



PASTIX and MaPHYS overview

INRIA, Bordeaux

X. Lacoste, S. Nakov

X. Lacoste
S. Nakov
LaBRI – Inria Bordeaux Sud-Ouest

1

PaStiX

PaStiX

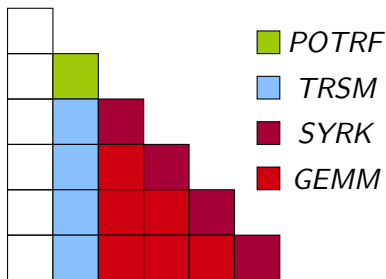
About PaStiX

- ▶ Direct sparse linear solver;
- ▶ Hybrid MPI and P-Thread implementation.

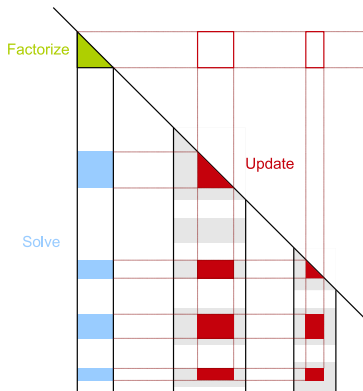
Current works

- ▶ Use GPUs to accelerate factorization;
 - ▶ dedicated sparse GEMM kernel;
 - ▶ implementation on top of generic runtimes (STARPU and PARSEC).
- ▶ Optimized schur complement computation;
- ▶ Finite element assembly API;
- ▶ H-Matrix.

Tasks structure



(a) Dense tile task decomposition



(b) Decomposition of the task applied while processing one panel

Supernodal sequential algorithm

```

forall the Supernode  $S_1$  do
  panel ( $S_1$ );
  /* update of the panel */
  forall the extra diagonal block  $B_i$  of  $S_1$  do
     $S_2 \leftarrow$  supernode_in_front_of ( $B_i$ );
    gemm ( $S_1, S_2$ );
    /* sparse GEMM  $B_{k,k \geq i} \times B_i^T$  subtracted from
        $S_2$  */
  end
end

```

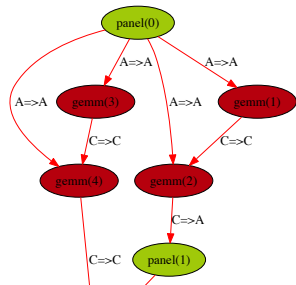
STARPU Tasks submission

```

forall the Supernode  $S_1$  do
  submit_panel ( $S_1$ );
  /* update of the panel                                     */
  forall the extra diagonal block  $B_i$  of  $S_1$  do
     $S_2 \leftarrow$  supernode_in_front_of ( $B_i$ );
    submit_gemm ( $S_1, S_2$ );
    /* sparse GEMM  $B_{k,k \geq i} \times B_i^T$  subtracted from  $S_2$ 
      */
  end
  wait_for_all_tasks ();
end

```

PARSEC's parameterized task graph



Task Graph

```

1  panel(j)
2
3  /* Execution Space */
4  j = 0 .. cblknbr-1
5
6  /* Task Locality (Owner Compute) */
7  :A(j)
8
9  /* Data dependencies */
10 RW A <- ( leaf ) ? A(j) : C gemm( lastbrow )
11     -> A gemm(firstblock+1 .. lastblock)
12     -> A(j)

```

Panel Factorization in JDF Format

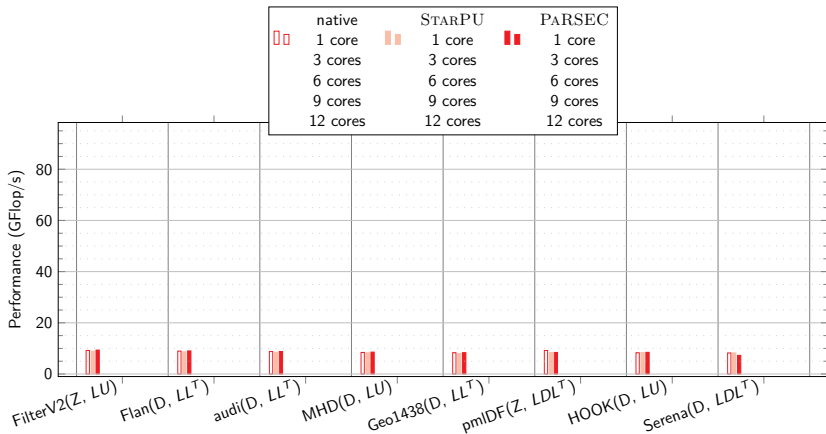
Matrices and Machines

Matrix	Prec	Method	Size	nnz _A	nnz _L	TFlop
FilterV2	Z	<i>LU</i>	0.6e+6	12e+6	536e+6	3.6
Flan	D	<i>LL^T</i>	1.6e+6	59e+6	1712e+6	5.3
Audi	D	<i>LL^T</i>	0.9e+6	39e+6	1325e+6	6.5
MHD	D	<i>LU</i>	0.5e+6	24e+6	1133e+6	6.6
Geo1438	D	<i>LL^T</i>	1.4e+6	32e+6	2768e+6	23
Pmldf	Z	<i>LDL^T</i>	1.0e+6	8e+6	1105e+6	28
Hook	D	<i>LU</i>	1.5e+6	31e+6	4168e+6	35
Serena	D	<i>LDL^T</i>	1.4e+6	32e+6	3365e+6	47

