



Improving mapping for sparse direct solvers: A trade-off between data locality and load balancing

Changjiang Gou, Ali Al Zoobi, Anne Benoit, Mathieu Faverge, Loris Marchal, Grégoire Pichon, Pierre Ramet

► To cite this version:

Changjiang Gou, Ali Al Zoobi, Anne Benoit, Mathieu Faverge, Loris Marchal, et al.. Improving mapping for sparse direct solvers: A trade-off between data locality and load balancing. EuroPar 2020 - 26th International European Conference on Parallel and Distributed Computing, Aug 2020, Warsaw / Virtual, Poland. pp.1-16. hal-02973315

HAL Id: hal-02973315

<https://inria.hal.science/hal-02973315>

Submitted on 21 Oct 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Improving mapping for sparse direct solvers: A trade-off between data locality and load balancing

Changjiang Gou^{1,2}, Ali Al Zoobi³, Anne Benoit¹, Mathieu Faverge⁴,
Loris Marchal¹, Grégoire Pichon¹, and Pierre Ramet⁴

¹ Univ Lyon, EnsL, UCBL, CNRS, Inria, LIP

² East China Normal University, Shanghai, China

³ Inria Sophia-Antipolis

⁴ Univ. Bordeaux and Inria

Abstract In order to express parallelism, parallel sparse direct solvers take advantage of the elimination tree to exhibit tree-shaped task graphs, where nodes represent computational tasks and edges represent data dependencies. One of the pre-processing stages of sparse direct solvers consists of mapping computational resources (processors) to these tasks. The objective is to minimize the factorization time by exhibiting good data locality and load balancing. The proportional mapping technique is a widely used approach to solve this resource-allocation problem. It achieves good data locality by assigning the same processors to large parts of the elimination tree. However, it may limit load balancing in some cases. In this paper, we propose a dynamic mapping algorithm based on proportional mapping. This new approach, named STEAL, relaxes the data locality criterion to improve load balancing. In order to validate the newly introduced method, we perform extensive experiments on the PASTIX sparse direct solver. It demonstrates that our algorithm enables better static scheduling of the numerical factorization while keeping good data locality.

1 Introduction

For the solution of large sparse linear systems, we design numerical schemes and software packages for direct parallel solvers. Sparse direct solvers are mandatory when the linear system is very ill-conditioned for example [5]. Therefore, to obtain an industrial software tool that must be robust and versatile, high-performance sparse direct solvers are mandatory, and parallelism is then necessary for reasons of memory capability and acceptable solution time. Moreover, in order to solve efficiently 3D problems with several million unknowns, which is now a reachable challenge with modern supercomputers, we must achieve good scalability in time and control memory overhead. Solving a sparse linear system by a direct method is generally a highly irregular problem that provides some challenging algorithmic problems and requires a sophisticated implementation scheme in order to fully exploit the capabilities of modern supercomputers.

There are two main approaches in direct solvers: the multifrontal approach [2,7], and the supernodal one [10,17]. Both can be described by a computational tree

whose nodes represent computations and whose edges represent transfer of data. In the case of the multifrontal method, at each node, some steps of Gaussian elimination are performed on a dense frontal matrix and the remaining Schur complement, or contribution block, is passed to the parent node for assembly. In the case of the supernodal method, the distributed memory version uses a right-looking formulation which, having computed the factorization of a supernode corresponding to a node of the tree, then immediately sends the data to update the supernodes corresponding to ancestors in the tree. In a parallel context, we can locally aggregate contributions to the same block before sending the contributions. This can significantly reduce the number of messages. Independently of these different methods, a static or dynamic scheduling of block computations can be used. For homogeneous parallel architectures, it is useful to find an efficient static scheduling.

In order to achieve efficient parallel sparse factorization, we perform the three sequential preprocessing phases:

1. The ordering step, which computes a symmetric permutation of the initial matrix such that the factorization process will exhibit as much concurrency as possible while incurring low fill-in.
2. The block symbolic factorization step, which determines the block data structure of the factorized matrix associated with the partition resulting from the ordering phase. From this block structure, one can deduce the weighted elimination quotient graph that describes all dependencies between column-blocks, as well as the supernodal elimination tree.
3. The block scheduling/mapping step, which consists in mapping the resulting blocks onto the processors. During this mapping phase, a static optimized scheduling of the computational and communication tasks, according to models calibrated for the target machine, can be computed.

When these preprocessing phases are done, the computation on the actual data, that is the numerical factorization, can start.

The optimization problem that needs to be solved at the scheduling/mapping stage is known to be NP-hard, and is usually solved using a proportional mapping heuristic [15]. This mono-constraint heuristic induces idle times during the numerical factorization. In this paper, we extend the proportional mapping and scheduling heuristic to reduce these idle times. We first detail in Section 2 proportional mapping heuristic with its issues and related work, before describing the original application in the context of the PASTIX solver [11] in Section 3. Then, in Section 4, we explain the introduced solution before studying its impact on a large set of test cases in Section 5. Conclusion and future work directions are presented in Section 6.

2 Problem statement and related work

Among different mapping strategies that are used by both supernodal and multifrontal sparse direct solvers, the subtree to subcube mapping [8] and the proportional mapping [15] are the most popular. These approaches consist of tree

partitioning techniques, where the set of resources mapped on a node of the tree are split among disjoint subsets, each mapped to a child subtree.

