

GSoFA: Scalable Sparse Symbolic LU Factorization on GPUs

Anil Gaihre, Xiaoye Sherry Li and Hang Liu

Abstract—Decomposing a matrix \mathbf{A} into a lower matrix \mathbf{L} and an upper matrix \mathbf{U} , which is also known as LU decomposition, is an essential operation in numerical linear algebra. For a sparse matrix, LU decomposition often introduces more nonzero entries in the \mathbf{L} and \mathbf{U} factors than in the original matrix. A *symbolic factorization* step is needed to identify the nonzero structures of \mathbf{L} and \mathbf{U} matrices. Attracted by the enormous potentials of the Graphics Processing Units (GPUs), an array of efforts have surged to deploy various LU factorization steps except for the symbolic factorization, to the best of our knowledge, on GPUs. This paper introduces GSoFA, the first GPU-based symbolic factorization design with the following three optimizations to enable scalable LU symbolic factorization for *nonsymmetric pattern* sparse matrices on GPUs. First, we introduce a novel fine-grained parallel symbolic factorization algorithm that is well suited for the *Single Instruction Multiple Thread* (SIMT) architecture of GPUs. Second, we tailor supernode detection into a SIMT friendly process and strive to balance the workload, minimize the communication and saturate the GPU computing resources during supernode detection. Third, we introduce a three-pronged optimization to reduce the excessive space consumption problem faced by multi-source concurrent symbolic factorization. Taken together, GSoFA achieves up to $31\times$ speedup from 1 to 44 Summit nodes (6 to 264 GPUs) and outperforms the state-of-the-art CPU project, on average, by $5\times$. Notably, GSoFA also achieves up to 47% of the peak memory throughput of a V100 GPU in Summit.

Index Terms—Sparse linear algebra, sparse linear solvers, LU decomposition, static symbolic factorization on GPU

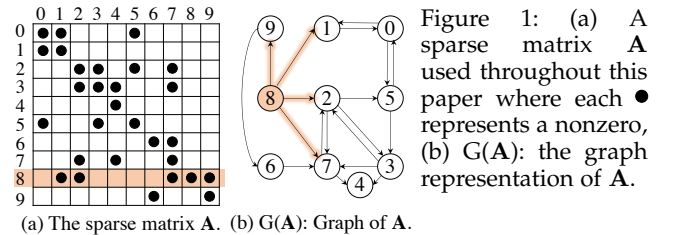
arXiv:2007.00840v4 [cs.DC] 9 May 2021

1 INTRODUCTION

Many scientific and engineering problems require solving large-scale linear systems, i.e., $\mathbf{Ax} = \mathbf{b}$. Solving this problem with direct methods [1] often involves LU factorization [2], [3], that is, decomposing the original matrix \mathbf{A} into lower and upper triangular matrices \mathbf{L} and \mathbf{U} , respectively, where $\mathbf{A} = \mathbf{LU}$. Since LU decomposition of a sparse matrix typically introduces new nonzeros, also known as *fill-ins*, in the \mathbf{L} and \mathbf{U} factored matrices, *symbolic factorization* [2] is used to compute the locations of the fill-ins for both \mathbf{L} and \mathbf{U} matrices. This information is needed to allocate the compressed sparse data structures for \mathbf{L} and \mathbf{U} , and for the subsequent numerical factorization.

Symbolic factorization is a graph algorithm that acts on the adjacency graph of the corresponding sparse matrix. Figure 1 shows a sparse matrix \mathbf{A} and its adjacency graph $G(\mathbf{A})$. In the graph, each row of a matrix corresponds to a vertex in the directed graph with a nonzero in that row corresponding to an out edge of that vertex. For instance, row 8 of \mathbf{A} in Figure 1(a) has nonzeros in columns $\{1, 2, 7, 8, 9\}$. We thus have those corresponding out edges for vertex 8 in $G(\mathbf{A})$ of Figure 1(b). For brevity, we do not represent the self edges corresponding to the diagonal elements. We will use this example throughout this paper.

The motivation to parallelize the sparse symbolic factorization on GPUs is threefold. Firstly, the newer generation



GPUs have a relatively large amount of device memory, which allows the entire sparse matrix to reside on the GPUs. Secondly, more and more exascale application codes are moving to GPUs, which demands the underlying linear solvers to reside on GPUs to eliminate CPU-GPU communication. Ultimately, we will move all the workflows of the sparse LU solver to GPUs. Third, many GPU-based efforts have gone to deploying numerical factorization and triangular solution on GPUs [4], [5] due to (i) numerical factorization and triangular solution take more time, and (ii) symbolic factorization consumes excessive memory space. We argue that since sparse symbolic factorization detects the fill-ins and supernodes for the numerical steps, it is imperative to move sparse symbolic factorization on GPUs.

Existing symbolic factorization algorithms are difficult to deploy on GPUs, as well as scale-up in distributed settings, due to stringent data dependency, lack of fine-grained parallelism, and excessive memory consumption. There mainly exist three types of algorithms to perform symbolic factorization, that is, fill1 and fill2 algorithms by Rose and Tarjan [6] and Gilbert-Peierls algorithm [7], [8]. They are based on traversing either the graphs $G(\mathbf{L})$ and $G(\mathbf{U})$ or graph $G(\mathbf{A})$. In particular, fill1 works on $G(\mathbf{L})$ and $G(\mathbf{U})$ for fill-ins detection while Gilbert-Peierls algorithm does that on $G(\mathbf{L})$. Fill2 algorithm is analogous to Dijkstra's algo-

- A. Gaihre and H. Liu are with Department of Electrical and Computer Engineering, Stevens Institute of Technology, NJ.
E-mail: {agaihre, hliu77}@stevens.edu
- X.S. Li is with Lawrence Berkeley National Laboratory, Berkeley, CA.
E-mail: xslil@lbl.gov

Manuscript revised May 7, 2021.

rithm [9]. Consequently, fill1 and Gilbert-Peierls algorithms are sequential in nature as the fill-in detection of larger rows or columns has to wait until the completion of the smaller rows or columns. Fill2 algorithm lacks fine-grained parallelism (i.e., in each source) due to strict priority requirement, although coarse-grained parallelism (i.e., at source level) exists. Furthermore, symbolic factorization needs to identify supernode [10], which is used to expedite the numerical steps. But supernode detection presents cross-source data dependencies. Third, we need to perform parallel symbolic factorization for a collection of rows to saturate the computing resources on GPUs. This results in overwhelming space consumption, which is similar to multi-source graph traversal [11].

In this paper, we choose the fill2 algorithm as a starting point stemming from the fact that fill2 directly traverses the original instead of the filled graphs, allowing independent traversals from multiple sources. However, we have to revamp this algorithm to expose a higher degree of fine-grained parallelism to match GPUs’ SIMT nature and reduce the excessive memory requirements during the multi-source symbolic factorization to, again, fit GPUs, which often equip limited memory space. To summarize, we make the following three major contributions.

First, to the best of our knowledge, we introduce the first fine-grained parallel symbolic factorization that is well suited for the SIMT architecture of GPUs. Particularly, instead of processing one vertex at a time, we allow all the neighboring vertices to be processed in parallel, such that the computations align well with the SIMT nature. *Since this relaxation enables traversal to some vertices before their dependencies are satisfied, we further allow revisitation of these vertices to ensure correctness.* We also introduce optimizations to reduce the reinitialization overhead and random memory access to the algorithmic data. Finally, we use multi-source concurrent symbolic factorization to saturate each GPU and achieve intra-GPU workload balance.

Second, we not only tailor supernode detection into a SIMT friendly process but also strive to *balance the workload, minimize the communication and saturate GPU computing resources* during supernode detection. Particularly, we break supernode detection into two phases to expose massive parallelism for GPUs. Further, we assign a chunk of continuous sources to one computing node to avoid inter-node communication and interleave the sources of each chunk across GPUs in each node to balance the workload. Eventually, we investigate the configuration space of unified memory and propose a design that can significantly reduce the page faults during supernode detection.

Third, we introduce a three-pronged optimization to combat the excessive space consumption problem faced by multi-source concurrent symbolic factorization: (i) We propose the external GPU frontier management because the space requirement of frontier-related data structures is very dynamic, and their access pattern is predictable. (ii) We identify and remove the “bubbles” in vertex status-related data structures; and (iii) since various data structures present dynamic space requirements with respect to different sources, we propose allocating single memory space for all these data structures and dynamically adjust their capacities to reduce the frequency of copying frontier-related data

structures between CPU and GPU memories.

The rest of this paper is organized as follows: Section 2 introduces the background of this work. Section 3 presents the design challenges. Sections 4, 5 and 6 present the fine-grained symbolic factorization algorithm design, parallel efficient supernode detection, and space consumption optimization techniques. We evaluate GSOFA in Section 7, study the related work in Section 8, and conclude in Section 9.

2 BACKGROUND

2.1 Sparse LU Factorization

Factorizing a sparse matrix \mathbf{A} into the \mathbf{L} and \mathbf{U} matrices often involves several major steps, such as preprocessing, symbolic factorization, and numerical factorization.

Matrix preprocessing performs either row (partial pivoting) or both row and column permutations (complete pivoting) in order to *improve numerical stability and reduce the number of fill-ins* in the \mathbf{L} and \mathbf{U} matrices. For numerical stability, existing methods aim to find a permutation so that the pivots on the diagonal are big, for instance, by maximizing the product of diagonal entries and make the matrix diagonal dominant [12]. Supplemental file illustrates how matrix preprocessing works with matrix \mathbf{A} . Regarding fill-in reduction, the popular strategies are minimum degree algorithm [13], Reverse Cuthill-McKee [14] and nested dissection [15], or a combination of them, such as METIS [16].

Symbolic factorization. Symbolic factorization can be viewed as performing Gaussian elimination on a sparse matrix with either 0 or 1 entries. At each elimination step, the pivoting row is multiplied by a nonzero scalar and updated into another row below, which has nonzero in the pivoting column; this is called an elementary row operation. After a sequence of row operations, the original matrix is transformed into two triangular matrices \mathbf{L} and \mathbf{U} . Considering row 8 in Figure 2(a), since row 8 has nonzeros at columns 1 and 2, both rows 1 and 2 will perform row operations on row 8. Using row operation from row 2 as an example, it produces a fill-in (8, 3), which further triggers row operations from rows 3 and 5 that produce fill-ins (8, 4) and (8, 5). We will detail this process shortly. The following fill-path theorem succinctly characterizes the fill-in locations reflecting the above elimination process.