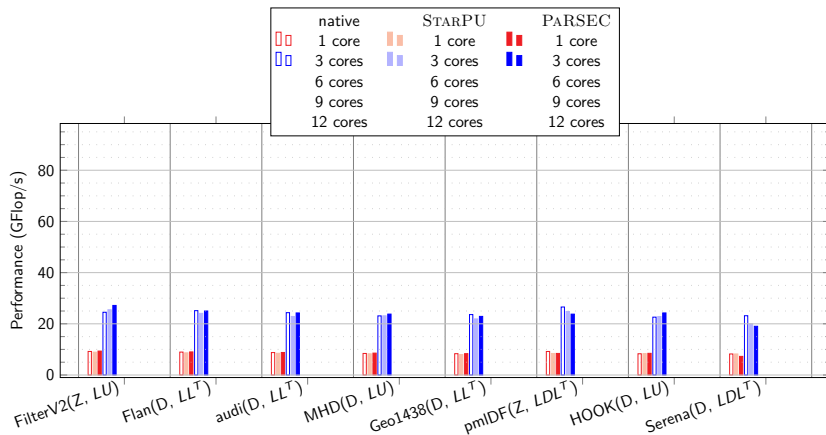
Table: Matrix description (Z: double complex, D: double).

Machine	Processors	Frequency	GPUs	RAM
Mirage	Westmere Intel Xeon X5650 (2 × 6)	2.67 GHz	Tesla M2070 (×3)	36 GB

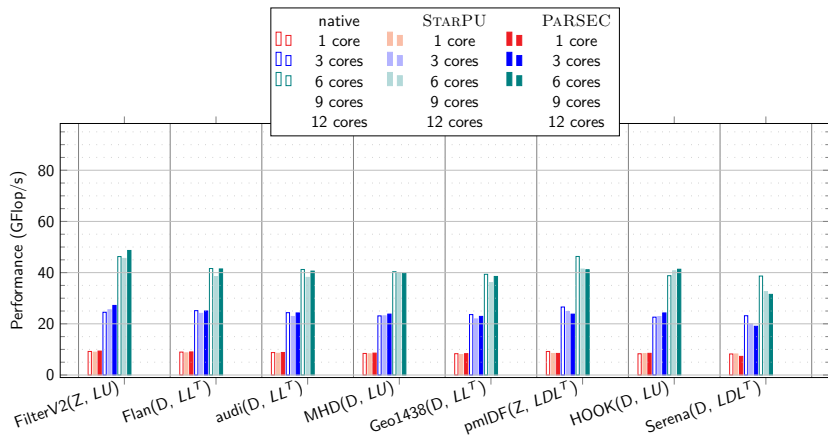
CPU scaling study: GFlop/s during numerical factorization



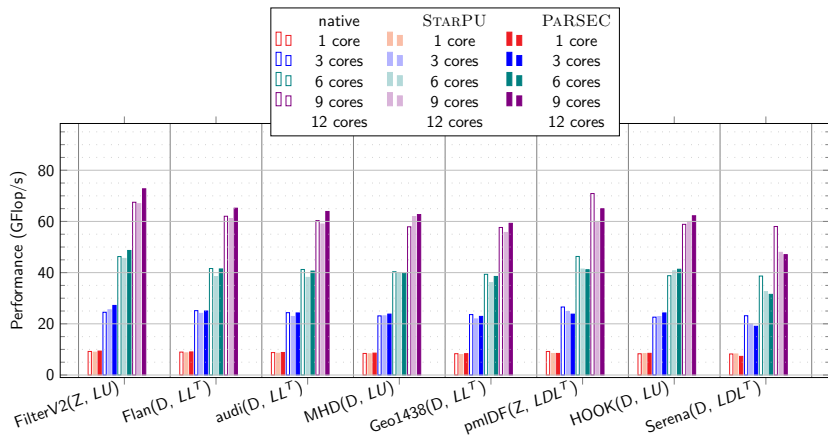
CPU scling study: GFlop/s during numerical factorization



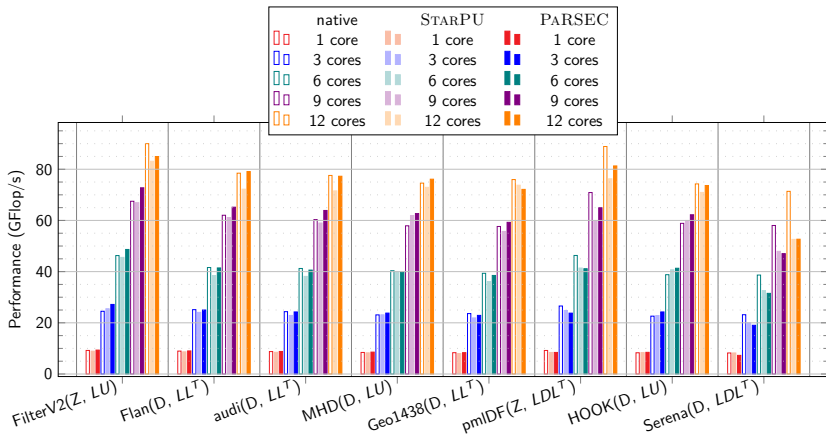
CPU scling study: GFlop/s during numerical factorization