The proportional mapping method performs a top-down traversal of the elimination tree, during which each node is assigned a set of computational resources. All the resources are assigned to the root node, which performs the last task. Then, the resources are split recursively following a balancing criterion. The set of resources dedicated to a node are split among its children, proportionally to their weight or any other balancing criterion. This recursive process ends at the leaves of the tree, or when entire subtrees are mapped onto a single resource.

The original version of the proportional mapping [15] computes the splitting of resources depending on the workload of each subtree, but more sophisticated metrics can also be used. In [16], a scheduling strategy was proposed for tree-shaped task graphs. The time for computing a parallel task (for instance at the root node of the elimination tree) is considered as proportional to the length of the task and to a given parallel efficiency. This method was proven efficient in [3] for a multifrontal solver. The proportional mapping technique is widely used because it helps reducing the volume of data transfers due to its data locality. In addition, it allows us to exhibit both tree and node parallelism.

Note that alternative solutions to the proportional mapping have been proposed, such as the 2D block-cyclic distribution of SUPERLU [13], or the 1D cyclic distribution of SYMPACK [12]. In the latter, the non load-balanced solution is compensated by a complex and advanced communication scheme that balances the computations in the nodes to get good performance results out of this mapping strategy.

As stated earlier, sparse direct solvers commonly use the proportional mapping heuristic to distribute supernodes (a full set of columns, i.e., 1D distribution that share the same row pattern) onto the processors. Note that each supernode can be split into smaller nodes to increase the level of parallelism, which modifies the original supernodes tree structure as shown in Figure 2. This heuristic provides a set of candidate processors for each supernode, which is then refined dynamically when going up the tree, as in MUMPS [1] or PASTIX [11], with a simulation stage that affects a single processor among the candidates, while providing a static optimized scheduling. The proportional mapping stage, by its construction, may however introduce idle time in the scheduling. This is illustrated on Figure 1. The ten candidate processors of the root supernode are distributed among the two sons of weight respectively 4 and 6. The Gantt diagram points out the issue of considering a single criterion heuristic to set the mapping: no work is given to processor p_9 due to the low level of parallelism of the right supernode, whereas it could benefit to the left supernode.

A naive way to handle this issue is to avoid the proportional mapping stage, and consider only the scheduling stage with all processors as candidates for each node of the tree. The drawback of this method is that 1) it does not preserve the data locality, and 2) it drastically increases the complexity of the scheduling step. This solution has been implemented in the PASTIX solver for comparison, and it is referred to as ALL2ALL, since all processors are candidates to all nodes.

3 Description of the application

At a coarse-grain level, the computation can be viewed as a tree T whose nodes (or vertices) represent supernodes of the matrix, and where the dependencies are directed towards the root of the tree. Because sparse matrices usually represent physical constraints and thanks to the nested dissection used to order the matrix, supernodes at the bottom of the tree are usually small and supernodes at the top are much larger. Each supernode is itself a small DAG (Directed Acyclic Graph) of tasks as illustrated on Fig. 2. A more refined view shows that the dependencies between two supernodes consist of dependencies between tasks of these supernodes. Another way to put it is that the computation is described as a DAG of tasks, tasks are partitioned into supernodes, and the quotient graph of supernodes is the tree T (with some transitivity edges). Note that with 1D distribution, as targetted here, the DAG within can also be seen as a tree with dependencies toward the roots. Thus, in this paper, we will use either nodes or supernodes to denote the vertices of the tree T as they can be used interchangeably.

This structure in two levels allows us to both reduce the cost of the analysis stage by considering only the first level (supernodes), while increasing the parallelism level (nodes) during the numerical factorization with finer grain computations.

We denote by $root(T)$ the node at the root of tree T , and by w_i the computational weight of the node i , for $1 \leq i \leq n$: this is the total number of operations of all tasks within node i . Also, $parent(i)$ is the parent of node i in the tree (except for the root), and $child(i)$ are the children nodes of i in the tree. Given a subtree T_i of T (rooted in $root(T_i)$), $W_i = \sum_{j \in T_i} w_j$ is the computational weight of this subtree.

As stated above, each node i of the tree is itself made of $n_i \geq 1$ tasks i_1, \dots, i_{n_i} , whose dependencies follow a directed acyclic graph (DAG). Each of these tasks is a linear algebra kernel (such as matrix factorization, triangular solve or matrix product) on block matrices. Hence, given a node i and its par-

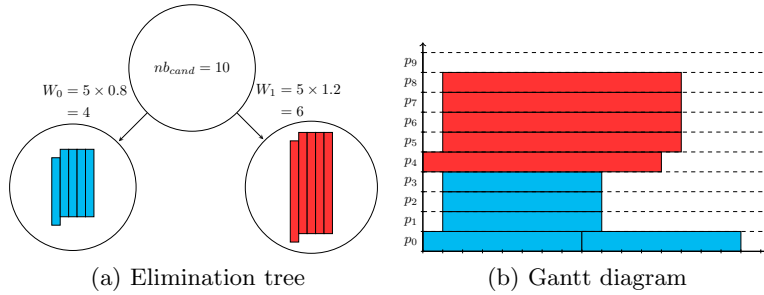


Figure 1: Illustration of proportional mapping: elimination tree on the left, and Gantt diagram on the right.

ent $j = \text{parent}(i)$ in the tree, only some of the tasks of i need to be completed before j is started, which allows some pipelining in the processing of the tree.