Theorem 1. A fill-in at index (i, j) is introduced if and only if there exists a directed path from i to j , with the intermediate vertices being smaller than both i and j [6].

Theorem 1 can be applied on either the (partially) filled graph $G(\mathbf{U})$ by the fill1 algorithm or the original graph $G(\mathbf{A})$ by the fill2 algorithm [6].

Figure 3(a) shows how fill1 works for the src -th row in three steps: (i) lines 3-9 initialize the fill(:) and frontierQueue(:). While the usages of frontierQueue(:) and newFrontierQueue(:) are straightforward, fill(:) is a flag array that tracks which vertex is visited by src vertex. (ii) If the neighbor vertex is first-time visited (i.e., lines 13-16), we mark this neighbor as visited by setting fill(neighbor) = src at line 13, and add *neighbor* to either $\mathbf{L}(src, :)$ or $\mathbf{U}(src, :)$ based upon the location of $(src, neighbor)$. (iii) If *neighbor* is smaller than src , we add it to *newFrontierQueue(:)* for next iteration traversal. Figure 2(b) shows how row 8 works

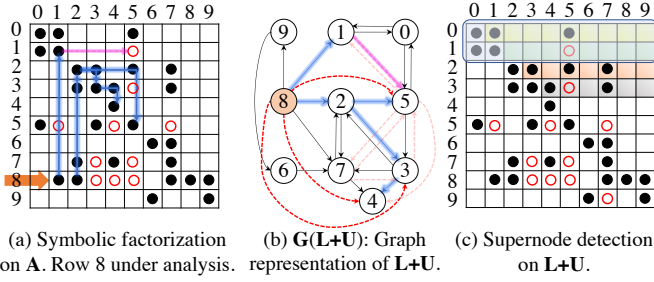


Figure 2: (a) Symbolic factorization on sparse matrix with row 8 under analysis \mathbf{A} , (b) the Graph representation of filled matrix $(\mathbf{L}+\mathbf{U})$ with row 8 under analysis and (c) Supernode detection on filled matrix. The \circ is the new nonzeros (fill-ins) in the \mathbf{L} and \mathbf{U} factors. And $\dashrightarrow/\dashrightarrow$ indicates the new edges in the $G(\mathbf{L}+\mathbf{U})$ graph. The dark edges with blue or pink glows in (a) and (b) show the traversal paths that lead to new fill-ins.

under fill1 algorithm. That is, $8 \rightarrow 2 \rightarrow 3$ leads to $(8, 3)$, $8 \rightarrow 2 \rightarrow 3 \rightarrow 4$ leads to $(8, 4)$, $8 \rightarrow 1 \dashrightarrow 5$ leads to $(8, 5)$. Note that the third path goes through a new fill-in edge $(1,5)$.

Figure 3(b) presents the fill2 algorithm. Similar to fill1, it uses $\text{fill}(\cdot)$ array to indicate an already visited vertex by setting $\text{fill}(\text{neighbor}) = \text{src}$. The algorithm differs from fill1 mainly on two points. (i) It uses the original graph $G(\mathbf{A})$ for traversal (line 12). Second, during traversal, this algorithm only permits traversing one vertex (i.e., threshold) at a time, starting from the smallest one (line 9 - 10). For each threshold that is treated as a frontier, this algorithm checks its neighbors, updates the statuses of the neighbors in $\text{fill}(\cdot)$, and adds new fills to either $\mathbf{L}(\text{src}, \cdot)$ or $\mathbf{U}(\text{src}, \cdot)$, as well as to the $\text{newFrontierQueue}(\cdot)$ if this neighbor obeys Theorem 1. This process continues until the vertices that are smaller and connected to the threshold vertex are exhausted, i.e., $\text{frontierQueue}(\cdot)$ is empty. Subsequently, fill2 will proceed to the next threshold vertex in line 9. Considering row 8 in Figure 2(b) again, the three fill-ins are due to the three paths going through only existing edges: $8 \rightarrow 2 \rightarrow 3$ leads to $(8, 3)$, $8 \rightarrow 2 \rightarrow 3 \rightarrow 4$ leads to $(8, 4)$, $8 \rightarrow 2 \rightarrow 5$ leads to $(8, 5)$.

A decade later, [8] introduces the Gilbert-Peierls algorithm to find the fill-in structures, which is a simpler way of interpreting fill1 algorithm. Specifically, this approach determines the nonzero structures column by column. For clarity, we define that $\mathbf{L}(i : j, m : n)$ and $\mathbf{U}(i : j, m : n)$ denote the block from rows i to j , and columns m to n in \mathbf{L} and \mathbf{U} matrices, respectively. For column k , it traverses the graph $\mathbf{L}(:, 0 : k - 1)^T$ in a Depth-First Search (DFS) manner. The vertex that is reachable by the vertices in column k results in a fill-in at column k . The graph used by [8] is similar to that of fill1 except that [8] applies a transpose to the \mathbf{L} matrix. Further, the vertices that can be reached from the source vertices will automatically satisfy Theorem 1 because only the graph of $\mathbf{L}(:, 0 : k - 1)^T$ is used.

It is worth noting that the nonsymmetric-pattern sparse LU symbolic factorization is much harder than the symmetric-pattern counterpart: (i) In the symmetric case, the transitive reduction of the filled graph $G(\mathbf{L})$ is a tree, called elimination tree; symbolic factorization using the elimination tree can be done in time $O(\text{nnzeros}(\mathbf{L}))$, linear to the output size. (ii) However, for the nonsymmetric cases, the transitive reduction of the filled graphs $G(\mathbf{L})$ and $G(\mathbf{U})$ are Directed Acyclic Graphs (DAGs), called elimination DAGs. Computing these DAGs is expensive, and

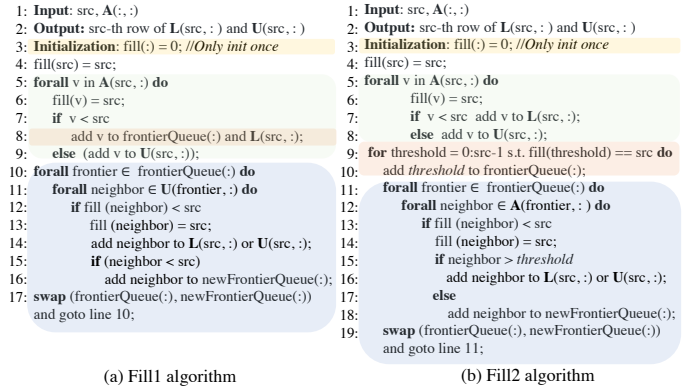


Figure 3: Fill1 and fill2 algorithms. We use MATLAB format to denote both arrays and matrices. The snippet with the same background colors in the algorithms represent equivalent sections of these two algorithms.

all variants of this method on the nonsymmetric symbolic factorization algorithms take asymptotically longer than linear time [6], [8].

Numerical factorization. After the structure of the fill-ins is determined, the solvers perform numerical factorization to calculate the values of the \mathbf{L} and \mathbf{U} matrices. Popular numerical factorization methods are *left-looking* [8] and *right-looking* [2]. We explain how numerical factorization works on matrix \mathbf{A} in supplemental file.

Triangular solution. Once LU factorization arrives at $\mathbf{A} = \mathbf{L}\mathbf{U}$, we can solve $\mathbf{A}\mathbf{x} = \mathbf{b}$ in two steps with triangular solver, given $\mathbf{A}\mathbf{x} = \mathbf{b}$ becomes $\mathbf{L}\mathbf{U}\mathbf{x} = \mathbf{b}$. First, triangular solver can solve $\mathbf{L}\mathbf{y} = \mathbf{b}$ to derive \mathbf{y} . Second, triangular solver derives \mathbf{x} by solving $\mathbf{U}\mathbf{x} = \mathbf{y}$.

2.2 Supernode

Symbolic factorization also needs to identify the supernodes in $G(\mathbf{L} + \mathbf{U})$, which is important to improve the performance for both numerical factorization and triangular solver. Particularly, a *supernode* is a range of rows or columns with the same nonzero structures, such as rows 1 and 2 in Figure 2(c), so that these rows or columns can be treated as a dense matrix [10]. During numerical factorization, a dense matrix format allows us to use Level 3 Basic Linear Algebra Subprograms (BLAS) operations such as matrix-matrix multiplication, which is often faster than the lower level BLAS operations like Level 2 BLAS operations. In this paper, gSOFA supports T3, one of the most popular types of supernodes among five types of supernodes [10]. Definition 1 defines a T3 supernode.

Definition 1. Supposing a T3 supernode begins at row r and extends through row $s - 1$, row s belongs to T3 if and only if $\text{nnz}(\mathbf{U}(s, \cdot)) = \text{nnz}(\mathbf{U}(s - 1, \cdot)) - 1$ and $\mathbf{L}(s, r) \neq 0$ according to [10].

Definition 1 states that for a supernode that already contains rows r through $s - 1$, two requirements are needed for the next row s to be included in this supernode: (i) The number of nonzeros in row s of \mathbf{U} shall be one fewer than that of row $s - 1$. (ii) There shall be a nonzero at (s, r) in the filled matrix \mathbf{L} . If either of the requirements is not met, row s will not belong to the supernode. Intuitively, (i) a nonzero at $\mathbf{L}(s, r)$ ensures that all the nonzero patterns of row r in

\mathbf{U} are mapped into row s during symbolic factorization. (ii) Further, because the nonzero count of row s is one fewer than that of row $s - 1$ in \mathbf{U} , we can conclude rows $s - 1$ and s follow the same nonzero pattern in \mathbf{U} . Thus, they belong to the same supernode.

Now, we use the example in Figure 2(c) to illustrate how to detect a supernode in $G(\mathbf{L} + \mathbf{U})$. Starting from row 0, we check whether (i) the number of nonzeros in row 1 of \mathbf{U} matrix (i.e. 2) is one fewer than row 0 (i.e., 3); and (ii) there is a nonzero at (1, 0). Since both requirements are met, row 1 is included in the supernode starting at row 0. Applying a similar process to row 2, we find that row 2 has neither one less nonzero than the previous row (i.e., row 1) nor a nonzero at (2, 0). Hence, row 2 starts a new supernode.

2.3 Graphics Processing Units

This section discusses general-purpose GPUs with recent NVIDIA V100 GPU [17] as an example.