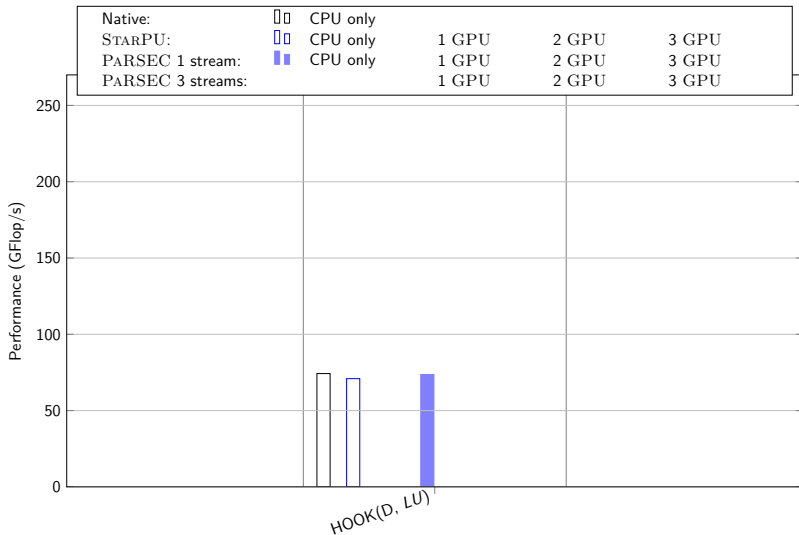
CPU scling study: GFlop/s during numerical factorization



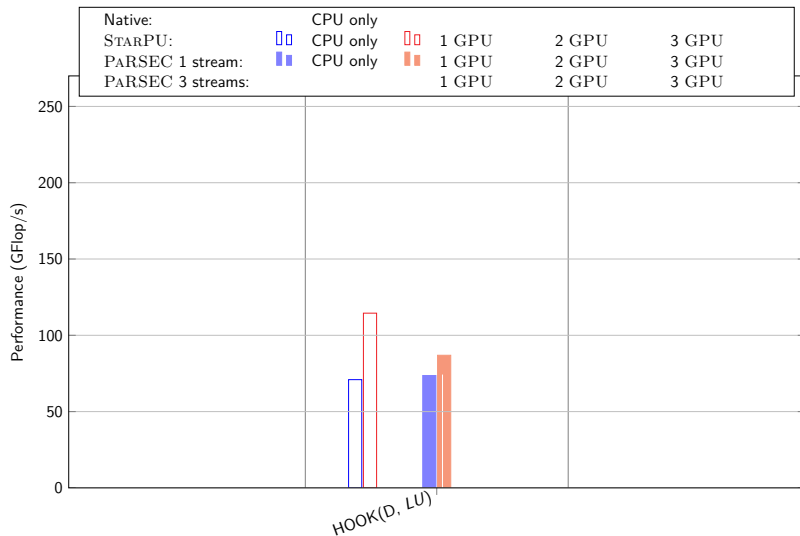
CPU scling study: GFlop/s during numerical factorization