When running on a parallel platform with a set P of p processors, nodes and tasks are distributed among available processing resources (processors) in order to ensure a good load-balancing. If node i is executed on $\text{alloc}[i] = k$ processors, its execution time is $f_i(k)$; this time depends on w_i and on the structure of the DAG of tasks.

Following the structure of the application, the mapping is done in two phases: the first phase, detailed in Section 3.1, consists in using the Proportional Mapping algorithm [15] to compute a mapping of nodes to subsets of processors. The second phase, detailed in Section 3.2, refines this mapping by allocating each task of a node i to a single processor of the subset allocated to i in the first step.

3.1 Coarse-grain load balancing using proportional mapping

The proportional mapping process follows the sketch of Algorithm 1. First, all processors are allocated to the root of the tree. Then, we compute the total weight of its subtrees (i.e., the sum of the weight of their nodes), and allocate processors to subtrees so that the load is balanced. Then, we recursively apply the same procedure on each subtree.

Apart from balancing the load among branches of the tree, the proportional mapping is known for its good data locality: a processor is allocated to nodes of a single path from a leaf to the root node, and only to nodes on this path. Thus, the data produced by a node and used by its parents mostly stay on a single processor, and no data transfer is made except for the necessary redistribution of data in the upper levels of the tree. This is particularly interesting in a distributed context, where communications among processors are costly.

We can wonder if Algorithm 1 really optimizes load-balancing, as subtrees with similar total weight W_i may exhibit different levels of parallelism, and thus end up with a different completion time, as illustrated with the example of Fig. 1. The formula $W_i/|P_i|$ correctly computes the duration of the subtree processing

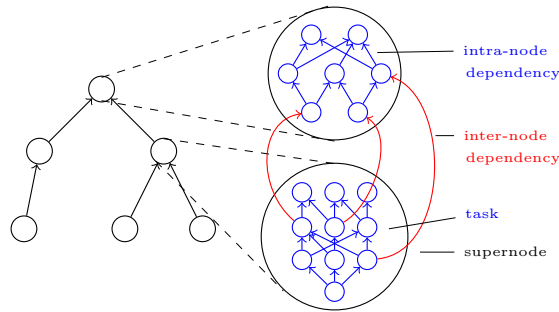


Figure 2: Structure of the computation: tree of supernodes, each supernodes being made of several tasks.

only for perfect parallelism. We propose here another mapping algorithm that optimizes the total computation time, under the constraint of perfect data locality. It iteratively adds processors to the root, and recursively to the subtree with the largest completion time (see Algorithm 2). In this mapping algorithm, $alloc[i]$ represents the number of processors allocated to node i , and $endTime[i]$ represents the completion time of task i . We assume that the function $f_i(k)$, that gives the duration of node i on k processors, is non-increasing with k and is known to the algorithm.

Theorem 1. *The GreedyMappingInt algorithm (Algorithm 2) computes an allocation with minimum total completion time under the constraint that each processor is only allocated to nodes on a path from a leaf to the root.*

Note that this result, proven in the companion research report [9], does not require a particular speed function for tasks: it is valid when the processing time of a task does not increase with the number of processors allocated to the task.

However, both previous mapping algorithms suffer a major problem when used in a practical context, because they forbid allocating processors to more than one child of a node. First, some nodes, especially leaves, have very small weight and several of them should be mapped on the same processor. Second, allocating integer numbers of processors to nodes creates unbalanced workloads, for example, when three processors have to be allocated to two identical subtrees. All implementations of the proportional mapping tackle this problem (including the first one in [15]). For example, the actual implementation in PASTIX, as sketched in Algorithm 3, allows “border processors” to be shared among branches, and keeps track of the occupation of each processor to ensure load-balancing. It first computes the total time needed to process the whole tree, and sets the initial availability time of each processor to an equal share of this total time. Whenever some (fraction of a) node is allocated to a processor, its availability time is reduced. Hence, if a processor is shared on two subtrees T_1, T_2 , the work allocated by T_1 is taken into account when allocating resources for T_2 . Note also that during the recursive allocation process, the subtrees are sorted by non-increasing total weights before being mapped to processors. This allows us to group small subtrees together in order to map them on a single processor, and to avoid unnecessary splitting of processors.

Algorithm 1 Proportional mapping with integer number of processors

function *PropMapInt*(tree T , set P of processors):
 Allocate all processors in P to the root of tree T
 For each subtree T_i of T , compute its total weight W_i
 Find subsets of processors P_i such that $\max(W_i/|P_i|)$ is minimal and $\sum |P_i| = |P|$
 For each subtree T_i of T , call *PropMapInt*(T_i, P_i)

Algorithm 2 Greedy mapping with integer number of processors

```

function GreedyMappingInt(tree  $T$ , number of processors  $p$ ):
   $alloc[1, 2, \dots, n] = [0, \dots, 0]$ 
   $endTime[1, 2, \dots, n] = [\infty, \dots, \infty]$ 
  for  $k = 1, \dots, p$  do
    Call AddOneProcessor( $root(T)$ )
  end for

  function AddOneProcessor(task  $i$ ):
     $alloc[i] \leftarrow alloc[i] + 1$ 
    if  $i$  is a leaf then
       $endTime[i] \leftarrow f_i(alloc[i])$  (duration of node  $i$  on  $alloc[i]$  processors)
    else
      Let  $j$  be the child of  $i$  with largest  $endTime[j]$ 
      AddOneProcessor( $j$ )
       $endTime[i] \leftarrow \max_{j \in child(i)} (endTime[j]) + f_i(alloc[i])$ 
    end if