Streaming processors and threads. The V100 GPU is designed with NVIDIA Volta architecture. V100 is powered by 80 Streaming Multiprocessors (SMX). Each SMX features 64 CUDA cores, resulting in a total of 5,120 CUDA cores. During execution, a GPU thread runs on one CUDA core. A SMX schedules a group of 32 consecutive threads known as a warp in a SIMT manner. A collection of consecutive warps further formulate a Cooperative Thread Array (CTA), or a block. All the CTAs together in one kernel are called a grid.

Memory architecture. V100 comes with two memory capacities, that is, 16 GB and 32 GB with peak bandwidth up to 900 GB/s. Each SMX has 96 KB on-chip fast memory that is shared by the configurable shared memory and L1 cache. All the SMXs share a L2 cache at a size of 6,144 KB. Each thread block can use up to 65,536 32-bit registers.

Fine-grained parallelism. GPUs favor fine-grained parallelism. Particularly, GPUs can only achieve the aforementioned ideal computing and memory throughput when a warp of threads is working on the same instruction and fetching data from consecutive memory addresses. Otherwise, GPUs might suffer from either warp divergence or uncoalesced memory access issues, resulting in order of magnitude performance degradation [17].

Unified memory [18] can unite the memory space of all GPUs and CPUs into a single virtual address space. During execution, any process or thread can access the data from this virtual address space. In the background, the required data is either transferred from where the data resides to where the data is needed implicitly or directly accessed remotely. This is different from the traditional explicit method, which requires programmers to explicitly call `cudaMemcpy()` (or similar functions) in order to transfer the data. Consequently, unified memory bests explicit transfer when the requested data is either small (in terms of size) or randomly accessed, or both. Further, unified memory can be used without terminating the kernel, which is not possible from the explicit method.

2.4 Dataset

Table 1 presents the datasets that are used to evaluate GSOFA. They are available from Suite Sparse Matrix Collection [19]. This dataset collection includes a variety of

Matrix (\mathbf{A})	Abbr.	Order (\mathbf{A})	nnz (\mathbf{A})	Struct. symm.	nnz(\mathbf{A}) / Order(\mathbf{A})	#Fill-in nnz(\mathbf{A})
BBMAT	BB	38,744	1,771,722	0.53	45.7	18.29
BCSSTK18	BC	11,948	149,090	1	12.47	6.52
EPB2	EP	25,228	175,027	0.67	6.93	9.28
G7JAC200SC	G7	59,310	717,620	0.03	12.1	24.51
LHR71C	LH	70,304	1,528,092	0	21.7	3.10
MARK3JAC	MK	64,089	376,395	0.07	5.9	28.59
RMA10	RM	46,835	2,329,092	1	49.729	3.14
AUDIKW_1	AU	943,695	77,651,847	1	82.28	31.43
DIELFILTER	DI	1,157,456	48,538,952	1	41.93	22.39
HAMRLE3	HM	1,447,360	5,514,242	0	3.8	32.63
PRE2	PR	659,033	5,834,044	0.33	8.8	20.70
STOMACH	ST	213,360	3,021,648	0.85	14.2	25.77
TWOTONE	TT	120,750	1,206,265	0.24	10	6.07

Table 1: Dataset specifications. Note, in graph terminology, order (\mathbf{A}) and nnz(\mathbf{A}) represent $|V|$ and $|E|$ of the graph $G(\mathbf{A})$, respectively, where \mathbf{A} is the matrix of interest.

applications, such as circuit simulation (HM, PR, and TT), structural problems (BC, AU), computational fluid dynamics (BB and RM), thermal problems (EP), economic modeling (G7 and MK), chemical engineering (LH), electromagnetics (DI) and bioengineering problems (ST). Besides, this dataset collection covers a wide range of variations in both structural symmetry and sparsity (the fifth and sixth columns in Table 1). This table arranges the datasets into smaller (upper) and larger (lower) collections with respect to order (\mathbf{A}). The larger matrices are used in space complexity analysis in Sections 6 and 7. Following SuperLU_DIST [20], we use ParMETIS [21] library to preprocess the matrix, and adopt Compressed Sparse Row (CSR) format to represent the matrices [22], [23].

3 DESIGN CHALLENGES

Challenge #1. Existing symbolic factorization algorithms present limited fine-grained parallelism.

Comparing the fill1 and fill2 algorithms, we can see that fill1 is severely limited in parallelism due to data dependency, because, to find the nonzero structure of the current row, we need to wait for the completion of detecting all the nonzero structures of the previous rows in \mathbf{U} . In comparison, the fill2 algorithm exhibits a high degree of parallelism: because graph $G(\mathbf{A})$ is static, we can perform parallel independent traversals from all vertices in $G(\mathbf{A})$. Therefore, an algorithm variant based on fill2 is more favorable for a massively parallel device like GPU.

Despite that fill2 presents more parallelism, fill2 also faces the limited *fine-grained parallelism issue* because line 9 - 10 of fill2 in Figure 3(b) only allows processing one threshold vertex at a time. It is important to note that *simply allowing fill2 to process multiple threshold vertices in parallel will not warrant the correctness* because the threshold value has to be strictly incremented sequentially. This is due to the constraint that each vertex can only be visited once in fill2. An example of this argument is discussed in Section 4.

To mitigate this bottleneck, GSOFA allows all frontiers to be processed in parallel, such that a warp of GPU threads can work on consecutive tasks. *Given this relaxation, some frontiers might be worked on earlier than they are supposed to,*

we introduce new data structures and control logics to allow revisitations in order to ensure correctness.

Challenge #2. Parallel supernode detection needs a good trade-off between load balance and communication reduction.

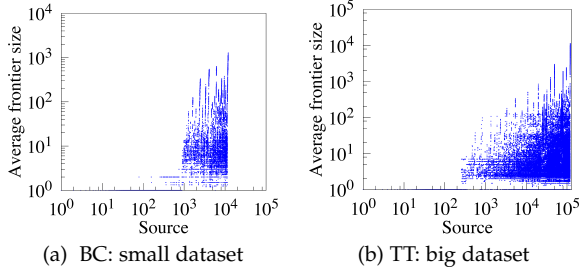


Figure 4: Average frontier size per source for BC and TT. Note, we only present two datasets for brevity.

Theorem 1 indicates that the amount of workload per source generally increases with the increase of source vertex ID, because more intermediate vertices will be smaller than the source ID when the source is larger. Our experimental results in Figure 4 corroborate with this. The general trend is that the workload soars with the increasing source ID. For instance, the workload ratios between the smallest and largest sources are $1,265\times$ and $6,230\times$ for BC and TT datasets, respectively.

We note that the optimization to balance the workload contradicts the optimization for supernode detection. On the one hand, as suggested by the workload dynamic of different sources in Figure 4, a proper workload balancing strategy would require to interleave the sources across GPUs. On the other hand, supernode detection has to check a continuous range of sources, implying that assigning a continuous range of sources to one GPU would minimize the communication cost. These contradictory goals require novel source scheduling and system optimizations.

In this paper, we first transform supernode detection into a massively parallel process that fits GPUs. Subsequently, we introduce a trade-off to both balance the workload and significantly reduce the communication cost during supernode detection with the help of a judicious source scheduling mechanism and an interesting unified memory-based data sharing design.

Challenge #3. The auxiliary data structures may consume overwhelming memory space.

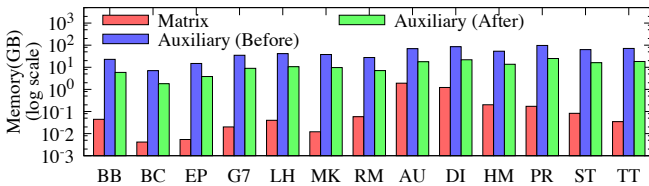


Figure 5: GSoFA memory consumption of the original matrix vs the auxiliary data structures on a Summit compute node.

While the fast runtime is important for symbolic factorization, so is the small space consumption. As shown in Figure 3(b), the original fill2 algorithm requires three auxiliary arrays, `fill(:)`, `frontierQueue(:)` and `newFrontierQueue(:)`.

```

1: Input: maxId(:), src, A(:,:)
2: Output: src-th row of L(src,:) and U(src,:)
3: Initialization: maxId(:)=maxVal; //Initialize every  $\frac{\text{maxVal}}{|V|}$  number of src
4: maxId(src) = 0;
5: for all v in A(src,:) fill(v) = 0, maxId(v) = 0, add v to either L(src,:) or U(src,:);
6: for all v in A(src,:) such that v < src, add src to frontierQueue(:);
7: for all frontier in frontierQueue(:) in parallel do //Fine-grained parallelism
8:   newMaxId = Maximum of frontier and maxId(frontier);
9:   for all neighbor in A(frontier,:) in parallel do //Fine-grained parallelism
10:    if ((atomicMin of maxId(neighbor) and newMaxId) > newMaxId)
11:      if neighbor > newMaxId
12:        if ((atomicMax of fill(neighbor) and src) < src)
13:          //if not detected as fill-in before
14:          add neighbor to either L(src,:) or U(src,:);
15:        else continue; //Avoid re-insertion to frontierQueue(:)
16:        if (neighbor < src) atomicAdd neighbor to newFrontierQueue(:)
16: swap (frontierQueue(:), newFrontierQueue(:));

```

Line 9.5: if fill(neighbor) == src continue; //Enable line 9.5 when there exist many fill-ins in the graph.

Figure 6: GSoFA parallel algorithm. The data structures are in MATLAB format. `L(src,:)` and `U(src,:)` represent the L and U structures of the row `src`.

When we improve the fill2 algorithm with combined frontier queue traversal, we need to add two new arrays: `tracker(:)` and `newTracker(:)`. In addition, we need `maxId(:)` to allow the revisitation mechanism to avoid false negatives. Altogether, each source needs six nontrivial auxiliary data structures, each of which consumes memory space of size $\mathcal{O}(|V|)$. Further, since we often need tens to thousands of concurrent sources to saturate a GPU (discussed in Section 4), the memory requirement of the initial version of GSoFA for certain matrices becomes significantly large in comparison to the available memory on GPUs. Figure 5 presents the GPU memory consumption for the sparse matrix and the auxiliary data structures (before and after space optimization) mentioned in Table 2 on one Summit node. One can observe that the memory requirement for the auxiliary data structures is orders of magnitude larger than that for the matrices. In particular, maximum and minimum memory consumption ratios for auxiliary data structure and matrices are $3121:1$ in MK and $36:1$ in AU, respectively.

Section 6 presents a series of optimizations to reduce the auxiliary data structures memory consumption while maintaining similar performance. As shown in Figure 5, these optimizations reduce the memory consumption of the auxiliary data structure by an average factor of $4\times$. The maximum and minimum ratios of memory consumption for auxiliary data structure and matrices are reduced to $801:1$ in MK and $9:1$ in AU.