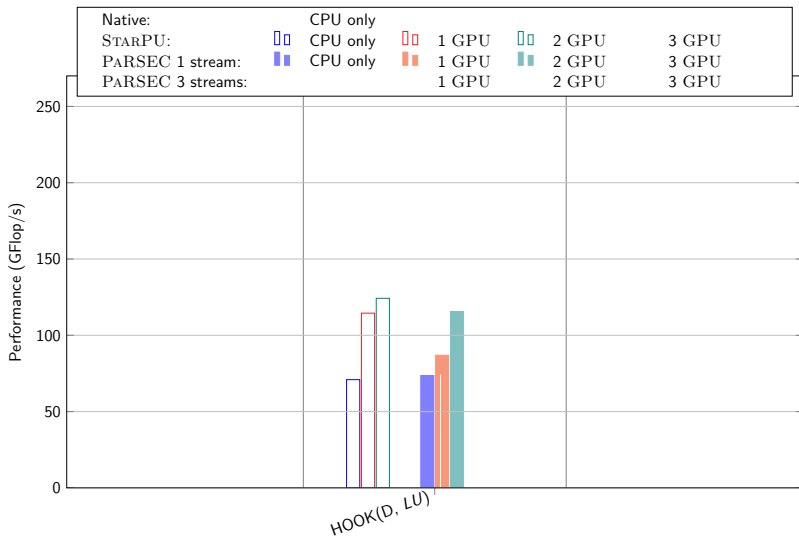
GPU scaling study



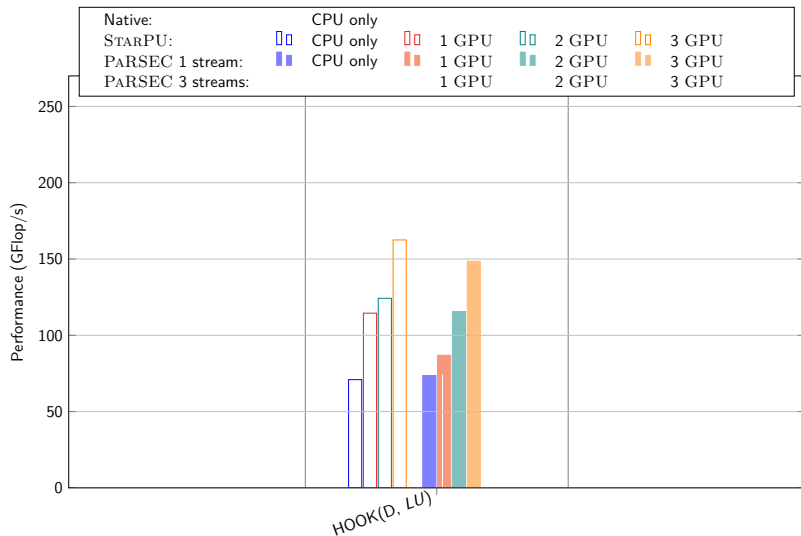
GPU scaling study



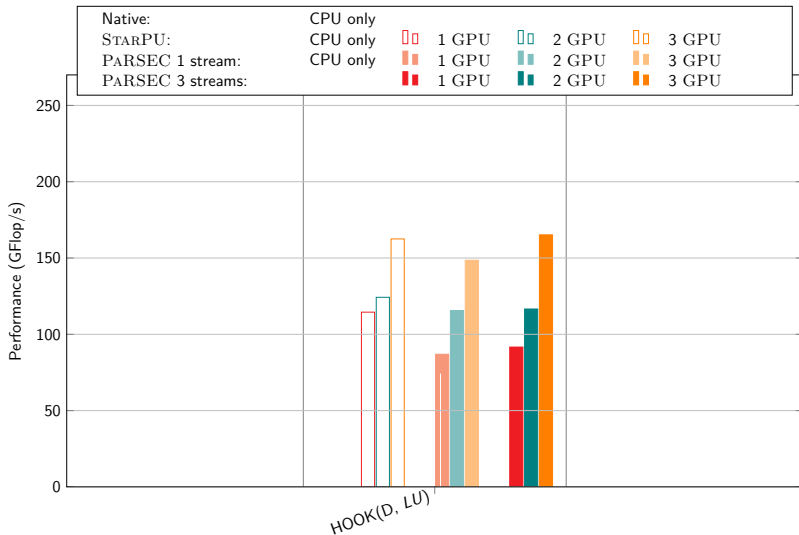
GPU scaling study



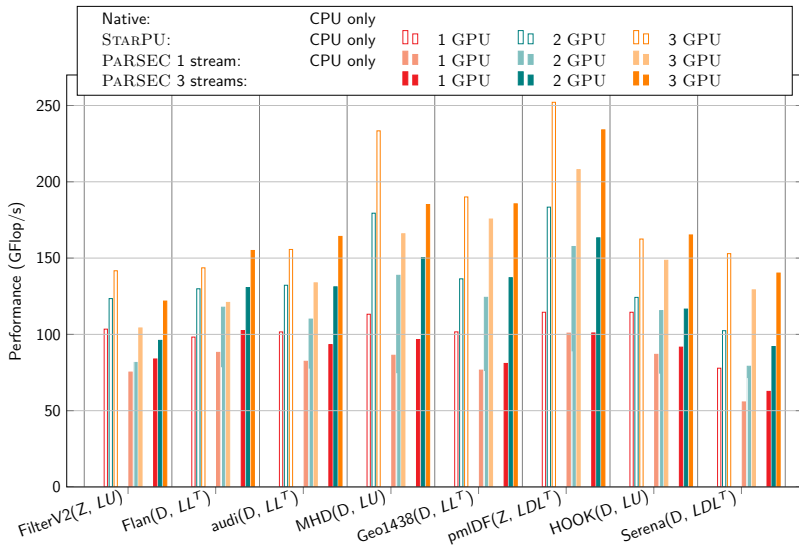
GPU scaling study



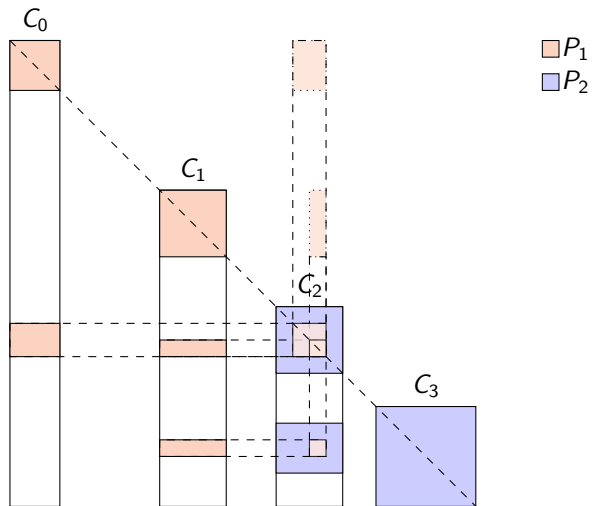
GPU scaling study



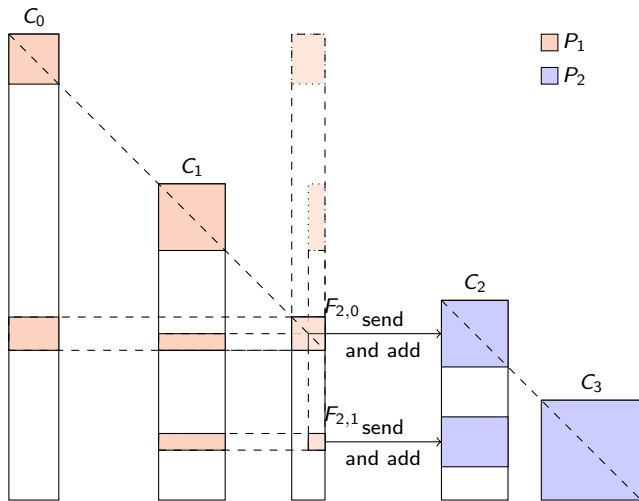
GPU scaling study



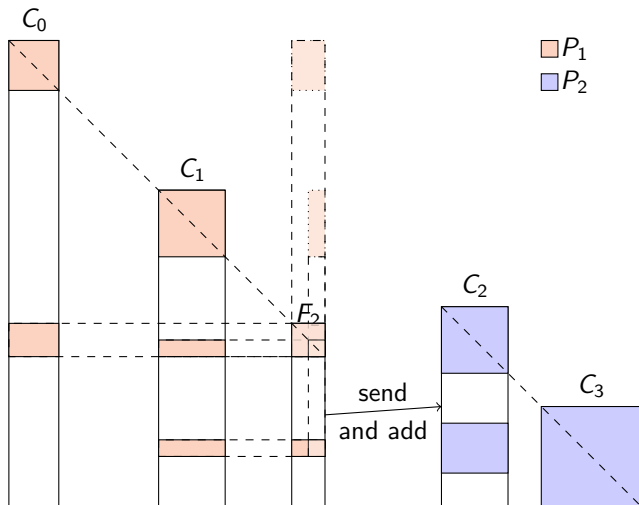
Distributed implementation



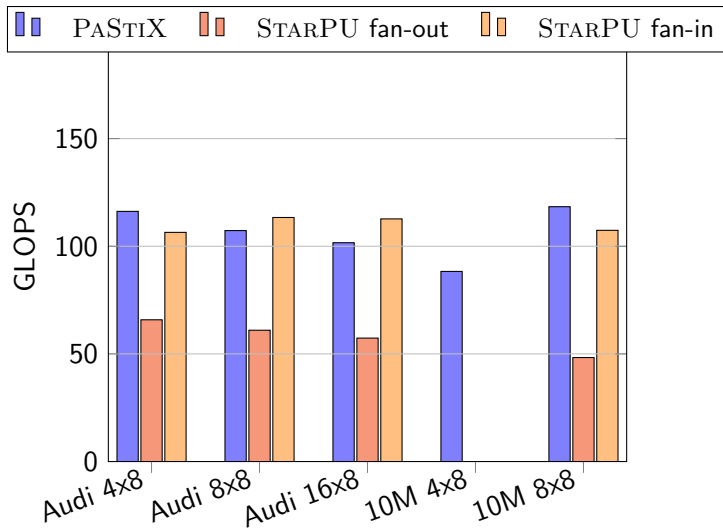
Fanin version (PaStiX)



Fanin version (STARPU)



Distributed preliminary results



Schur complement computation in PASTIX

1. User provides a list of schur unknowns;
2. PASTIX moves this unknown to the last block of the matrix which is centralized;
3. PASTIX perform the factorization and stop before Schur block factorization;
4. PASTIX returns the computed Schur complement;
5. User can call PASTIX to solve using non Schur unknowns.

Improvements on Schur complement

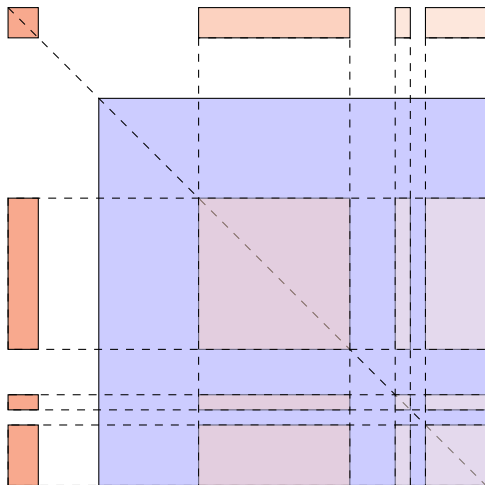
Next release

- ▶ Temporary buffer size optimization:
 - ▶ Do not take into account last block (doesn't produce GEMM);
 - ▶ Split Schur facing blocks as if Schur was splitted;
- ▶ Multi-threaded update of the Schur complement.

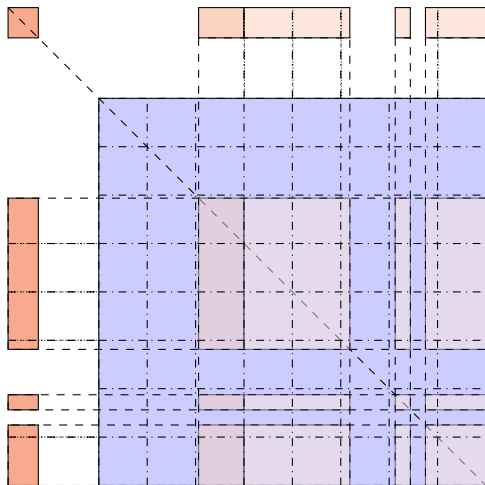
To do

- ▶ Distributed schur complement;
- ▶ Integrate Astrid Casadei work on memory consumption:
 - ▶ Allocate data only when required;
 - ▶ Drop coupling data:
