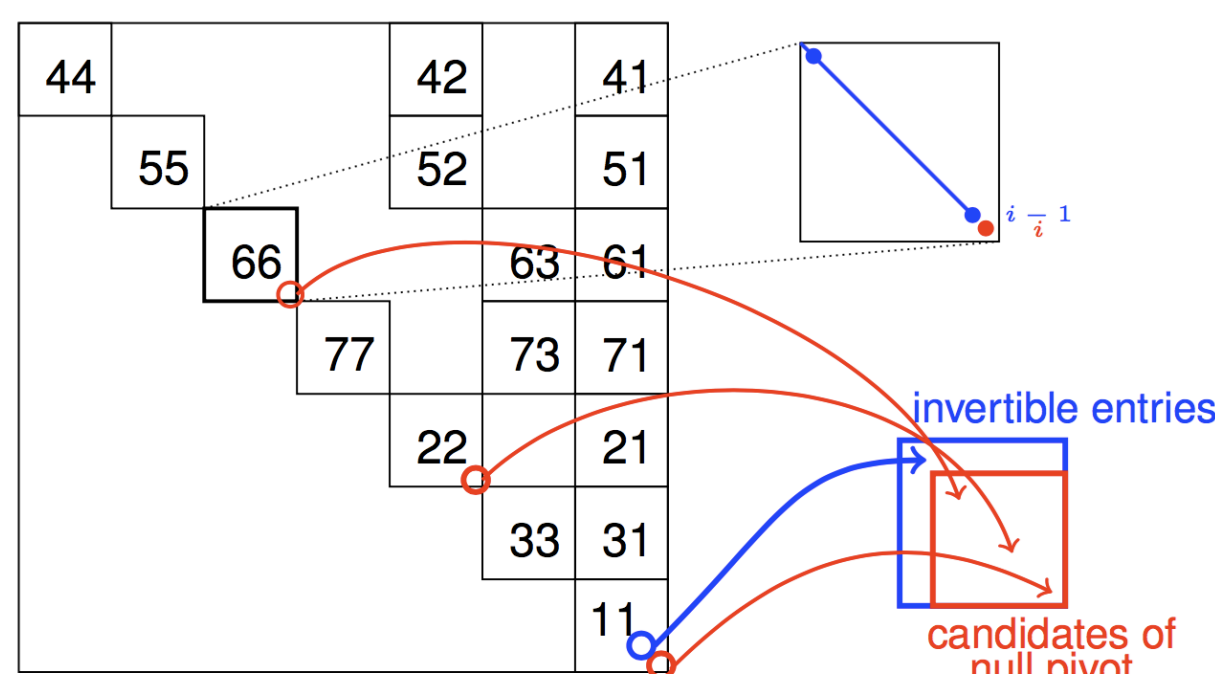
do $k = 1, \dots, N$
 find $k < l \leq n$ $\max |A(l, k)|$,
 exchange rows and columns : $A(k, *) \leftrightarrow A(l, *)$, $A(*, k) \leftrightarrow A(*, l)$.
 dscal $A(k, j) / A(k, k)$ $k < j \leq N$,
 dscal $A(i, k) / A(k, k)$ $k < i \leq N$,
 dger $A(i, j) - A(i, k)A(k, k)^{-1}A(k, j)$ $k < i, j \leq N$.