4 A NEW FINE-GRAINED PARALLEL SYMBOLIC FACTORIZATION ALGORITHM

GSoFA addresses the first challenge, i.e., limited fine-grained parallelism within traversal for a single row, by *i) allowing parallel processing of the frontiers, (ii) adding a new data structure `maxId(:)` array of size $|V|$, and (iii) allowing revisitation of vertices*. Below we explain the data structure and algorithm design separately. Formally, from the source vertex `src` to a specific vertex `v`, there often exist multiple paths, and each path has a maximum numbered vertex. GSoFA uses `maxId(v)` to store the minimum of all the maximums of the paths from `src` to `v`.

GSoFA allows repeated updates to `maxId(:)` along the traversal. And the traversal terminates once the `maxId` values of all the vertices converge to their minimal value. As shown in Figure 6, in the beginning, the neighbors of

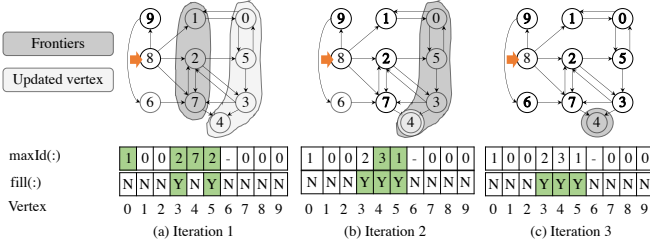


Figure 7: GSOF A traversing the graph from Figure 1(b). The dark, and light gray regions, respectively, represent the current frontier and vertices with $\text{maxId}(\cdot)$ updated. Vertex 8 is the source. The $\text{maxId}(\cdot)$ and $\text{fill}(\cdot)$ are at size of $|V|$. Note, Y and N in $\text{fill}(\cdot)$ are respectively used to show if the corresponding vertex is detected as fill-in or not.

src that are smaller than src are eligible for continuing traversal along these neighbors. Hence they are inserted into $\text{frontierQueue}(\cdot)$ at line 6. During traversal, the algorithm allows all the frontiers to explore the graph in parallel. With initial setting of 0 for every vertex, $\text{fill}(v)$ of a vertex v in the filled graph. The value of $\text{fill}(\cdot)$ helps avoid the re-detection of fill-ins at line 12. For a vertex v , $\text{fill}(v) < \text{src}$ means that vertex v is not yet in the filled structures of row src .

GSOF A fulfills two tasks in lines 9 - 15. The first task is to check whether a neighbor introduces a fill. The second task is to decide whether a vertex can become a frontier. For both tasks, we need the new path from the frontier to the current neighbor to change the maxId of this neighbor, that is, $\text{maxId}(\text{neighbor})$ should be updated by newMaxId at line 10. The function atomicMin updates the $\text{maxId}(\text{neighbor})$ by the minimum of $\text{maxId}(\text{neighbor})$ and newMaxId atomically at line 10. For the first task, this neighbor further needs to satisfy Theorem 1. And this neighbor should not have introduced a fill before, i.e., line 12. Otherwise, we face the issue of re-insertion of this neighbor into either L or U. For the second task, a neighbor further needs to meet another criterion in order to become a frontier, that is, this neighbor is smaller than the source (line 15). Otherwise, this neighbor either cannot be a frontier.

Example. To aid the understanding of GSOF A, Figure 7 demonstrates how GSOF A traverses from $\text{src} = 8$ on the graph of Figure 1(b). After line 5, the values in $\text{maxId}(\cdot)$ corresponding to the vertices $\{1, 2, 7, 8, 9\}$ become 0, while the rest remain as $|V|$ (i.e., ‘-’ in this case).

At iteration 1, the frontiers $\{1, 2, 7\}$, in parallel, traverse their neighbors $\{0\}$, $\{3, 5, 7\}$ and $\{2, 4\}$, arriving at updated $\text{maxId}(\cdot)$ and $\text{fill}(\cdot)$ array at the bottom of Figure 7(a). Subsequently, we obtain 1, 2, 7 and 2 as the maxId along the paths $8 \rightarrow 1 \rightarrow 0$, $8 \rightarrow 2 \rightarrow 3$, $8 \rightarrow 7 \rightarrow 4$ and $8 \rightarrow 2 \rightarrow 5$. Clearly, these paths introduce fills at 3 and 5 as shown in the $\text{fill}(\cdot)$. Proceeding to iteration 2, frontiers $\{0, 3, 4, 5\}$ generate new paths as $8 \rightarrow 1 \rightarrow 0 \rightarrow 5$, $8 \rightarrow 2 \rightarrow 3 \rightarrow 4$ and $8 \rightarrow 2 \rightarrow 5 \rightarrow 3$. Even path $8 \rightarrow 1 \rightarrow 0 \rightarrow 5$ updates $\text{maxId}(5)$ at line 10, it will not introduce a fill due to line 12, thus 5 will not be enqueued into $\text{newFrontierQueue}(\cdot)$. Path $8 \rightarrow 2 \rightarrow 5 \rightarrow 3$ stops at line 10 because the stored $\text{maxId}(3) = 2$ is smaller than the $\text{newMaxId} = 5$. Finally, path $8 \rightarrow 2 \rightarrow 3 \rightarrow 4$ qualifies line 10 since the stored $\text{maxId}(4) = 7$ is greater than $\text{newMaxId} = 3$, as well as introduces a new fill due to line 12 is true. Afterwards, only 4 is in the $\text{newFrontierQueue}(\cdot)$ at iteration 3 which will not introduce

any new frontiers.

It is important to note that the parallel traversal to all the neighbors of a frontier is not possible by simply allowing to process multiple threshold vertices in parallel, as mentioned in Challenge #1 (lines 9 - 10 of fill2 algorithm). Now, we discuss how the detection of the fill-in (8,4) from Figure 7 may not occur when we allow multiple threads to traverse the neighbors $\{1, 2, 7\}$ in parallel in fill2. At line 6 in fill2 algorithm, the $\text{fill}(\cdot)$ entries of all the neighbors $\{1, 2, 7\}$ of source 8 are marked as 8. Assuming multiple threads can work in parallel at line 9, implying the thresholds will be 1, 2 and 7 from three different threads, respectively. In that case, all the neighbors 1, 2 and 7 will be inserted into the $\text{frontierQueue}(\cdot)$. When accessing the neighbor of the frontier at line 12, if a thread working on frontier 7 explores vertex 4 before any other threads, fill2 algorithm will update $\text{fill}(4) = 8$ without detecting a fill-in at $(8, 4)$ because the path $8 \rightarrow 7 \rightarrow 4$ does not satisfy Theorem 1. Later, the thread that could detect a fill at $(8, 4)$ because of path $8 \rightarrow 2 \rightarrow 3 \rightarrow 4$ fails to do so because this thread cannot enter the branch at line 13 since $\text{fill}(4)$ is already updated as 8. Hence, the fill-in $(8, 4)$ remains undetected if we seek fine-grained parallelism in straightforward fill2 algorithm at Figure 3(b).

Optimizing the initialization of $\text{maxId}(\cdot)$. Maintaining a $\text{maxId}(\cdot)$ array for every traversal would require an excessive amount of space that becomes impractical for large datasets considering limited memory space on modern GPUs. Hence reusing $\text{maxId}(\cdot)$ array for different sources becomes essential. But this reuse also comes with reinitialization overhead, that is, making $\text{maxId}(\cdot)$ to maxVal before traversal. And this overhead is nontrivial, e.g., it takes 22% of the total time for PR dataset to re-initialize the $\text{maxId}(\cdot)$.

To reduce the reinitialization overhead, we propose to divide the value range $[0, \text{maxVal}]$ of $\text{maxId}(\cdot)$ into smaller ranges, i.e., $[0, |V|]$ range for each source. Note that maxVal is 2^{32} for a 32-bit integer type. During traversal, different sources can work on their respective value range of $\text{maxId}(\cdot)$ without reinitialization. For instance, for the first source, the traversal updates the maxId of the vertices in the range of $[\text{maxVal} - |V|, \text{maxVal}]$. Moving to the next source, we regard the range of $[\text{maxVal} - 2 \cdot |V|, \text{maxVal} - |V|]$ as valid for maxId . In this context, any maxId value greater than the upper bound, which is $\text{maxVal} - |V|$ in this case, is treated as initialized. Note, the value of maxId will not be smaller than $\text{maxVal} - |V|$ before execution. This optimization helps skip the $\text{maxId}(\cdot)$ initialization for a total of $\frac{\text{maxVal}}{|V|}$ sources. For instance, it helps reduce the ratio of reinitialization over the total time consumption from 22% to 0.082% for PR dataset.

Optimizing the access to $\text{maxId}(\cdot)$ and $\text{fill}(\cdot)$. It is important to mention that the accesses to both $\text{maxId}(\cdot)$ and $\text{fill}(\cdot)$ array are random thus time consuming. One can either access $\text{maxId}(\cdot)$ first to reduce the follow-up access to $\text{fill}(\cdot)$ or vice versa (i.e., adding line 9.5 to the pseudocode). In both cases, we avoid repeated frontier enqueueing for vertices that are already detected as fill-ins at line 14.

Putting $\text{maxId}(\cdot)$ access before $\text{fill}(\cdot)$, which is the pseudocode in Figure 6, will avoid the access to $\text{fill}(\cdot)$ array when the new path fails to update the maxId of existing paths, i.e., line 10 is evaluated as false. Consequently, this path avoids the access to follow-up $\text{fill}(\cdot)$ at line 12. Path $8 \rightarrow 7 \rightarrow 4$ from Figure 7(a) falls in this case.

Adding line 9.5 to the pseudo-code in Figure 6 will avoid unnecessarily lowering the `maxId` of an already detected fill-in. Note, we do not need to do so because this vertex will always propose its own vertex ID as the `maxId(:)` for the paths that come across this vertex. Using Figure 7(a) as an example, this logic avoids lowering `maxId[5]` from 2 to 1 when the path $8 \rightarrow 1 \rightarrow 0 \rightarrow 5$ attempts to do so because (8, 5) is already a fill due to $8 \rightarrow 2 \rightarrow 5$. And the continuation of the paths from 8 through 5 will surely use 5 as the `maxId`.

Including line 9.5 or not is a graph dependent option. Particularly, adding line 9.5 is beneficial for graphs that have relatively larger number of fill-ins. But for graphs with relatively smaller number of fill-ins, the condition of line 9.5 will become false for most of the time, resulting in higher overhead than benefits.