 - ▶ Adapted data structure;
 - ▶ Mix left/right-looking algorithm.

Temporary buffer optimization



Temporary buffer optimization



2

MaPHYS - MPI+threads parallelization

Motivations

Goal: solving $\mathcal{A}x = b$, where A is sparse



Usual trades off

Direct

- ▶ Robust/accurate for general problems
- ▶ BLAS-3 based implementations
- ▶ Memory/CPU prohibitive for large 3D problems

Iterative

- ▶ Problem dependent efficiency / accuracy
- ▶ Sparse computational kernels
- ▶ Less memory requirements and possibly faster

Motivations

Goal: solving $Ax = b$, where A is **sparse**



Usual trades off

Direct

- ▶ Robust/accurate for general problems
- ▶ BLAS-3 based implementations
- ▶ Memory/CPU prohibitive for large 3D problems

Iterative

- ▶ Problem dependent efficiency / accuracy
- ▶ Sparse computational kernels
- ▶ Less memory requirements and possibly faster

Motivations

Goal: solving $\mathcal{A}x = b$, where A is sparse



Usual trades off

Direct

- ▶ Robust/accurate for general problems
- ▶ BLAS-3 based implementations
- ▶ Memory/CPU prohibitive for large 3D problems

Iterative

- ▶ Problem dependent efficiency / accuracy
- ▶ Sparse computational kernels
- ▶ Less memory requirements and possibly faster

Motivations

Goal: solving $\mathcal{A}x = b$, where A is sparse



Usual trades off

Direct

- ▶ Robust/accurate for general problems
- ▶ BLAS-3 based implementations
- ▶ Memory/CPU prohibitive for large 3D problems