```

3.2 Refined mapping

After allocating nodes of the tree to subsets of processors, a precise mapping of each task to a processor has to be computed. In PASTIX, this is done by simulating the actual factorization, based on the prediction of both the running times of tasks and of the time needed for data transfers. The refined mapping process is detailed in Algorithm 4. Thanks to the previous phase, we know that each task can run on a subset of processors (the subset associated to the node it belongs to), called *candidate processors* for this task. We associate to each processor a ready queue, containing tasks whose predecessors have already completed, and a waiting queue, with tasks that still have some unfinished predecessor. At the beginning of the simulation, each task is put in the waiting queue of all its candidate processors (except tasks without predecessors, which are put in the ready task queue of their candidate processors). Queues are sorted by decreasing depth of the tasks in the graph (tasks without predecessors are ordered first). The depth considered here is an estimation of the critical path length from the task to the root of the tree T .

A ready time is associated both to tasks and processors:

- The ready time $RP[k]$ of processor k is the completion time of the current task being processed by k (initialized with 0).
- The ready time $RT[i]$ of task i is the earliest time when i can be started, given its input dependencies. This is at least equal to the completion time of each of its predecessors, but also takes into account the time needed for data movement, in case a predecessor of i is not mapped on the same processor as i . The ready time of tasks with non-started predecessor is set to $+\infty$.

Algorithm 3 Proportional mapping with shared processors among subtrees

```

function ProportionalMappingShared(tree  $T$ , number of processors  $p$ ):
  for each processor  $k = 1, \dots, p$  do
     $avail\_time[k] = \sum_{i \in T} w_i / p$ 
  end for
  Call PropMapSharedRec( $T, 1, p$ )

function PropMapSharedRec(subtree  $T$ , indices  $first\_proc, last\_proc$ ):
  if  $last\_proc = first\_proc$  then
    Map all nodes in subtree  $T$  to processor  $first\_proc$ 
     $avail\_time[first\_proc] = avail\_time[first\_proc] - \sum_{i \in T} w_i$ 
  else
    Map node  $r = root(T)$  to all processors in  $first\_proc, \dots, last\_proc$ 
    for each  $k = first\_proc, \dots, last\_proc$  do
       $avail\_time[k] = avail\_time[k] - w_r / (last\_proc - first\_proc)$ 
    end for
     $next\_proc \leftarrow first\_proc$ 
    Sort the subtrees of  $T$  by non-increasing total weight
    for each subtree  $T_i$  in this order do
       $cumul\_time \leftarrow 0$ 
       $w_{subtree} \leftarrow \sum_{j \in T_i} w_j$ 
       $first\_proc\_for\_subtree \leftarrow next\_proc$ 
      while  $cumul\_time < w_{subtree}$  do
         $new\_time\_share \leftarrow \min(w_{subtree} - cumul\_time, avail\_time[next\_proc])$ 
         $cumul\_time \leftarrow cumul\_time + new\_time\_share$ 
         $avail\_time[next\_proc] \leftarrow avail\_time[next\_proc] - new\_time\_share$ 
        if  $avail\_time[next\_proc] = 0$  then  $next\_proc \leftarrow next\_proc + 1$ 
      end while
      PropMapSharedRec( $T_i, first\_proc\_for\_subtree, next\_proc$ )
    end for
  end if

```

4 Proposed mapping refinement

Our objective is to correct the potential load imbalance (and thus idle times) created by the proportional mapping, as outlined in Section 2, but without impacting too much the data locality. We propose a heuristic based on work stealing [4] that extends the refined mapping phase (see Algorithm 4) using simulation (see Algorithm 5). Intuitively, we propose that if the simulation predicts that a processor will be idle, this processor tries to steal some tasks from its neighbors.

In the proposed refinement, we replace the update of the ready and waiting queues of the last line in Algorithm 4 by a call to *UpdateQueuesWithStealing* (Algorithm 5). For each processor k , we first detect if k will have some idle time, and we compute the duration d of this idle slot. This happens in particular when the ready time of the first task in its waiting queue is strictly larger than the ready time of the processor ($RT[i] > RP[k]$) and ready queue is empty. Whenever both queues are empty, the processor will be idle forever, and thus

Algorithm 4 Precise scheduling and mapping using simulation

```

for all task  $i$  do
  If  $i$  is a leaf, put  $i$  in the ready queue of every processor in  $candidate(i)$ , otherwise
  put it in the waiting queue.
end for
while all tasks have not been mapped do
  For each processor  $k$ , consider the triplet  $\langle i, k, t \rangle$  where  $i$  is the first task in
  the ready queue of processor  $k$  and  $t$  is the starting time of  $i$  on  $k$  ( $t =$ 
 $\max(RT[i], RP[k])$ )
  Consider  $F$ , the set of all such triplets
  Select the triplet  $\langle i, k, t \rangle$  in  $F$  with the smallest  $t$  (if ties, choose the one with
  largest depth)
  Schedule task  $i$  on processor  $k$  at time  $t$ 
  Update the ready times of processor  $k$  and of the successors of  $i$  on all their
  candidate processors
  Update the ready queue and waiting queue of processor  $k$ , as well as of candidates
  processors of successors of  $i$ 
end while

```

d is set to a large value. Then, if an idle time is detected (the ready queue is empty and d is a positive value), a task is stolen from a neighbor processor using function *StealTask*. Otherwise, the ready and waiting queues are updated as previously: the tasks of the waiting queue that will be freed before the processor becomes available are moved to the ready queue.