Multi-source symbolic factorization with combined frontierQueue(:) executes multiple sources of the Algorithm in Figure 6 concurrently on a single GPU in order to saturate the computing resources. To balance the workload across threads, we combine the frontiers of various concurrent traversals into a single `frontierQueue(:)`, which is analogous to recent multi-source graph traversal projects [11]. We further make three revisions to this design. First, we need to use `maxId(:)` instead of a single bit to track the status of each vertex. Furthermore, GSOFA relies upon `tracker(:)` and `newtracker(:)` to identify the source of each frontier. Last but not least, we directly use atomic operations to enqueue frontiers into the combined frontier queue, which is faster than prefix-sum based approach, according to [24].

5 PARALLEL EFFICIENT SUPERNODE DETECTION

This section further introduces a *massively parallel supernode detection algorithm* which can *balance the workload, minimize the communication* and *saturate the GPU resources* with judicious source scheduling and unified memory assistance.

GPU-friendly supernode detection with shared memory optimization. We introduce a two-phase supernode detection design to match the SIMT nature of GPUs. First, we use one GPU thread to examine whether the nonzero count of a row in \mathbf{U} is one fewer than the previous row (i.e., requirement i), and use a bitmap in shared memory to keep track of this information. If requirement (i) is not met, it means the current row starts a new supernode. We will use a queue in shared memory to keep track of this leading row. Second, each GPU thread dequeues a leading row of a supernode, examines the bitmap, also assesses the requirement (ii) of a supernode in Definition 1. Only when both requirements are met, the current row is considered to be part of this supernode. *The novelty of this two-phase design is that the first phase helps break the chunkSize of rows into independent subranges so that our second phase can work on these subranges in parallel, which fits the SIMT nature of GPUs.* In addition, the first phase is also massively parallel.

Figure 8 demonstrates how to, in parallel, detect supernodes for row ranges 0 - 3 of matrix \mathbf{U} . In phase I, all threads in parallel check whether the `nnz` of current row is 1 fewer than the prior row. Since only row 1 satisfies the condition, the bitmap is set to "0100" (①). Further, the leading rows of supernode are {0, 2, 3} in queue (②). In phase II, each thread is assigned to one leading row entry,

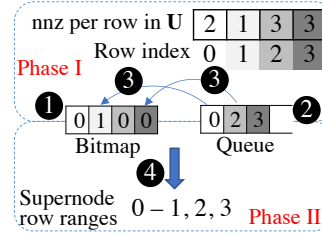


Figure 8: Parallel supernode detection for matrix \mathbf{U} in Figure 2(c), where Phase I confirms $nnz(\mathbf{U}(s, :)) = nnz(\mathbf{U}(s - 1, :)) - 1$ while Phase II checks $\mathbf{L}(s, r) \neq 0$.

and checks whether the bit following the leading row is 1 (③), as well as the corresponding entry in $\mathbf{L}(s, r) \neq 0$ (④). If both conditions are met, the supernode grows. We will continue this process until no supernode grows. Since row 1 satisfies both conditions while row 2 does not, supernode of row 0 grows to include only row 1.

Communication- and saturation- aware inter-node source scheduling. Since inter-node communication is significantly more expensive than intra-node counterpart, we restrict the supernode detection inside of each computing node. We use the following equation to perform inter-node source scheduling:

$$\begin{aligned} chunkSize \times numConcurrChunksPerNode \\ = \#C \times numGPUsPerNode, \end{aligned} \quad (1)$$

where *chunkSize* is identical to the size of the user defined maximum supernode. Note, we want the size of a chunk to be neither a fraction nor multiple of the maximum supernode size. On the one hand, since we interleave consecutive chunks to different nodes, making *chunkSize* a fraction of the maximum supernode size will introduce inter-node communication during supernode detection. On the other hand, making *chunkSize* a multiple of supernode size will make chunk size too large, potentially lead to worse workload imbalance. Further, $\#C$ is the preferred number of sources to saturate one GPU. In short, all the three items in Equation (1), i.e., *chunkSize*, $\#C$, and *numGPUsPerNode* are known, one can derive the value for *numConcurrChunksPerNode*. For example, the PR dataset on one Summit node would need 6,144 sources to saturate the GPUs. If the defined supernode size is 128, *numConcurrChunksPerNode* would be 48.

Unified memory assisted intra-node source scheduling. Once the number of concurrent chunks for a computing node is decided, we assign the sources of these chunks to various GPUs of this node in an interleaved fashion in order to balance the workload. However, since supernode detection requires the nonzero count of consecutive rows, communicating the fill information between GPUs is needed. Specifically, for a supernode spanning from row r to $s - 1$, one needs to communicate the *nonzero count* of row s in \mathbf{U} , and whether there is a *nonzero* at $\mathbf{L}(s, r)$. Since both data sizes to communicate are small and we cannot predict which nonzero count and fill location are required beforehand, traditional *cudaMemcpy* based data transfer is not suitable for such a kind of data sharing.

To this end, we choose the unified memory option over explicit data transfer. In this context, GSOFA uses a single CPU process to manage all the GPUs in each node, where the kernel launches are performed asynchronously. It is also important to note that unified memory introduces several

configurations which concern the performance.

First, we choose remote access over page migration when accessing the unified memory. Since we interleave the sources assignment in a fine-grained manner, various GPUs might compete to migrate the same page that stores the fill(:) or nonzero count information from the source GPU. This will hurt the performance. During implementation, our key guideline is *keeping this data in the GPU that modifies this data*. Particularly, only one GPU will detect the fill information for a specific row (i.e., modify that fill(:)), so does that for the supernode detection. This implies that only one GPU needs to keep the fill(:) information of that row locally while the remaining GPUs will access that information remotely. In implementation, we use the `cudaMemAdviseSetPreferredLocation` flag to advise the unified memory driver to keep that data in the same GPU which modifies this data [25]. For the remaining GPUs, we use `cudaMemAdviseSetAccessedBy` flag to instruct them to access that data remotely.

Second, GSOFA allocates separate fill(:) memory for different GPUs instead of using one virtual space spanning across all GPUs in one node. This design is inspired by a key observation that there exist a significant number of page faults when all the GPUs share one fill(:) address. This is caused by the fact that this single memory space is not perfectly aligned from one GPU to another. That is, one page could span across two GPUs. In this case, once both GPUs need to modify that page, these GPUs will compete for that shared page for writing. This will lead to frequent page migration.

6 SPACE OPTIMIZATION

Continuing our discussion in Challenge #3 of Section 3, this section will rigorously quantify the *space crisis faced by multi-source concurrent symbolic factorization*. Table 2 presents the space complexity of the six major data structures used by GSOFA, namely, `frontierQueue(:)` & `newFrontierQueue(:)`, `tracker(:)` & `newTracker(:)`, `maxId(:)`, and `fill(:)`. When the number of concurrent sources is $\#C$, the total space consumption would be around $6 \cdot |V| \cdot \#C$ entries. Apparently, space consumption immediately becomes a key problem for large graphs with relatively big number of concurrent sources. However, when it comes to performance, GSOFA prefers a larger $\#C$ which will provide more workload to better saturate the GPU computing resources.

Data structure	Space complexity
<code>frontierQueue(:)</code> & <code>newFrontierQueue(:)</code>	$2 \cdot V \cdot \#C$
<code>tracker(:)</code> & <code>newTracker(:)</code>	$2 \cdot V \cdot \#C$
<code>maxId(:)</code>	$ V \cdot \#C$
<code>fill(:)</code>	$ V \cdot \#C$

Table 2: The data structures used by multi-source concurrent GSOFA, where $|V|$ and $\#C$ are the number of vertices in the graph, and concurrent sources, respectively.

Dataset	Average usage (%)	Peak usage (%)
AU	0.12	8.2
DI	0.04	4.55
HM	0.01	1.7
PR	0.11	25.0
ST	0.1	1.0
TT	0.1	11.0

Table 3: Percentage of usage of allocated frontier queues.

In this section, we propose three interesting optimizations based upon the access pattern and usage of various

data structures to combat the high space complexity. Particularly, we introduce external GPU frontier management, bubble removal in `maxId(:)`, dynamic space allocation to dynamically assign memory space to various data structures. Note, this dynamic space allocation also allows GSOFA to support configurable space consumption.

External GPU frontier management is motivated by the observation in Table 3 where the average usage of frontier-related data structures remains low for the large datasets, i.e., AU, DI, HM, ST and TT. However, the peak usage can rise to as high as 25% for PR. This observation implies that we can allocate relatively smaller memory space to hold frontier-related data structures because the usage remains low for vast majority of the iterations. Once the usage goes beyond the allocated space, we resort to our external-GPU option which is presented below.

We propose to only allocate a fraction of the required space on GPU for the four frontier-related data structures, i.e., `frontierQueue(:)`, `newFrontierQueue(:)`, `tracker(:)` and `newTracker(:)` and write the extra frontiers out of GPU. Then, we load these external GPU data structures in GPU for computation. Figure 9 demonstrates this design. For brevity, Figure 9 only uses a single thread to traverse the graph for source 8 without loss of generality. The size of the allocated `newFrontierQueue(:)` is only three. At iteration 1, frontiers 1 and 2 exhaust the `newFrontierQueue(:)` by their neighbors 0, 5 and 3. Subsequently, we copy these three frontiers to CPU memory to still have available space for the incoming frontiers, 7 in this case. Proceeding to the next iteration that comes after swapping the queues, we will finish frontier 7 in the `frontierQueue(:)` first and load the frontiers from CPU in GPU for further computation. It is worthy of mentioning that, instead of directly storing the source ID in `tracker(:)` and `newTracker(:)` in multi-source concurrent traversal, we propose to store the index of the source for each frontier to reduce the space further.

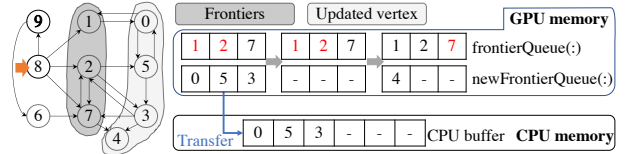


Figure 9: External GPU frontier management. Note, the traversal iteration is same as in Figure 7(a).

Bubble removal in `maxId(:)` is supported by the key observation that a source vertex is not allowed to traverse vertices that are larger than the source. Consequently, we can remove the “bubbles” in the `maxId(:)` that are larger than the source. Here, “bubble” means the allocated space that is not used in `maxId(:)`. Assuming the source vertex is v , according to Theorem 1, we will never update the `maxId` of the vertices that are larger than v .