Iterative

- ▶ Problem dependent efficiency / accuracy
- ▶ Sparse computational kernels
- ▶ Less memory requirements and possibly faster

Motivations

Goal: solving $\mathcal{A}x = b$, where A is sparse



Usual trades off

Direct

- ▶ Robust/accurate for general problems
- ▶ BLAS-3 based implementations
- ▶ Memory/CPU prohibitive for large 3D problems

Iterative

- ▶ Problem dependent efficiency / accuracy
- ▶ Sparse computational kernels
- ▶ Less memory requirements and possibly faster

Sparse hybrid (direct/iterative) Schur complement solvers based on domain decomposition

The main idea behind these methods is to permute the linear problem $\mathcal{A}x = b$ in a form:

$$\begin{pmatrix} \mathcal{A}_{II} & \mathcal{A}_{I\Gamma} \\ \mathcal{A}_{\Gamma I} & \mathcal{A}_{\Gamma\Gamma} \end{pmatrix} \begin{pmatrix} x_I \\ x_\Gamma \end{pmatrix} = \begin{pmatrix} b_I \\ b_\Gamma \end{pmatrix}$$

where \mathcal{A}_{II} is a block diagonal matrix. Eliminating x_I from the second block row leads to the reduced system

$$\mathcal{S}x_\Gamma = f$$

where

$$\mathcal{S} = \mathcal{A}_{\Gamma\Gamma} - \mathcal{A}_{\Gamma I} \mathcal{A}_{II}^{-1} \mathcal{A}_{I\Gamma} \quad \text{and} \quad f = b_\Gamma - \mathcal{A}_{\Gamma I} \mathcal{A}_{II}^{-1} b_I.$$

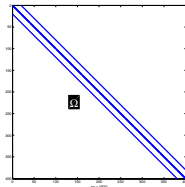
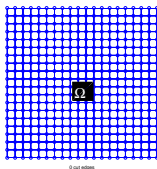
Method used in MaPHYS

- ▶ Partitioning the global matrix in several local matrices

- ▶ Local factorization

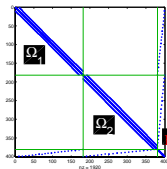
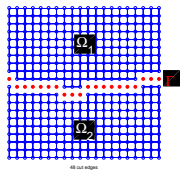
- ▶ Constructing of the preconditioner

- ▶ Solving the reduced system



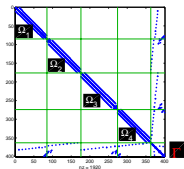
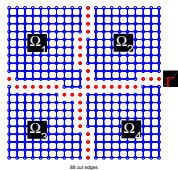
Method used in MaPHYS

- ▶ Partitioning the global matrix in several local matrices
- ▶ Local factorization
- ▶ Constructing of the preconditioner
- ▶ Solving the reduced system



Method used in MaPHYS

- ▶ Partitioning the global matrix in several local matrices
- ▶ Local factorization
- ▶ Constructing of the preconditioner
- ▶ Solving the reduced system

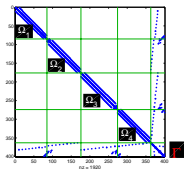
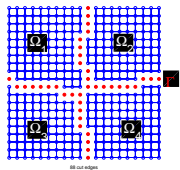


Method used in MAPHYS

- ▶ Partitioning the global matrix in several local matrices
 - ▶ METIS [G. Karypis and V. Kumar]
 - ▶ SCOTCH [Pellegrini and al.]
- ▶ Local factorization

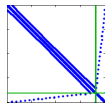
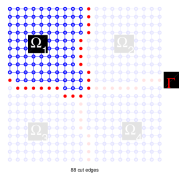
- ▶ Constructing of the preconditioner

- ▶ Solving the reduced system



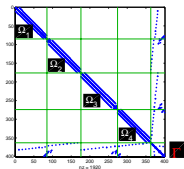
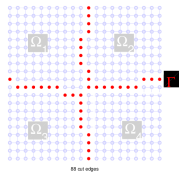
Method used in MAPHYs

- ▶ Partitioning the global matrix in several local matrices
 - ▶ METIS [G. Karypis and V. Kumar]
 - ▶ SCOTCH [Pellegrini and al.]
- ▶ Local factorization
 - ▶ MUMPS [P. Amestoy and al.] (with Schur option)
 - ▶ PASTIX [P. Ramet and al.] (with Schur option)
- ▶ Constructing of the preconditioner
- ▶ Solving the reduced system



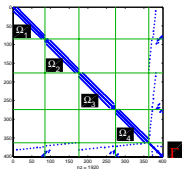
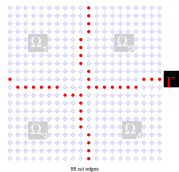
Method used in MaPHYS

- ▶ Partitioning the global matrix in several local matrices
 - ▶ METIS [G. Karypis and V. Kumar]
 - ▶ SCOTCH [Pellegrini and al.]
- ▶ Local factorization
 - ▶ MUMPS [P. Amestoy and al.] (with Schur option)
 - ▶ PASTIX [P. Ramet and al.] (with Schur option)
- ▶ Constructing of the preconditioner
 - ▶ MKL library
- ▶ Solving the reduced system



Method used in MAPHYs

- ▶ Partitioning the global matrix in several local matrices
 - ▶ METIS [G. Karypis and V. Kumar]
 - ▶ SCOTCH [Pellegrini and al.]
- ▶ Local factorization
 - ▶ MUMPS [P. Amestoy and al.] (with Schur option)
 - ▶ PASTIX [P. Ramet and al.] (with Schur option)
- ▶ Constructing of the preconditioner
 - ▶ MKL library
- ▶ Solving the reduced system
 - ▶ CG/GMRES/FGMRES [V.Fraysse and L.Giraud] using MKL library for the reduced system