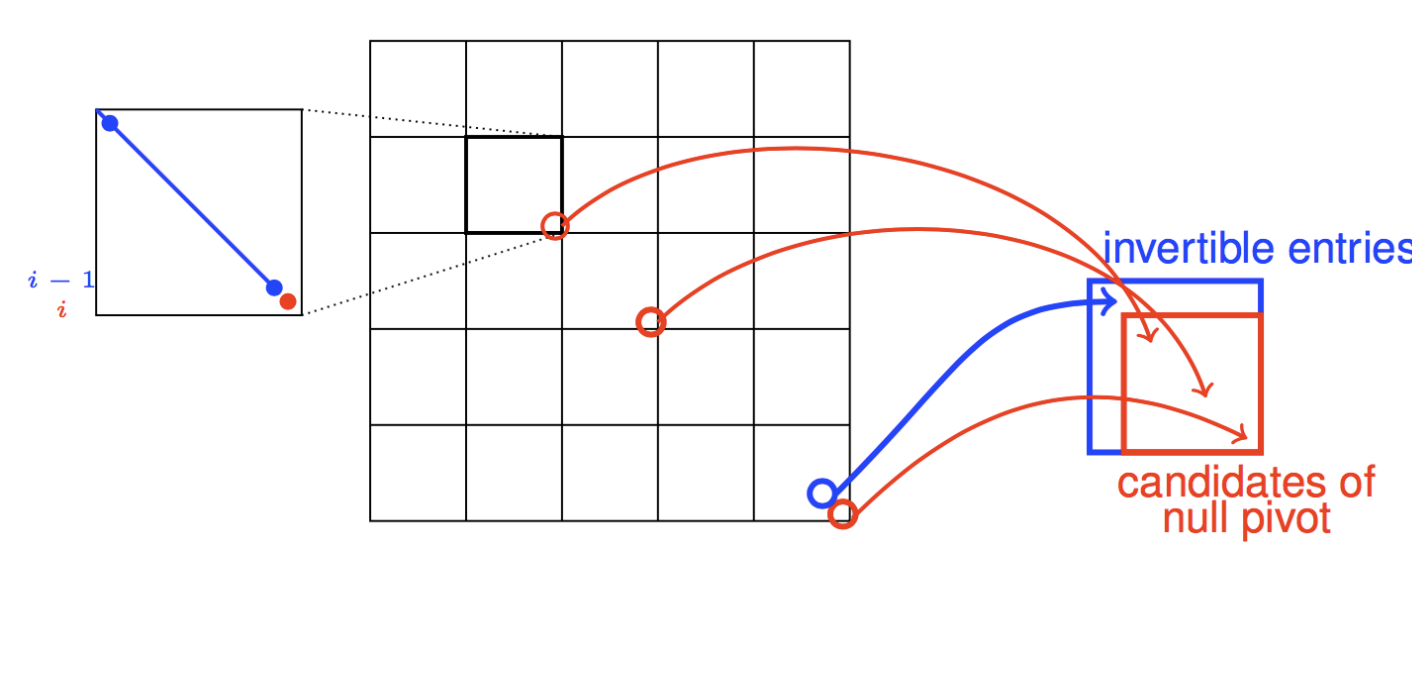
symmetric pivoting in a block

multi-frontal and block factorization with threshold pivoting

τ : given threshold for null pivot $\approx 10^{-2}$
 $|A(i, i)| / |A(i-1, i-1)| < \tau \Rightarrow |A(i, i)|$ is null pivot.