When stealing tasks, we distinguish between two cases, depending whether we use shared or distributed memory. In shared memory, the two possible victims of the task stealing operation are the two neighbors of processor k , considering that processors are arranged in a ring. In the case of distributed memory, we first try to steal from two neighbor processors *within the same cluster*, that is, within the set of processors that share the same memory. Stealing to a distant processor is considered only when clusters are reduced to a single element. Once steal victims are identified (set S), we consider the first task of their ready queues and select the one that can start as soon as possible. If the task is able to start during the idle slot of processor k (and thus reduce its idle time), it is then copied into its ready queue.

5 Experimental results

Experiments were conducted on the *Plafrim*⁵ supercomputer, and more precisely on the *miriel* cluster. Each node is equipped with two INTEL XEON E5-2680v3 12-cores running at 2.50 GHz and 128 GB of memory. The INTEL MKL 2019 library is used for sequential BLAS kernels. Another shared memory experiment was performed on the *crunch* cluster from the LIP⁶, where a node is equipped

⁵ <https://www.plafrim.fr>

⁶ <http://www.ens-lyon.fr/LIP/>

Algorithm 5 Update ready and waiting queues with task stealing

```

function UpdateQueuesWithStealing(nb. of proc.  $p$ , switch  $IsSharedMem$ ):
  for  $k = 1$  to  $p$  do
    if  $waiting\_queue_k \neq \emptyset$  then
      Let  $i$  be the first task in  $waiting\_queue_k$ 
       $d \leftarrow RT[i] - RP[k]$ 
    else
       $d \leftarrow +\infty$ 
    end if
    if  $ready\_queue_k = \emptyset$  and  $d > 0$  then
      StealTask( $k, p, d, IsSharedMem$ )
    else
      Let  $i$  be the first task in  $waiting\_queue_k$ 
      while  $RT[i] \leq RP[k]$  do
        Move task  $i$  from  $waiting\_queue_k$  to  $ready\_queue_k$ 
        Let  $i$  be the first task in  $waiting\_queue_k$ 
      end while
    end if
  end for

function StealTask(proc.  $k$ , proc. nb.  $p$ , idle time  $d$ , switch  $IsSharedMem$ ):
  set  $S \leftarrow \emptyset$ 
  if  $IsSharedMem = false$  then
    set  $S_k \leftarrow \{k - 1, k + 1, k - 2, k + 2\}$ 
    for  $j = 1$  to 4 do
      if  $S_k[j] \geq 0$ ,  $S_k[j] < p$ ,  $S_k[j]$  is in the same cluster as  $k$  and  $|S| < 3$  then
        add  $S_k[j]$  to  $S$ 
      end if
    end for
  end if
  if  $IsSharedMem = true$  or  $S$  is empty then
    set  $S \leftarrow \{k - 1 \pmod{p}, k + 1 \pmod{p}\}$ 
  end if
  Build the set  $O$  with the first element of each ready queue of processors in  $S$ 
  Let  $o$  be the task of  $O$  with minimum  $RT[o]$ 
  if  $RT[o] < RP[k] + d$ , then insert  $o$  into  $ready\_queue_k$ 

```

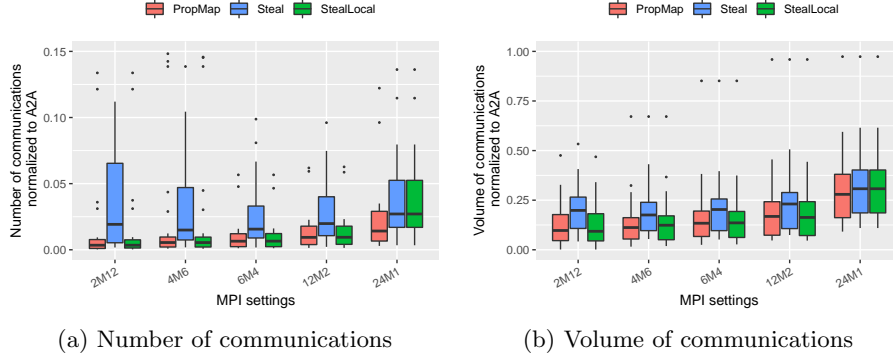


Figure 3: MPI communication number (*left*) and volume (*right*) for the three methods: PROPMAP, STEAL, and STEALLOCAL, with respect to ALL2ALL.

with four INTEL XEON E5-4620 8-cores running at 2.20 GHz and 378 GB of memory. On this platform, the INTEL MKL 2018 library is used for sequential BLAS kernels. The PASTIX version used for our experiments is based on the public git repository⁷ version at the tag 6.1.0.