Dynamic space allocation across data structures is motivated by two facts. First, `maxId(:)` and `fill(:)` accesses are more random than frontier-related data structure. Particularly, the accesses to `maxId(:)` and `fill(:)` are determined by the frontier’s neighbors which often have random vertex IDs. Therefore, GSOFA needs to put the entire `maxId(:)` and `fill(:)` arrays in GPU memory in order to achieve desirable performance. Second, the space requirement `maxId(:)` is

dynamic, that is, smaller sources need smaller `maxId(:)` and vice versa. Therefore, we can dynamically adjust the `maxId(:)` space so that the frontier-related data structures can have more GPU resources when possible.

Towards this end, we first allocate a big chunk of memory, in contrast to allocate separate memory spaces for various data structures discussed in Table 2. Subsequently, this memory chunk is dynamically divided among the data structures with priority given to `maxId(:)` and `fill(:)`. For a given number of concurrent sources in a traversal, after the space reduction strategy of `maxId(:)`, the amount of memory required for `maxId(:)` remains low for smaller sources. In this scenario, we can allocate more space for other data structures. Once we start working on larger sources, the `maxId(:)` space requirements begin to climb. In this context, we prioritize the space requirement for `maxId(:)` along with `fill(:)` so that a large number of concurrent sources can execute together with better, at least sustained, performance.

GSOFA with configurable space budget. After putting all the major data structures of GSOFA into one continuous memory space, we enable a new feature, i.e., configurable memory budget for GSOFA. This optimization will first conduct bubble removal, dynamic space allocation and external GPU frontier-related data structure management. If GSOFA still suffers from space shortage, GSOFA will judiciously reduce the number of concurrent sources in order to restrict the memory space consumption in the given budget.

7 EVALUATION

We implement GSOFA with $\sim 2,500$ lines of C++/CUDA code, and compile the source code with NVIDIA CUDA 10.1 Toolkit with the optimization flag set to be O3. We use IBM Spectrum MPI 10.3.0.0 for inter-node communication. The evaluation platform is the Summit supercomputer at Oak Ridge National Laboratory [26], where each computing node is equipped with dual-socket IBM POWER9 CPU processors (i.e., 42 cores), and six NVIDIA V100 GPUs. All the GPUs on one Summit node are connected with NVIDIA’s high-speed NVLink. We use Traversed Edges Per Second (TEPS) to report the graph traversal performance, and take the average of three runs. We use NVIDIA nvprof profiling tool for Figures 13 and 18.

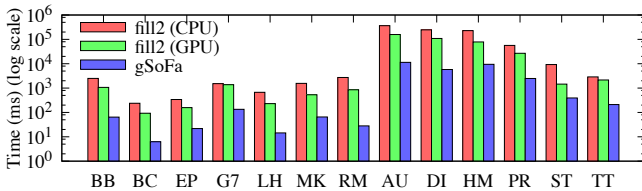


Figure 10: Performance of the parallel CPU and GPU versions of `fill2` and GSOFA on a Summit compute node.

Figure 10 compares the original `fill2` algorithms on CPU and GPU, and GSOFA on one Summit node. Both CPU and GPU versions of `fill2` that we are using are the straightforward parallel implementation of the `fill2` algorithm in Figure 3(b). The parallel CPU version employs all the 42 cores in a Summit node to perform symbolic factorization for 42 sources in parallel. The GPU version allows every allocated thread to work in a source as long as the GPU memory is sufficient. On average, GSOFA is $13.01\times$ faster than the

GPU-based parallel `fill2` algorithm, with the maximum and minimum speedups $30.41\times$ (RM) and $3.70\times$ (ST), respectively. When GSOFA is compared to CPU version, GSOFA enjoys even higher speedups, that is, $33.06\times$, on average, with $97.11\times$ (RM) and $11.32\times$ (G7) as the maximum and minimum speedups. This suggests that the original `fill2` algorithm is not suitable for massively parallel GPUs.

Comparison with the state-of-the-art CPU algorithms.

We compare GSOFA with the state-of-the-art sequential symbolic factorization in GLU3.0 [27] and the parallel symbolic factorization from SuperLU_DIST [2]. Note, both the GLU3.0 and SuperLU_DIST may use GPUs only during numerical factorization or later phases. However, both libraries use CPUs for symbolic factorization.

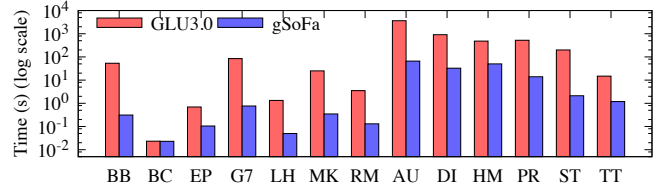


Figure 11: Performance comparison of GSOFA with the state-of-the-art CPU symbolic factorization in GLU3.0. For fair comparison, GSOFA uses one V100 GPU since GLU3.0 performs symbolic factorization sequentially.

Figure 11 demonstrates the time comparison of GSOFA with the symbolic factorization on GLU3.0 [27]. We limit the GSOFA to a single GPU because GLU3.0 can only perform single-threaded symbolic factorization on CPU. Note, GLU3.0 is based upon Gilbert-Peierls algorithm [8] which suffers from stringent data dependency problems if intended to implement in parallel (discussed in Section 2). In the figure, we can observe a maximum and minimum speedup of $171.1\times$ (BB) and $1.01\times$ (BC). The workload in BC dataset is too small to saturate the GPUs, i.e., see Figure 4(a). On average, the GSOFA is $50.1\times$ faster than symbolic factorization in GLU3.0.

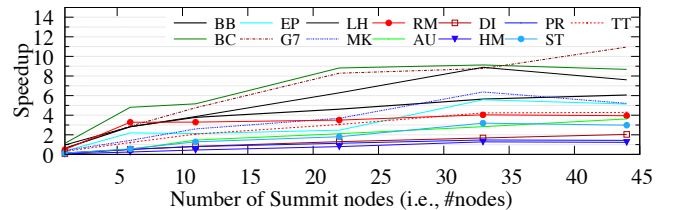


Figure 12: Speedup of GSOFA over the state-of-the-art CPU parallel symbolic algorithm in SuperLU_DIST.

Figure 12 compares the performance between GSOFA and the CPU parallel symbolic factorization in SuperLU_DIST. GSOFA starts worse than the CPU parallel algorithm for majority of the datasets. However, with more and more Summit nodes, GSOFA begins to outperform the CPU algorithm. Particularly, initially on one node, GSOFA is $1.2\times$ faster on BC and $1.6\times$, $2.8\times$, $1.6\times$, $1.03\times$, $2.6\times$, $2.0\times$, $13.8\times$, $11.3\times$, $19.3\times$, $6.7\times$, $7.4\times$ and $3.2\times$ slower on BB, EP, G7, LH, MK, RM, AU, DI, HM, PR, ST and TT, respectively on one node. When it goes to six Summit nodes, GSOFA bests CPU parallel symbolic factorization on majority of the datasets, i.e., BC, EP, BB, RM, G7, MK, LH and TT. GSOFA finally outperforms CPU parallel symbolic

factorization across all the datasets by $5\times$, on average, with 44 Summit nodes, at which the maximum speedup of $10.9\times$ is achieved for G7 matrix and minimum of $1.3\times$ is achieved for HM matrix. For G7, the speedup of GSOFA is high because G7 is highly non-symmetric which limits the benefits of symmetric pruning in SuperLU_DIST and its larger workload (i.e., $\frac{\text{nnz}(A)}{\text{Order}(A)}$ in Table 1). HM though has relatively low symmetry but it presents the lowest workload which causes the speedup to be poor.

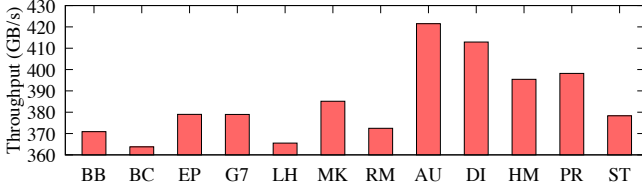


Figure 13: Throughput achieved by GSOFA on a V100 GPU.

Throughput. Figure 13 demonstrates the throughput achieved by the GSOFA on V100 GPU. Particularly, GSOFA achieves a maximum of 421.5 GB/s (47% of peak memory throughput) for AU dataset and a minimum of 363.8 GB/s (40% of peak memory throughput). On average, GSOFA achieves a throughput of 385.2 GB/s (43% of peak memory throughput) over all the dataset, which is rarely observed for graph traversal related applications [28].

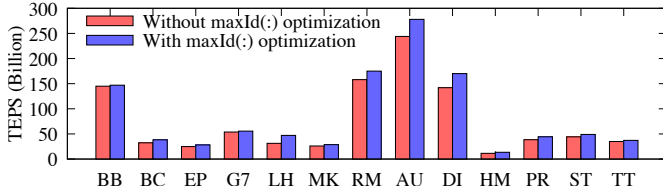


Figure 14: Impacts of maxId() initialization optimization.

Performance impact of maxId() optimization. This reduces the initialization overhead for maxId(). We observe noticeable performance gains across all datasets. On average, we achieve 14% speedup with this optimization, where the maximum impact comes from LH of 50% and a minimum of 2% increment in BB dataset.

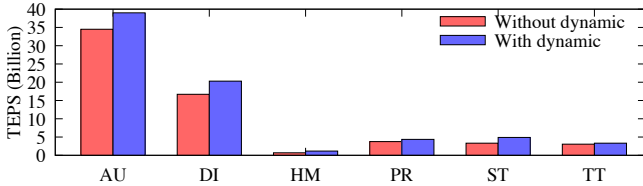


Figure 15: Performance impact of dynamic space allocation on one GPU with 1 GB space allocation.

Performance impact of dynamic space allocation. Figure 15 presents the effect of dynamic memory allocation. We allocate 1 GB of memory and study the performance of *with* versus *without* the dynamic memory optimization. As expected, this optimization yields performance gains for all the large datasets in Table 1. Particularly, we observe improvements of $1.2\times$, $1.2\times$, $1.7\times$, $1.2\times$, $1.5\times$ and $1.1\times$ on AU, DI, HM, PR, ST and TT respectively.

GSOFA space optimization with limited memory budget. This includes all the space optimizations, further along with an explicit restriction on the available memory space.