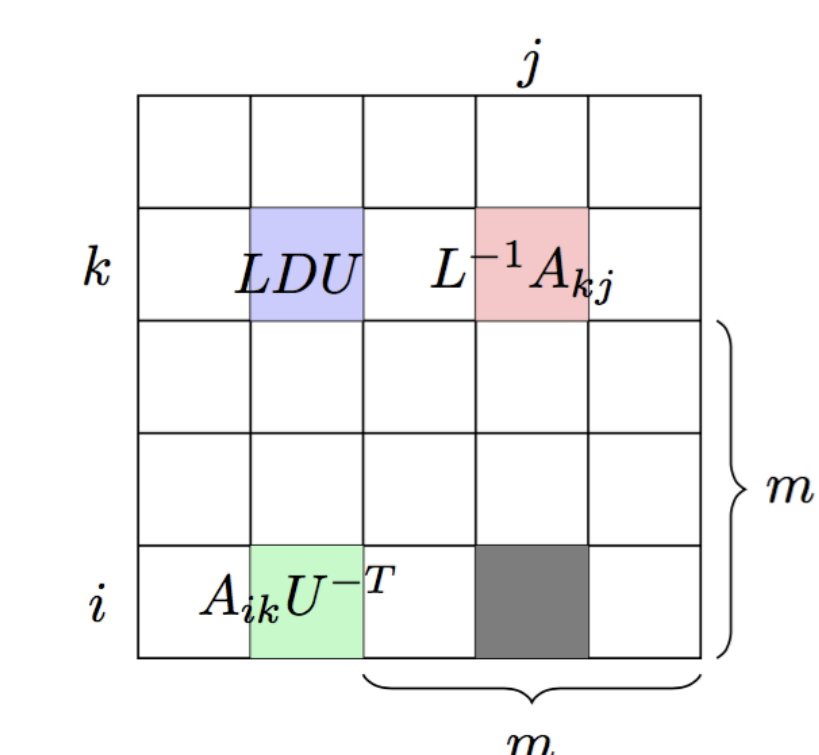


postponed pivots in multi-frontal structure

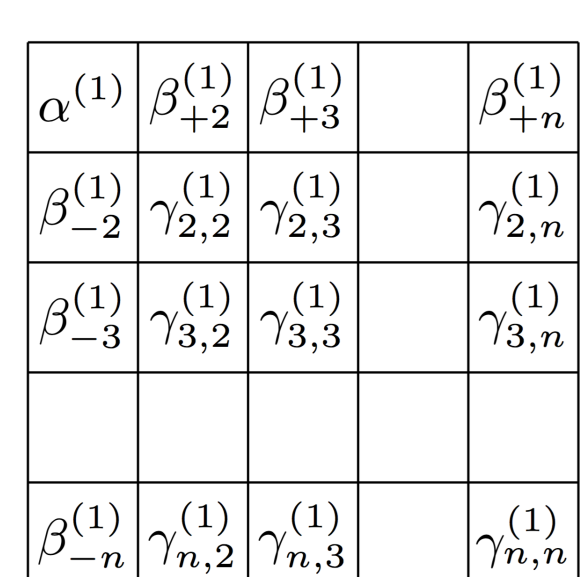


postponed pivots in block dense structure

tasks in block factorization and dependency analysis



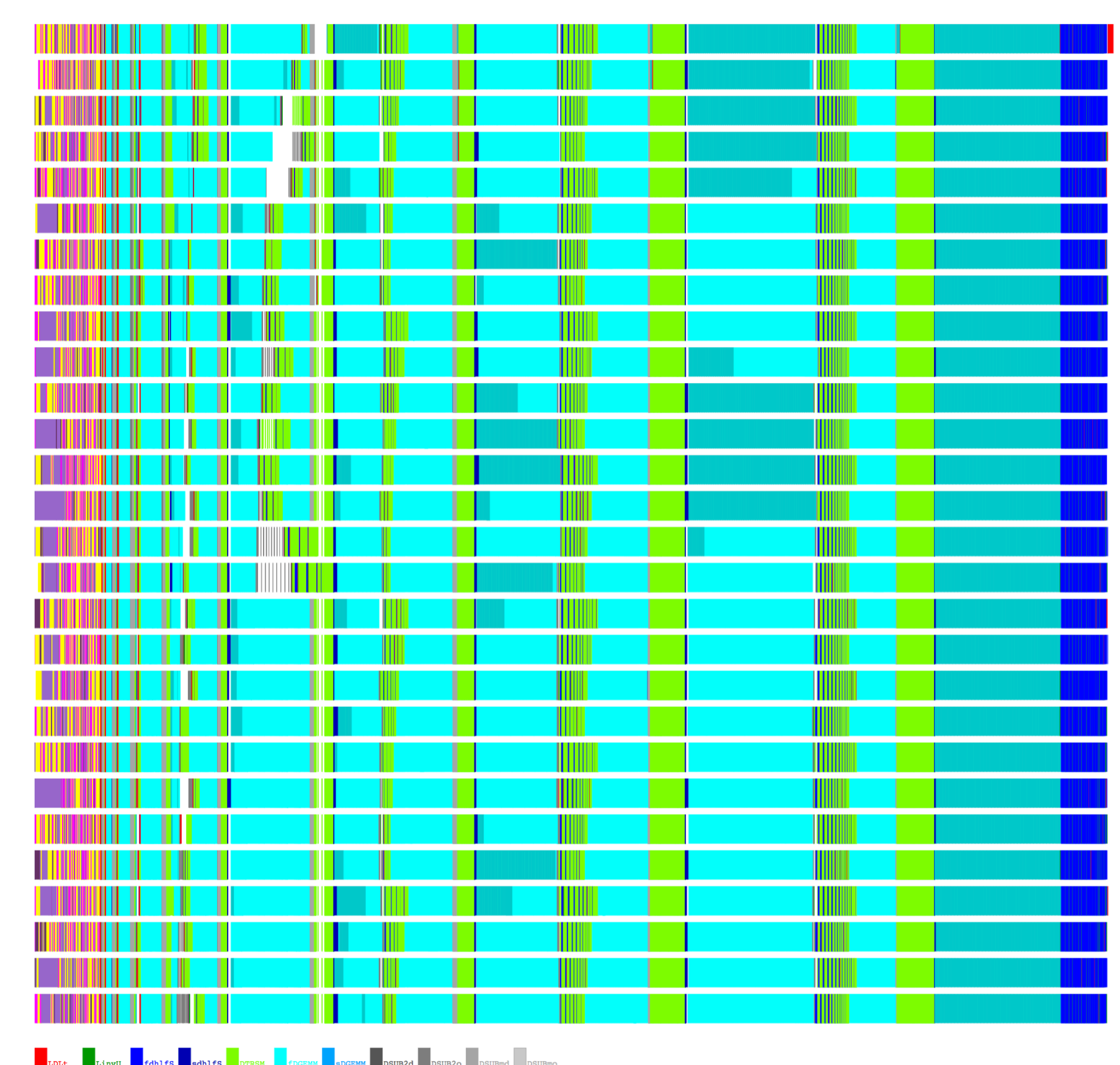
- ▶ α : LDU-factorization $A_{kk} = L_{kk}D_{kk}U_{kk}$
- ▶ β_+ : solve multiple-RHS $L_{kk}X_j = A_{kj}$, $1 \leq j \leq m$
- ▶ β_- : solve multiple-RHS $U_{kk}^T Y_i^T = A_{ik}^T$, $1 \leq i \leq m$
- ▶ γ : update Schur complement $S_{ij} = A_{ij} - Y_i^T(D_{kk}^{-1}X_j)$