In the following, the different methods used to compute the mapping are compared. All to All, referred to as ALL2ALL, and Proportional mapping, referred to as PROPMAP, are available in the PASTIX library, and the newly introduced method is referred to as STEAL. When the option to limit stealing tasks into the same MPI is enabled, we refer to it as STEALLOCAL. In all the following experiments, we compare these versions with respect to the ALL2ALL strategy, which provides the most flexibility to the scheduling algorithm to perform load balance, but does not consider data locality. The multi-threaded variant is referred to as SharedMem, while for the distributed settings, pMt stands for p MPI nodes with t threads each. All distributed settings fit within a single node.

In order to make a fair comparison between the methods, we use a set of 34 matrices issued from the SuiteSparse Matrix collection [6]. The matrix sizes range from 72K to 3M of unknowns. The number of floating point operations required to perform the LL^t , LDL^t , or LU factorization ranges from 111 GFlops to 356 TFlops, and the problems are issued from various application fields.

Communications. We first report the relative results in terms of communications among processors in different clusters (MPI nodes), which are of great importance for the distributed memory version. The number and the volume of communications normalized to ALL2ALL are depicted in Fig. 3a and Fig. 3b respectively. One can observe that all three strategies largely outperform the ALL2ALL heuristic, which does not take communications into account. The number of communications especially explodes with ALL2ALL as it mainly moves

⁷ <https://gitlab.inria.fr/solverstack/pastix>

around leaves of the elimination tree. This creates many more communications with a small volume. This confirms the need for a proportional-mapping-based strategy to minimize the number of communications. Both numbers and volumes of communications also confirm the need for the local stealing algorithm to keep it as low as possible. Indeed, STEAL generates 6.19 times more communications on average than PROPMAP, while STEALLOCAL is as good as PROPMAP. Note the exception of the 24M1 case where STEAL and STEALLOCAL are identical. No local task can be stolen. These conclusions are similar when looking at the volume of communication with a ratio reduced to 1.92 between STEAL and PROPMAP.

The distribution of the results is shown by boxplot. It shows five summary statistics (median, two hinges and two whiskers), and all "outlying" points individually. The lower and upper hinges correspond to the first and third quartiles. The upper whisker extends from the hinge to the largest value no further than $1.5 \times \text{IQR}$ from the hinge (where IQR is the inter-quartile range, or distance between the first and third quartiles). The lower whisker extends from the hinge to the smallest value at most $1.5 \times \text{IQR}$ of the hinge. Data beyond the end of the whiskers are called "outlying" points and are plotted individually [14].

Data movements. Fig. 4 depicts the number and volume of data movements normalized to ALL2ALL and summed over all the MPI nodes with different MPI settings. The data movements are defined as a write operation on the remote memory region of other cores of the same MPI node. Note that accumulations in local buffers before send, also called fan-in in sparse direct solvers, are always considered as remote write. This explains why all MPI configurations have equivalent number of data movements. As expected, proportional mapping heuristics outperform ALL2ALL by a large factor on both number and volume, which can have an important impact on NUMA architectures. Compared to PROPMAP, STEAL and STEALLOCAL are equivalent and have respectively 1.38x, and 1.32x, larger number of data movements on average respectively, which translates into 9%, and 8% of volume increase. Note that in the shared memory case, STEALLOCAL behaves as STEAL as there is only one MPI node.

Simulation cost. Fig. 5 shows the simulation cost in seconds (duration of the refined mapping via simulation) on the left, and that of PROPMAP, STEAL and STEALLOCAL with respect to ALL2ALL on the right. Fig. 6 shows the original simulated factorization time obtained with these heuristics and a normalized version. Note that, for the sake of clarity, some large outliers are removed from Fig. 6a. As stated in Section 2, the ALL2ALL strategy allows for more flexibility in the scheduling, hence it results in a better simulated time for the factorization in average. However, its cost is already 4x larger for this relatively small number of cores. Fig. 5a shows that the proposed heuristics have similar simulation cost to the original PROPMAP, while Fig. 6 shows that the simulated factorization time gets closer to ALL2ALL, and can even outperform it in extreme cases. Indeed, in the 24M1 case, STEAL outperforms ALL2ALL due to bad decisions taken by the latter at the beginning of the scheduling. The bad mapping of the leaves is then never recovered and induces extra communications that explain

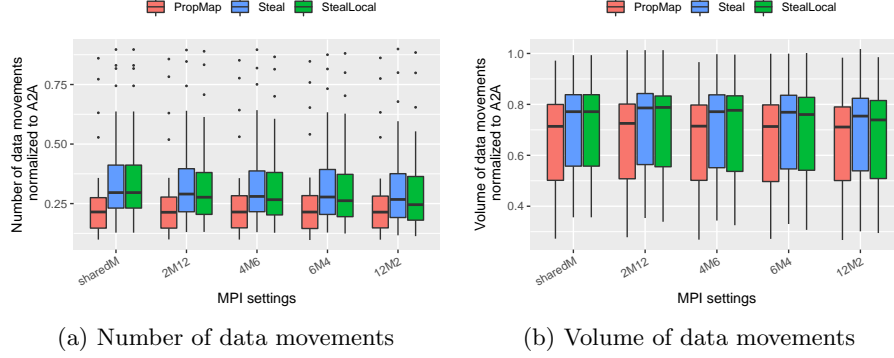


Figure 4: Shared memory data movements number (*left*) and volume (*right*) within MPI nodes for PROPMap, STEAL, and STEALLOCAL, with respect to ALL2ALL.

this difference. In conclusion, the proposed heuristic, STEALLOCAL, manages to generate better schedules with a better load-balance than the original PROPMap heuristic, while generating small or no overhead on the mapping algorithm. This strategy is also able to limit the volume of communications and data movements as expected.