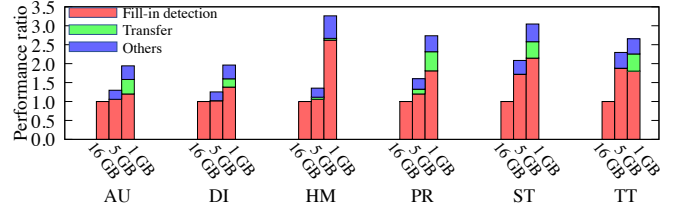


Figure 16: Performance study of GSOFA with different memory capacities on 44 Summit nodes. The allowed GPU memory for all the data structures in GSOFA varies from 1 to 16 GB.

Particularly, we enable external GPU frontierQueue(), newFrontierQueue(), tracker() and newtracker() management, bubble removal and dynamic allocation. For the single large array that is shared by all the data structures, we limit its size to be 16, 5 and 1 GBs to demonstrate the performance robustness of GSOFA. This optimization will involve transferring data between CPU and GPU memories and other overheads, such as, checking the condition of memory overflow.

Figure 16 shows the trade-off between the runtime and the space consumption. The general trend is that the performance drops with the decrease of allocated space. Particularly, with merely 1 GB memory budget, the GSOFA performance decreases by $1.9\times$, $1.9\times$, $3.3\times$, $2.7\times$, $3.0\times$ and $2.6\times$, respectively, on AU, DI, HM, PR, ST and TT datasets. The performance drops in fill-in detection are caused by the fact that a limited memory budget leads to the reduction of #C. Note, enabling external GPU GSOFA also results in the overhead of checking whether the frontier queue overflows, which is denoted as “GSOFA: Others” in Figure 16.

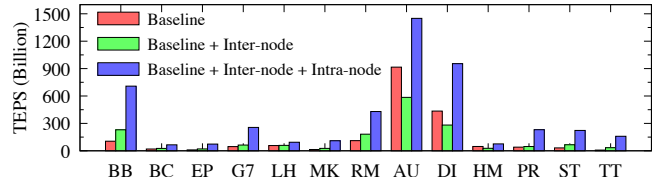


Figure 17: Performance impacts of supernode detection optimizations on six nodes (i.e., 36 GPUs). Here, the inter-node optimization interleaves the chunk of rows among the compute nodes and the intra-node optimization interleaves rows scheduling in fine-grained manner during fill-in detection to obtain optimum intra-node workload balancing.

Performance impact of supernode detection optimizations. In Figure 17 we use six nodes (i.e., 36 GPUs) in order to showcase the impacts of our inter-node workload balancing strategies. Specifically, the “baseline” version assigns a block of continuous chunks of sources to a node. The “inter-node” interleaves the chunks of sources across different nodes in a round-robin approach. And “intra-node” performs unified memory-assisted fine-grained source scheduling across GPUs in a node. On average, the inter-node scheduling optimization yields $1.6\times$ speedup over the baseline. The intra-node optimization further adds another $3.3\times$ speedup. The maximum gains of inter- and intra- node optimizations are $4.3\times$ (TT) and $4.9\times$ (PR), respectively. We also observe performance drops for the inter-node optimization on AU, DI and HM graphs because the new combination of chunks of sources might have relatively more non-uniform workload distribution.

Performance impacts of unified memory optimizations.

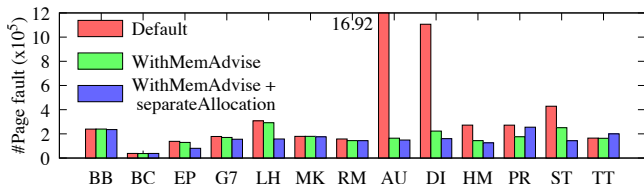


Figure 18: Performance impacts of unified memory optimization on one Summit node.

Figure 18 demonstrates the effects of `cudaMemAdvise` and data structure separation optimizations. We can observe a general trend of reduction in the number of page faults with each optimization’s addition. Particularly, `cudaMemAdvise` reduces the page fault by 2.2 \times . In particular, we see the maximum drop of page faults in AU by a factor of 10.3 \times . The data structure separation optimization, when added to the `cudaMemAdvise` version, further reduces the page faults by an average of 20%, with the maximum drop to be 85% in LH. Note, PR and TT experience more page faults due to the data structure separation optimization. This is potentially caused by the fact that the change of memory alignment could result in more page faults for the memory that is not at the boundary.

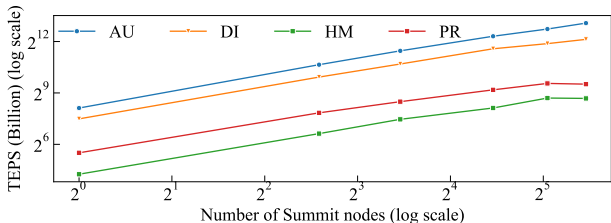


Figure 19: Scaling GSOFA to 44 Summit nodes (264 GPUs).

Strong scalability. Figure 19 demonstrates strong scaling of GSOFA up to 44 Summit nodes (264 GPUs) for relatively large datasets, i.e., AU, DI, PR, and HM. Particularly, GSOFA achieves speedups of 31.0 \times , 24.9 \times , 21.5 \times , and 16.1 \times , respectively, for AU, DI, PR, and HM. It is worth noting that GSOFA can effectively use a number of GPUs that is not necessarily a power-of-two, which provides great flexibility to the application code. We also notice that the HM and PR enjoy close to linear scalability from 1 to 33 nodes but flat out from 33 to 44 nodes because these two datasets do not have the adequate workload to saturate 44 nodes.

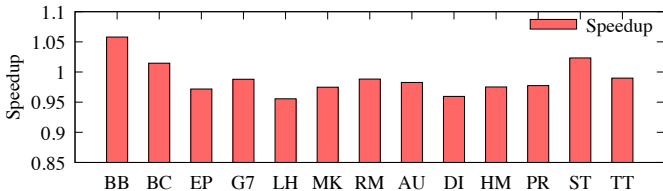


Figure 20: Impacts of using line 9.5 in algorithm in Figure 6.

Performance impacts of optimizing access to `maxId(:)` and `fill(:)`. Putting `fill(:)` ahead of `maxId(:)` introduces observable performance differences. Particularly, as shown in Figure 20 in a GPU, enabling line 9.5 in Figure 6 leads to performance gain in BB, BC and ST datasets with the maximum gain of 5.5% on BB dataset. This performance gain indicates that noticeable fill-ins are repeatedly detected. Likewise, among remaining datasets, LH dataset experiences a maximum performance drop of 4.6% because the number of `maxId(:)` update after fill detection of a vertex

is small. Given the majority of the datasets experience performance drop with the enabling of line 9.5, we disable line 9.5 for all the datasets in our evaluations.

GSOFA performance variation with `chunkSize`. Figure 21 presents the performance impact of `chunkSize` from Equation 1 to GSOFA. In general, one can observe that the increase of the `chunkSize` leads to longer time. On average, the performance degradation is 1.04 \times and 1.13 \times when we increase `chunkSize` from 64 to 128 and 256, respectively, with the maximum slowdown as 1.6 \times of MK for `chunkSize` = 64 to 256. The reason behind this trend is that the increase of the `chunkSize` typically leads to more imbalanced workload. In this work, we follow SuperLU_DIST to set `chunkSize` = 128 for GSOFA as the default configuration even though `chunkSize` = 64 is slightly faster.

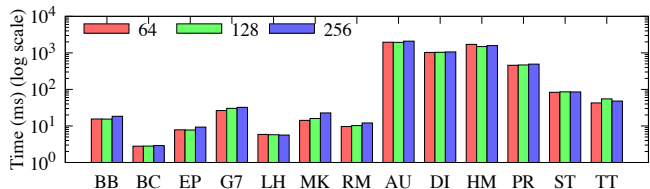


Figure 21: GSOFA performance variation with respect to `chunkSize` = 64, 128 and 256 in six Summit compute nodes.

8 RELATED WORK

Most of the major state-of-the-art sparse LU factorization codes, such as GLU3.0 [27], [29], [30] and SuperLU_DIST [20] adopt the Gilbert-Peierls algorithm [8] for symbolic factorization. GLU3.0, [29] and [30] directly use the sequential version of the Gilbert-Peierls algorithm [8]. SuperLU_DIST implements a parallel CPU-based symbolic factorization algorithm [2] based upon Gilbert-Peierls algorithm with symmetric pruning [31] and supernodal traversal [10]. To improve parallelism for symbolic factorization, SuperLU_DIST resorts to nested dissection to partition **A** into independent rows at a level of the separator tree [2]. The filled structures computed at a level of the separator tree are communicated among the required processes that perform symbolic factorization at the higher level of the separator tree. This leads to inter-process communications while the computation changes the level along the separator tree. To the best of our knowledge, GSOFA is the first GPU-based parallel symbolic factorization algorithm.

It is also worth mentioning that NVIDIA cuSOLVER library [32] provides solvers for sparse and dense linear systems with different approaches including LU decomposition. However, cuSOLVER does not yet support GPU version of sparse LU decomposition. Hence, we cannot compare GSOFA against the symbolic factorization of cuSolver.

Fill2 algorithm is similar to Dijkstra’s Single Source Shortest Path algorithm [33] in the sense that fill2 uses the maximum vertex ID on a path to represent the “distance” metric in Dijkstra’s algorithm. The salient difference between these two algorithms lies in that one does not need to reduce the “distance” further if a fill-in is already detected in symbolic factorization. We also would like to point out that the previously proposed Δ -step optimization [34] for Dijkstra’s algorithm is not effective here: because fill2 only allows vertices that are smaller than the source to be active, Δ -step will further restrict the parallelism. However, the

similarity between these two algorithms suggests that our design of fine-grained symbolic factorization, external GPU optimizations, and supernode detection can potentially provide performance enhancement and space saving technique to the multi-source Dijkstra’s algorithm [35].

9 CONCLUSION AND FUTURE WORK

This paper introduces GSOFA the first, to the best of our knowledge, GPU-based sparse LU symbolic factorization system. In particular, we revamp the fill2 algorithm to enable fine-grained symbolic factorization, redesign supernode detection to expose massive parallelism, balanced workload and minimal inter-node communication, and introduce a three-pronged space optimization to handle large sparse matrices. Taken together, we scale GSOFA up to 264 GPUs with unprecedented performance. As the future work, we plan to integrate GSOFA into the state-of-the-art SuperLU_DIST package.

10 ACKNOWLEDGEMENT

This work was in part supported by NSF CRII Award No. 2000722, CAREER Award No. 2046102 and the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration (NNSA). Any opinions, findings and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of DOE, NNSA or NSF.