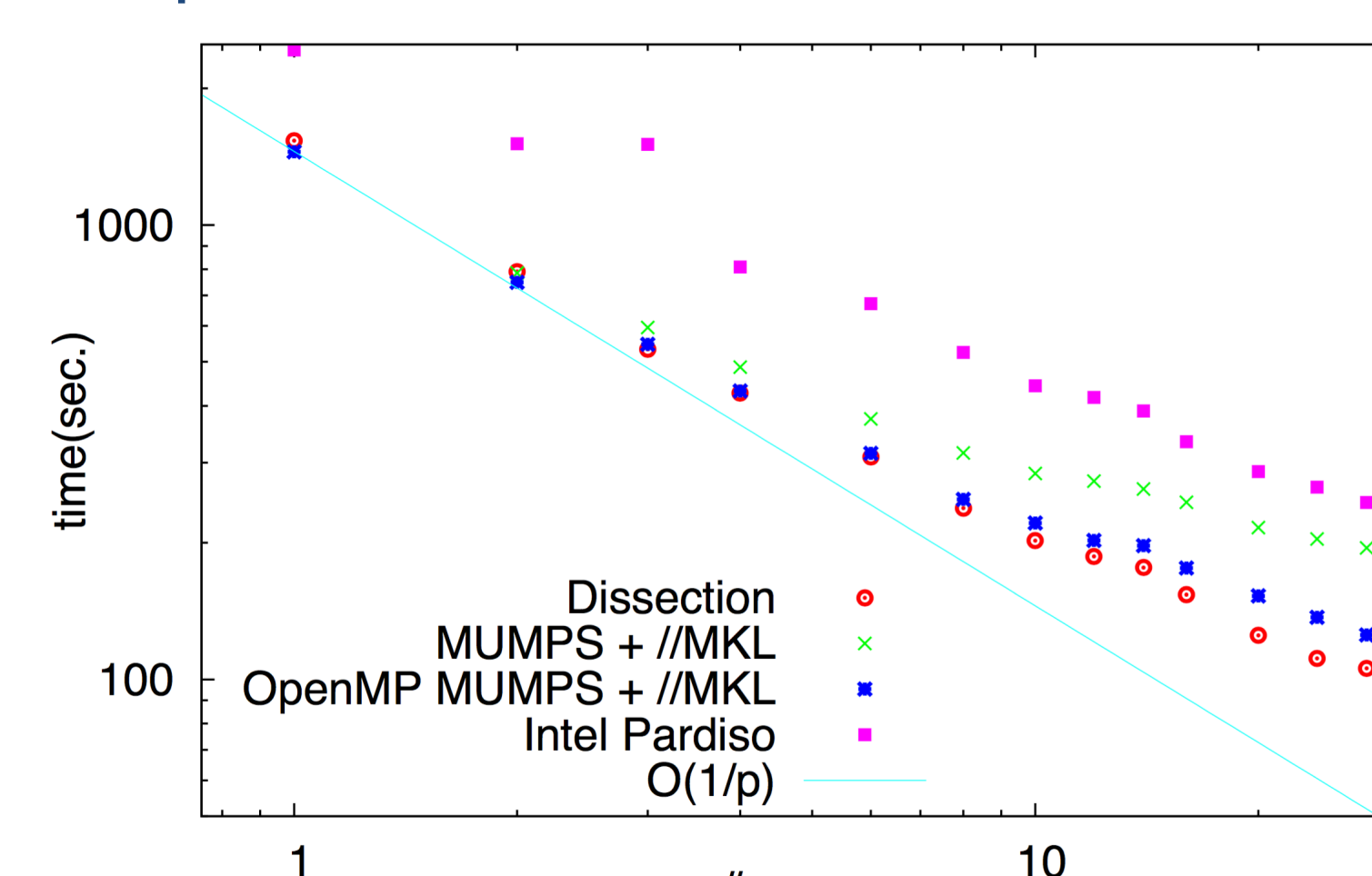
- $\alpha_1^{(1)} \leftarrow \{\beta_{+2}^{(1)}, \beta_{-2}^{(1)}, \gamma_{2,2}^{(1)}, \alpha_2^{(2)}, \beta_{+3}^{(1)}, \beta_{-3}^{(1)}, \beta_{+4}^{(1)}, \beta_{-4}^{(1)}, \dots, \beta_{+n}^{(1)}, \beta_{-n}^{(1)}\}$
- $\leftarrow \{\gamma_{2,3}^{(1)}, \gamma_{3,3}^{(1)}, \dots, \gamma_{3,2}^{(1)}, \dots, \gamma_{n,n}^{(1)}\}$
- $\leftarrow \{\beta_{+3}^{(2)}, \beta_{-3}^{(2)}, \gamma_{3,3}^{(2)}, \alpha_3^{(3)}, \beta_{+4}^{(2)}, \beta_{-4}^{(2)}, \dots, \beta_{+n}^{(2)}, \beta_{-n}^{(2)}\}$
- $\leftarrow \{\gamma_{3,4}^{(2)}, \dots, \gamma_{4,3}^{(2)}, \dots, \gamma_{n,n}^{(2)}\} \leftarrow \dots$
- $\leftarrow \beta_{+n}^{(n-1)}, \beta_{-n}^{(n-1)}, \gamma_{n,n}^{(n-1)}, \alpha_n^{(n)}$.

Dependency of tasks is analyzed and rearranged into set of groups, where all tasks in a group are independent and executed in parallel. Symbols ' \leftarrow ' to show dependence between tasks, ' \leftarrow ' to force sequential execution, and braces to show a group of tasks, are used in the above right figure.

example of task scheduling with 28 cores



parallel performance with FEM matrix of fluid dynamics,



N=1,032,183, nnz=97,961,089 by Xeon v3, 57GB mem