Factorization time for shared memory. Fig. 7 presents factorization time and its normalized version in a shared memory environment, on both *miriel* and *crunch* machines. Note that we present only the results for STEAL, as STEALLOCAL and STEAL behave similarly in shared memory environment. For the sake of clarity, some large outliers are removed from Fig. 7a. On *miriel*, with a smaller number of cores and less NUMA effects, all these algorithms have almost similar factorization time, and present variations of a few tens of GFlop/s over 500GFlop/s in average. STEAL slightly outperforms PROPMap, and both are slower than ALL2ALL respectively by 1% and 2% in average. On *crunch*, with more cores and more NUMA effects, the difference between STEAL and PROPMap increases in favor of STEAL. Both remain slightly behind ALL2ALL, respectively by 2% and 4%; indeed, ALL2ALL outperforms them since it has the greatest flexibility, and communications have less impact in a shared memory environment.

6 Conclusion

In this paper, we revisit the classical mapping and scheduling strategies for sparse direct solvers. The goal is to efficiently schedule the task graph corresponding to an elimination tree, so that the factorization time can be minimized. Thus, we aim at finding a trade-off between data locality (focus of the traditional

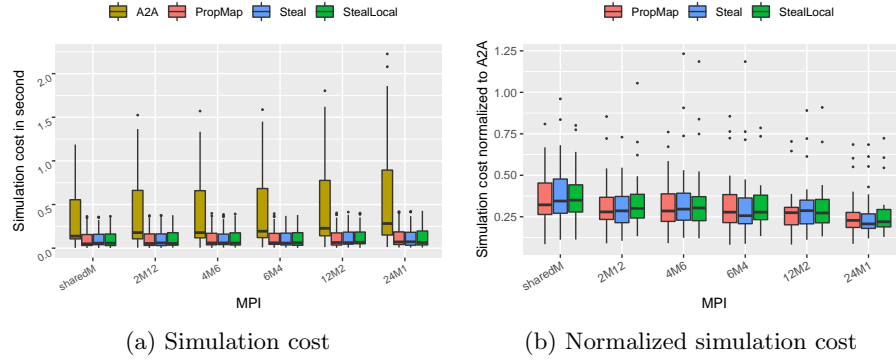


Figure 5: Final simulation cost in second (*left*) and simulation cost of PROPMap, STEAL and STEALLOCAL, normalized to ALL2ALL (*right*).

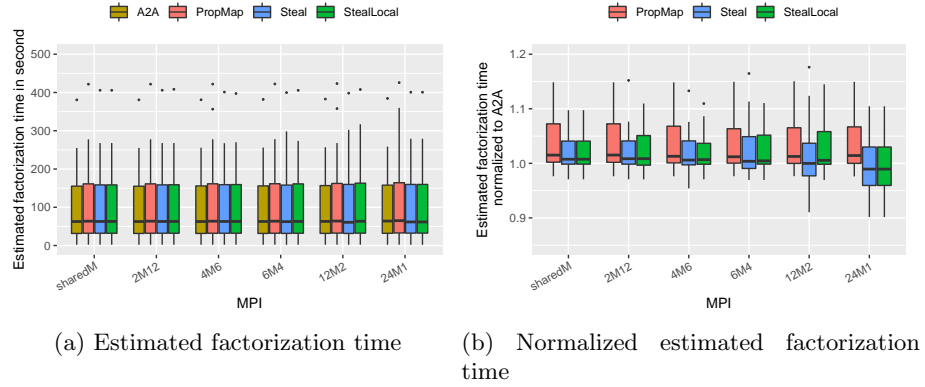


Figure 6: Estimated factorization time (*left*), and that of PROPMap, STEAL, and STEALLOCAL, normalized to ALL2ALL (*right*).

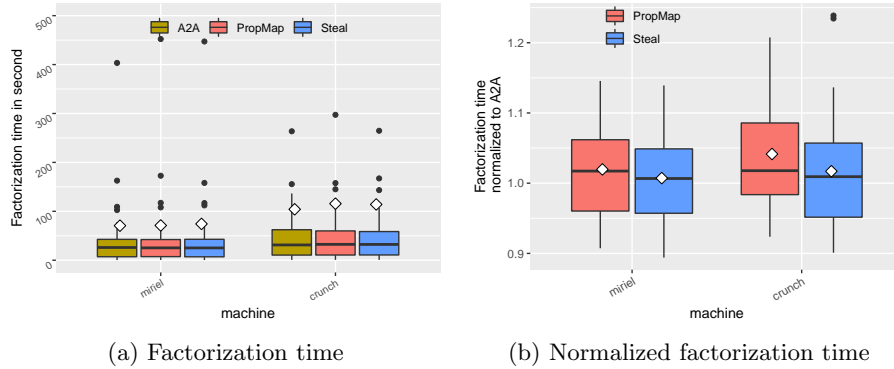


Figure 7: Factorization time (*left*), and that of PROPMAP and STEAL, normalized to ALL2ALL (*right*), on *miriel* and *crunch*. White diamonds represent mean values.