Software used in MAPHYs (before)

Partitioners

- ▶ SCOTCH
- ▶ METIS

Dense direct solver

- ▶ MKL library

Sparse direct solvers

- ▶ MUMPS
- ▶ PASTIX

Iterative Solvers

- ▶ CG/GMRES/FGMRES using MKL library

Software used in MAPHYS (this study)

Dense direct solver

- ▶ Multi-threaded MKL library

Sparse direct solvers

- ▶ MUMPS
- ▶ Multi-threaded PASTIX

Iterative Solvers

- ▶ CG/GMRES/FGMRES using multi-threaded MKL library

Software used in MAPHYS (this study)

Dense direct solver

- ▶ Multi-threaded MKL library

Sparse direct solvers

- ▶ MUMPS
- ▶ Multi-threaded PASTIX

Iterative Solvers

- ▶ CG/GMRES/FGMRES using multi-threaded MKL library
- ▶ Challenge
 - ▶ Composability
 - ▶ Performance

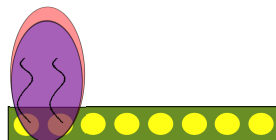
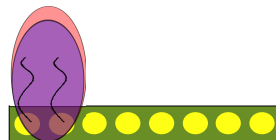
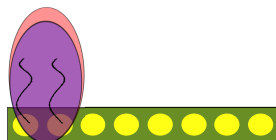
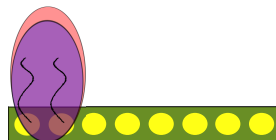
Scalability of MaPHYS on multicore nodes

● MPI process { Thread ● Domain



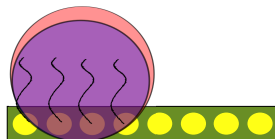
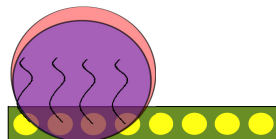
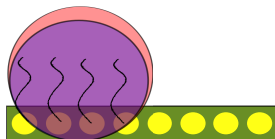
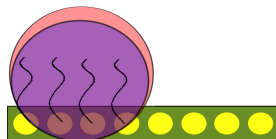
Scalability of MaPHYS on multicore nodes

● MPI process { Thread ● Domain



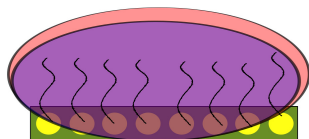
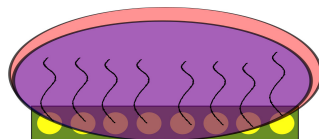
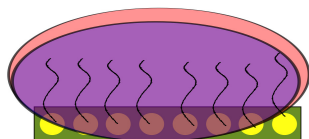
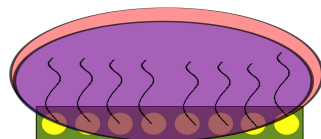
Scalability of MaPHYS on multicore nodes

● MPI process { Thread ● Domain



Scalability of MaPHYS on multicore nodes

● MPI process { Thread ● Domain



Experimental set up

Hardware (on each node)

- ▶ Two Quad-core Nehalem Intel[®] Xeon[®] X5550
- ▶ Memory: 24 GB GDDR3
- ▶ Double precision

Matrices

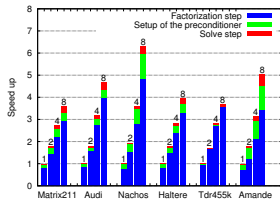
Matrix	Matrix211k	Audi	Nachos	Haltere	Amande
Id	1	2	3	4	5
N	801K	943K	1,120K	1,288K	6,994K
Nnz	129,4M	39,29M	39,9M	10,47M	58,47M

Table: Overview of sparse matrices used in this study.

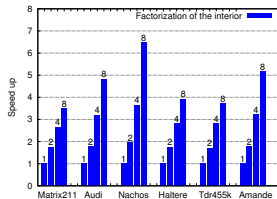
Scalability of MAPHYS on multicore nodes

Achieved performance on four nodes with dense preconditioner

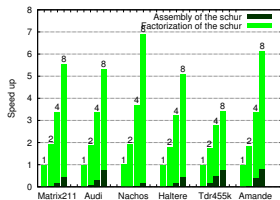
All computational steps



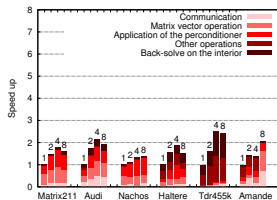
Factorization step



Preconditioning step

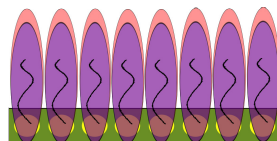
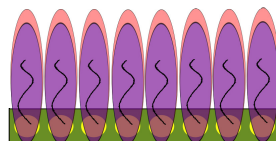
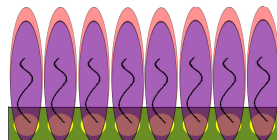
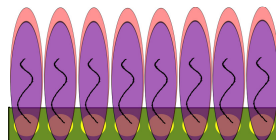


Solve step



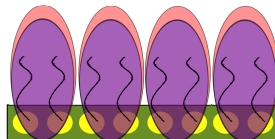
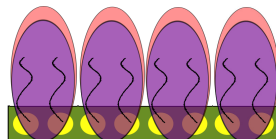
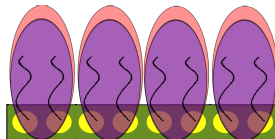
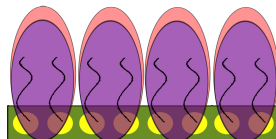
Flexibility to exploit entire multicore nodes