REFERENCES

- [1] Timothy A Davis, et al. A Survey of Direct Methods for Sparse Linear Systems. *Acta Numerica*, 2016.
- [2] Laura Grigori, et al. Parallel Symbolic Factorization for Sparse LU with Static Pivoting. *SIAM Journal on Scientific Computing*, 2007.
- [3] E Lezar et al. GPU-based LU Decomposition for Large Method of Moments Problems. *Electronics letters*, 2010.
- [4] Piyush Sao, et al. A Distributed CPU-GPU Sparse Direct Solver. In *Euro-Par*. Springer, 2014.
- [5] Weifeng Liu, et al. A Synchronization-free Algorithm for Parallel Sparse Triangular Solves. In *Euro-Par*. Springer, 2016.
- [6] Donald J Rose et al. Algorithmic Aspects of Vertex Elimination on Directed Graphs. *SIAP*, 1978.
- [7] John R Gilbert et al. Elimination Structures for Unsymmetric Sparse LU Factors. *SIMAX*, 1993.
- [8] John R Gilbert et al. Sparse Partial Pivoting in Time Proportional to Arithmetic Operations. *SISC*, 1988.
- [9] Edsger W Dijkstra. A Note on Two Problems in Connexion with Graphs. *Numerische mathematik*, 1959.
- [10] James W Demmel, et al. A Supernodal Approach to Sparse Partial Pivoting. *SIMAX*, 1999.
- [11] Hang Liu, et al. ibfs: Concurrent Breadth-first Search on GPUs. *SIGMOD*, 2016.
- [12] Iain S Duff et al. The Design and Use of Algorithms for Permuting Large Entries to the Diagonal of Sparse Matrices. *SIMAX*, 1999.
- [13] William F Tinney et al. Direct Solutions of Sparse Network Equations by Optimally Ordered Triangular Factorization. *Proceedings of the IEEE*, 1967.
- [14] Elizabeth Cuthill et al. Reducing the Bandwidth of Sparse Symmetric Matrices. In *24th National Conference*, 1969.
- [15] Brian W Kernighan et al. An Efficient Heuristic Procedure for Partitioning Graphs. *The Bell system technical journal*, 1970.
- [16] George Karypis et al. A Parallel Algorithm for Multilevel Graph Partitioning and Sparse Matrix Ordering. *Elsevier*, 1998.
- [17] Tesla NVIDIA. NVIDIA Tesla V100 GPU Architecture, 2017.
- [18] Lingda Li et al. Compiler Assisted Hybrid Implicit and Explicit GPU Memory Management Under Unified Address Space. In *SC*, 2019.
- [19] Suite Sparse Matrix Collection. Suite Sparse Matrix Collection. Retrieved from <https://sparse.tamu.edu/>, 2020. Accessed: 2019, November 15.

- [20] Xiaoye S Li et al. SuperLU_DIST: A Scalable Distributed-memory Sparse Direct Solver for Unsymmetric Linear Systems. *TOMS*, 2003.
- [21] Karypis Lab. ParMETIS. Retrieved from <http://glaros.dtc.umn.edu/gkhome/metis/parmetis/>. Accessed: 2020, June 02.
- [22] Hang Liu et al. Enterprise: Breadth-first Graph Traversal on GPUs. In *SC*, 2015.
- [23] Duane Merrill, et al. Scalable GPU Graph Traversal. *Acm Sigplan Notices*, 47(8):117–128, 2012.
- [24] Anil Gaihre, et al. XBFS: Exploring Runtime Optimizations for Breadth-first Search on GPUs. *HPDC*, 2019.
- [25] NVIDIA. Cuda Toolkit Documentation. Retrieved from <https://docs.nvidia.com/cuda/index.html>. Accessed: 2020, August 11.
- [26] Oak Ridge National Laboratory. Summit: America’s newest and smartest supercomputer. Retrieved from <https://www.olcf.ornl.gov/summit/>. Accessed: 2018, August 6.
- [27] Shaoyi Peng et al. GLU3.0: Fast GPU-based Parallel Sparse LU Factorization for Circuit Simulation. *IEEE Design & Test*, 2020.
- [28] Yangzihao Wang, et al. GUNROCK: A High-performance Graph Processing Library on the GPU. In *PPoPP*, 2016.
- [29] Ling Ren, et al. Sparse LU Factorization for Parallel Circuit Simulation on GPU. In *DAC*, 2012.
- [30] Kai He, et al. GPU-accelerated Parallel Sparse LU Factorization Method for Fast Circuit Analysis. *IEEE VLSI Systems*, 2015.
- [31] Stanley C Eisenstat et al. Exploiting Structural Symmetry in Unsymmetric Sparse Symbolic Factorization. *SIMAX*, 1992.
- [32] CUDA Toolkit Documentation. Available at <https://docs.nvidia.com/cuda/cusolver/>.
- [33] Edsger W Dijkstra et al. A Note on Two Problems in Connexion with Graphs. *Numerische mathematik*, 1959.
- [34] Ulrich Meyer et al. Δ -stepping: A Parallelizable Shortest Path Algorithm. *Journal of Algorithms*, 2003.
- [35] Piyush Sao, et al. A Supernodal All-pairs Shortest Path Algorithm. *PPoPP*, 2020.

Anil Gaihre is a Ph.D. candidate at the Department of Electrical and Computer Engineering, Stevens Institute of Technology. His research interests include graph theory, sparse linear algebra, high performance computing in multi-core computer architectures and blockchain. His prior work experience involves working as a lead Software Engineer at E&T Nepal Pvt. Ltd. that involved Research and Development on CFD simulations on GPUs.



Xiaoye Sherry Li is a Senior Scientist at Lawrence Berkeley National Laboratory. She has worked on diverse problems in high performance scientific computations, including parallel computing, sparse matrix computations, high precision arithmetic, and combinatorial scientific computing. She has (co)authored over 120 publications, and contributed to several book chapters. She is the lead developer of SuperLU, a widely-used sparse direct solver, and has contributed to the development of several other mathematical libraries, including ARPREC, LAPACK, PDSLin, STRUMPACK, and XBLAS. She is a SIAM Fellow and an ACM Senior Member.



Hang Liu is an Assistant Professor of Electrical and Computer Engineering at Stevens Institute of Technology. Prior to joining Stevens, he was an assistant professor at the Electrical and Computer Engineering Department of University of Massachusetts Lowell. He is on the editorial board for *Journal of BigData: Theory and Practice*, a program committee member for SC, HPDC, and IPDPS, and regular reviewer for TPDS and TC. He earned his Ph.D. degree from the George Washington University 2017. He is the Champion of the MIT/Amazon GraphChallenge 2018 and 2019, and one of the best papers awardee in VLDB '20.



11 EXAMPLES FOR LU DECOMPOSITION PHASES

Matrix preprocessing. Figure 22(a) demonstrates the matrix preprocessing on an example matrix. Particularly, we swap columns 2 and 7 so that column 2 will have 4 instead of 5 nonzeros. Further, we perform a swap between rows 2 and 3 for better numerical stability, i.e., larger numerical values are moved to diagonal.

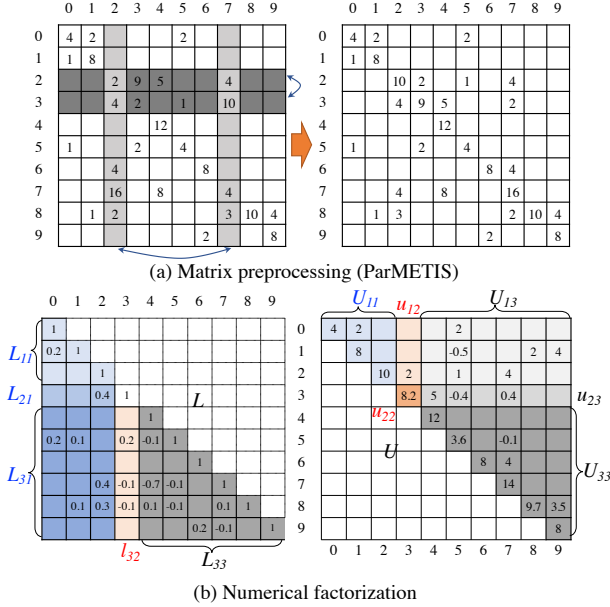


Figure 22: Other phases in LU decomposition.

Numerical factorization. Figure 22(b) shows a left-looking approach of numerical factorization to determine the k -th column of the filled matrix by using the computed results from columns 0 to $k-1$. Using column 3 of $(L+U)$ in Figure 22(b) as an example, we will use columns 0, 1 and 2 to solve this column. Representing A , L and U in block-matrix form, we can write LU factorization of A as follows.

$$\begin{bmatrix} L_{11} & & \\ l_{21} & 1 & \\ L_{31} & l_{32} & L_{33} \end{bmatrix} \begin{bmatrix} U_{11} & u_{12} & U_{13} \\ & u_{22} & u_{23} \\ & & U_{33} \end{bmatrix} = \begin{bmatrix} A_{11} & a_{12} & A_{13} \\ a_{21} & a_{22} & a_{23} \\ A_{31} & a_{32} & A_{33} \end{bmatrix}, \quad (2)$$

where $L_{11} = L(0 : 2, 0 : 2)$, $l_{21} = L(3 : 3, 0 : 2)$, $L_{31} = L(4 : 9, 0 : 2)$, $l_{32} = L(4 : 9, 3 : 3)$, $L_{33} = L(4 : 9, 4 : 9)$, $U_{11} = U(0 : 2, 0 : 2)$, $u_{12} = U(0 : 2, 3 : 3)$, $U_{13} = U(0 : 2, 3 : 9)$, $u_{22} = U(3 : 3, 3 : 3)$, $u_{23} = U(3 : 3, 4 : 9)$, and $U_{33} = U(4 : 9, 4 : 9)$. The texts in green, red and black colors represent the known, currently under solving, and unknown blocks, respectively. Through block matrix multiplication of Equation (2) towards a_{12} , a_{22} and a_{32} , we further obtain:

$$\begin{aligned} L_{11}u_{12} &= a_{12}, \\ l_{21}u_{12} + u_{22} &= a_{22}, \\ L_{31}u_{12} + l_{32}u_{22} &= a_{32}, \end{aligned} \quad (3)$$

where u_{12} , u_{22} and l_{32} are computed in order in Equation (3) so that the u_{22} and l_{32} can be calculated once u_{12} is resolved in the first equation of Equation (3).