PROPMAP strategy) and load balancing (focus of the ALL2ALL strategy). First, we improve upon PROPMAP by proposing a refined (and optimal) mapping strategy with an integer number of processors. Next, we design a new heuristic, STEAL, together with a variant STEALLOCAL, which predicts processor idle times in PROPMAP and assigns tasks to idle processors. This leads to a limited loss of locality, but improves the load balance of PROPMAP.

Extensive experimental and simulation results, both on shared memory and distributed memory settings, demonstrate that the STEAL approach generates almost the same number of data movements than PROPMAP, hence the loss in locality is not significant, while it leads to better simulated factorization times, very close to that of ALL2ALL, hence improving the load balance of the schedule.

PASTIX has only recently been extended to work on distributed settings, and hence we plan to perform further experiments on distributed platforms, in order to assess the performance of STEAL on the numerical factorization in distributed environments. Future working directions may also include the design of novel strategies to further improve performance of sparse direct solvers.

Acknowledgments. This work is supported by the Agence Nationale de la Recherche, under grant ANR-19-CE46-0009. Experiments presented in this paper were carried out using the PlaFRIM experimental testbed, supported by Inria, CNRS (LABRI and IMB), Université de Bordeaux, Bordeaux INP and Conseil Régional d’Aquitaine (<https://www.plafrim.fr/>).

References

1. Amestoy, P.R., Buttari, A., Duff, I.S., Guermouche, A., L’Excellent, J.Y., Uçar, B.: MUMPS. In: Padua, D. (ed.) *Encyclopedia of Parallel Computing*, pp. 1232–1238. Springer (2011). https://doi.org/10.1007/978-0-387-09766-4_204

2. Amestoy, P.R., Duff, I.S., L'Excellent, J.Y.: Multifrontal parallel distributed symmetric and unsymmetric solvers. *Computer Methods in Applied Mechanics and Engineering* **184**(2), 501 – 520 (2000)
3. Beaumont, O., Guermouche, A.: Task scheduling for parallel multifrontal methods. In: *European Conference on Parallel Processing*. pp. 758–766. Springer (2007)
4. Blumofe, R.D., Leiserson, C.E.: Scheduling multithreaded computations by work stealing. *Journal of the ACM* **46**(5), 720–748 (1999)
5. Davis, T.A.: *Direct Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematics (2006). <https://doi.org/10.1137/1.9780898718881>
6. Davis, T.A., Hu, Y.: The University of Florida Sparse Matrix Collection. *ACM Trans. Math. Softw.* **38**(1) (Dec 2011). <https://doi.org/10.1145/2049662.2049663>
7. Duff, I.S., Reid, J.K.: The multifrontal solution of indefinite sparse symmetric linear. *ACM Trans. Math. Softw.* **9**(3), 302–325 (Sep 1983)
8. George, A., Liu, J.W., Ng, E.: Communication results for parallel sparse Cholesky factorization on a hypercube. *Parallel Computing* **10**(3), 287–298 (1989)
9. Gou, C., Al Zoobi, A., Benoit, A., Faverge, M., Marchal, L., Pichon, G., Ramet, P.: Improving mapping for sparse direct solvers: A trade-off between data locality and load balancing. Research report RR-9328, Inria (2020), hal.inria.fr
10. Heath, M.T., Ng, E., Peyton, B.W.: Parallel algorithms for sparse linear systems. *SIAM Rev.* **33**(3), 420–460 (Aug 1991). <https://doi.org/10.1137/1033099>
11. Hénon, P., Ramet, P., Roman, J.: PaStiX: A High-Performance Parallel Direct Solver for Sparse Symmetric Definite Systems. *Parallel Computing* **28**(2), 301–321 (Jan 2002)
12. Jacquelin, M., Zheng, Y., Ng, E., Yelick, K.A.: An Asynchronous Task-based Fan-Both Sparse Cholesky Solver. *CoRR* (2016), <http://arxiv.org/abs/1608.00044>
13. Li, X.S., Demmel, J.W.: SuperLU-DIST: A Scalable Distributed-Memory Sparse Direct Solver for Unsymmetric Linear Systems. *ACM Trans. Math. Softw.* **29**(2), 110–140 (Jun 2003). <https://doi.org/10.1145/779359.779361>
14. McGill, R., Tukey, J.W., Larsen, W.A.: Variations of box plots. *The American Statistician* **32**(1), 12–16 (1978), <http://www.jstor.org/stable/2683468>
15. Pothén, A., Sun, C.: A mapping algorithm for parallel sparse Cholesky factorization. *SIAM Journal on Scientific Computing* **14**(5), 1253–1257 (1993)
16. Prasanna, G.S., Musicus, B.R.: Generalized multiprocessor scheduling and applications to matrix computations. *IEEE Transactions on Parallel and Distributed systems* **7**(6), 650–664 (1996)
17. Rothburg, E., Gupta, A.: An Efficient Block-Oriented Approach to Parallel Sparse Cholesky Factorization. In: *Proceedings of the 1993 ACM/IEEE Conference on Supercomputing*. pp. 503–512 (1993). <https://doi.org/10.1145/169627.169791>