● MPI process { Thread ● Domain



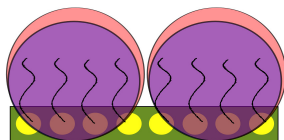
Flexibility to exploit entire multicore nodes

● MPI process { Thread ● Domain

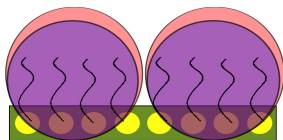


Flexibility to exploit entire multicore nodes

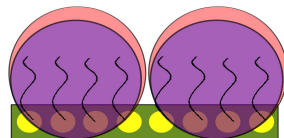
● MPI process { Thread ● Domain



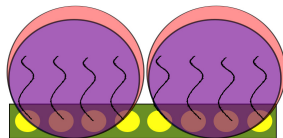
Node 1



Node 2



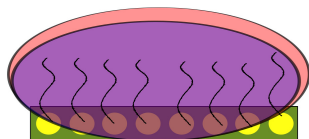
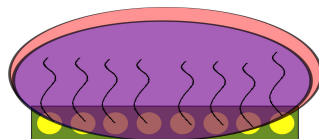
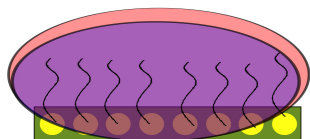
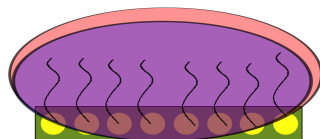
Node 3



Node 4

Flexibility to exploit entire multicore nodes

● MPI process } Thread ● Domain



Experimental set up

Hardware (on each node)

- ▶ Two Quad-core Nehalem Intel[®] Xeon[®] X5550
- ▶ Memory: 24 GB GDDR3
- ▶ Double precision

Matrices

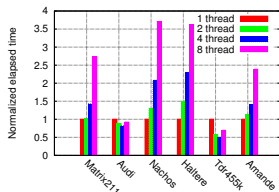
Matrix	Matrix211	Audi	Nachos	Haltere	Tdr455k	Amande
Nb.	1	2	3	4	5	6
N	801K	943K	1,120K	1,288K	2,738K	6,994K
Nnz	129,4M	39,29M	39,9M	10,47M	112,7M	58,47M
Nb_nodes	8	4	16	4	8	16
Preconditioner	dense	sparse03	sparse02	sparse03	dense	sparse03

Table: Overview of sparse matrices used in this study.

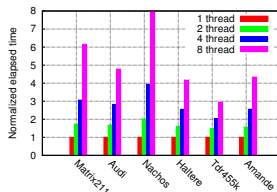
Flexibility to exploit entire multicore nodes

Achieved performance when all cores are used

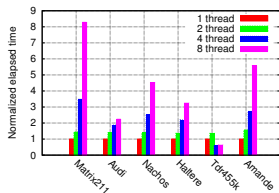
All computational steps



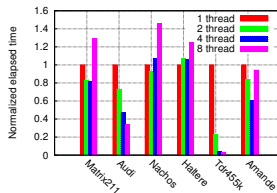
Factorization step



Preconditioning step



Solve step



Experimental set up

Hopper platform (Hardware)

- ▶ Two twelve-core AMD 'MagnyCours' 2.1-GHz
- ▶ Memory: 32 GB GDDR3
- ▶ Double precision

Matrices

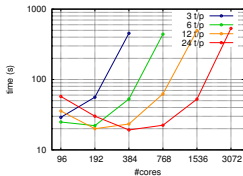
Matrix	Tdr455K	Nachos4M
Nb.	1	2
N	2,738K	4,147K
Nnz	112,7M	256,4M
Preconditioner	dense	sparse02

Table: Overview of sparse matrices used on the Hopper platform.

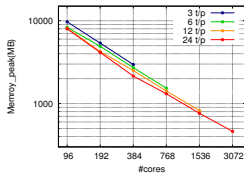
Results on the Hopper platform.

Achieved performance for the Tdr455K matrix

All computational steps

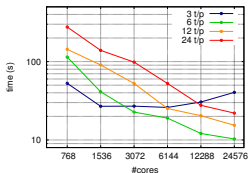


Memory used per node

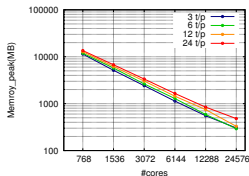


Achieved performance for the Nachos4M matrix.

All computational steps



Memory used per node



PAStIX

- ▶ Recent developpements:
 - ▶ Heterogeneous direct sparse linear solver using generic runtime;
 - ▶ Schur complement computation
 - ▶ Generic finite element oriented matrix assembly interface.
- ▶ Current developpements:
 - ▶ Distributed heterogeneous direct sparse linear solver;
 - ▶ Distributed Schur and optimizations;
 - ▶ Low-rank compression - H-Matrix.
 - ▶ Rewritten code for better upgradability, CMake compilation tools.

MAPHyS

- ▶ Recent study:
 - ▶ With the two-level parallelism we are able to efficiently exploits multicore architectures
 - ▶ Presentation of Julien Pedron (tomorrow)
- ▶ Current work:
 - ▶ Comparison with the PDSLIn solver
 - ▶ Further experiments in the collaboration with TOTAL
 - ▶ Full task-based hybrid prototype

Thanks !



Xavier LACOSTE
Stojce NAKOV
INRIA HiePACS team
July 10, 2014