# Parallel Symbolic Cholesky Factorization

Tobias Ribizel[1], Hartwig Anzt[1,2]

[1] Karlsruhe Institute of Technology (KIT), Karlsruhe, Germany
[2] Innovative Computing Laboratory (ICL), University of Tennessee, Knoxville, USA

# Sparse Direct Solvers

- Sparse linear system $Ax = b$, solved via factorization $A = L \cdot U$

- Factorization adds *fill-in* to sparsity pattern of $A$

- Fill-in reducing heuristics like Approximate Minimum Degree (AMD) and Nested Dissection reordering

- Reordered matrix $PAP^T$ has fewer fill-in entries → less storage, faster factorization

- Reordering → Symbolic Factorization → Numerical Factorization → Triangular Solve

# Sparse Direct Solvers

- Sparse linear system $Ax = b$, solved via factorization $A = L \cdot U$

- Factorization adds *fill-in* to sparsity pattern of $A$

- Fill-in reducing heuristics like Approximate Minimum Degree (AMD) and Nested Dissection reordering

- Reordered matrix $PAP^T$ has fewer fill-in entries $\rightarrow$ less storage, faster factorization

- Reordering $\rightarrow$ **Symbolic Factorization** $\rightarrow$ Numerical Factorization $\rightarrow$ Triangular Solve

# Sparse Factorizations: Symbolic Phase

- General criterion: For non pattern-symmetric $A$, $(L + U)_{ij}$ is nonzero if and only if there is a path $i \rightarrow k_1 \rightarrow \cdots \rightarrow k_n \rightarrow j, \ k_l < \min(i, j)$ through $A$

- Symmetric case: $L_{ij}$ nonzero if and only if there is a path $i \rightarrow k_1 \rightarrow \cdots \rightarrow k_n \rightarrow j, k_l < j$

- $L$ has a compact representation through the Elimination Tree $T$ (transitive reduction)

- $T$ can be computed in almost linear time (size and number of nonzeros of $A$)

- Sparsity pattern of row $L_{i*}$ consists of all pairwise paths between $A_{i*}$ through $T$

# Symbolic Cholesky: Enumerating Fill-in

**Theorem:** The entry $l_{ij}$ in $L$ is (symbolically) nonzero if and only if the row subtree of $i$ contains $j$, i.e. there is a nonzero $a_{ik}$ in $A$ such that $j$ lies on the path from $k$ to $i$ in the elimination tree $T$

$\rightarrow$ we need to identify lowest common ancestors (LCAs) between pairs $k_1, k_2$ belonging to nonzeros $a_{ik_1}, a_{ik_2}$

**Theorem:** In a post-ordered tree $T$, the LCA between any pair of nodes $u < v$ is the first node $w$ on the path from $u$ towards the root that fulfills $v \leq w$
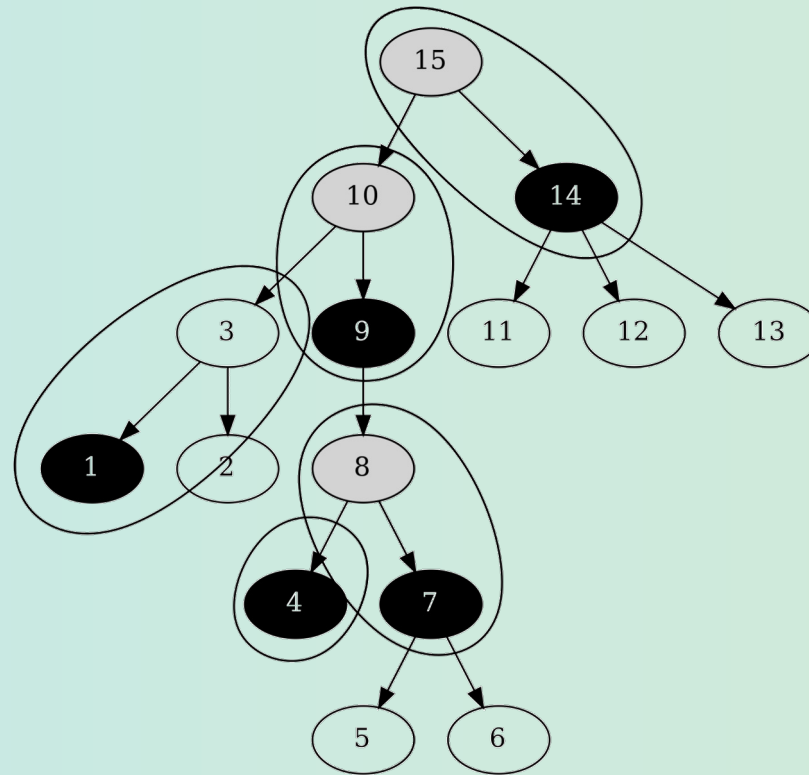
$\rightarrow$ we can limit ourselves to LCA search between consecutive nonzeros of $A$ in postorder

**Theorem:** After postordering the matrix, the ordered lower nonzeros of a row $a_{ij_k}$, $j_1 < \cdots < j_n = i$ give a path decomposition of the row subtree of row $i$ via $[j_k, LCA(j_k, j_{k+1}))$, plus the root $i$.

# Symbolic Cholesky: Enumerating Fill-in



Nonzeros in $A$

LCAs between nodes

# Symbolic Cholesky: Algorithmic Framework

- Copy $A$ to the CPU (optional)

- CPU: Compute Elimination Tree $T$

- CPU: Compute post-ordering of $T$

- Copy $T$ to the device (optional)

- Device: Reorder rows of $A$ with post-order column indices

- Device: Count number of nonzeros for each row

- Device: Allocate memory

- Device: Generate nonzeros for each row

- Device: Sort rows by column index (optional)

# Symbolic Cholesky: Computing the Elimination Tree

make_sets($n$);
**for** row $i = 0, \ldots, n$ **do**
  subtree_roots[$i$] = $i$;
  $i_{rep}$ = $i$; parent[$i$] = $n$;
  **for** column $j < i$ with $a_{ij} \neq 0$ **do**
    $j_{rep}$ = find($j$);
    $r$ = subtree_roots[$j_{rep}$];
    **if** parent[$r$] = $n$ and $r \neq i$ **then**
      parent[$r$] = $i$;
      $i_{rep}$ = union($i_{rep}, j_{rep}$);
      subtree_roots[$i_{rep}$] = $i$;
    **end**
  **end**
**end**
children = [$0, \ldots, n-1$];
child_ptrs, children =
  transpose(children, parents);

# Symbolic Cholesky: Enumerating

```
// Map column indices to postorder
```
**parfor** $j = 0, \ldots, m - 1$ **do**
$\quad$ $\text{post\_columns}[j] = \text{postorder}^{-1}[\text{columns}[j]]$;
**end**
```
// Sort postorder column indices
```
**parfor** $i = 0, \ldots, n - 1$ **do**
$\quad$ $\text{sort}(\text{post\_columns}[\text{row\_ptrs}[i], \ldots, \text{diag\_idx}[i]])$;
**end**
```
// Traverse row subtrees
```
**parfor** $i = 0, \ldots, n - 1$ **do**
$\quad$ **parfor** $j = \text{row\_ptrs}[i], \ldots, \text{diag\_idx}[i]$ **do**
$\quad\quad$ $u_{post} \leftarrow \text{post\_columns}[j]$;
$\quad\quad$ $v_{post} \leftarrow \text{post\_columns}[j + 1]$;
$\quad\quad$ **while** $u_{post} < v_{post}$ **do**
$\quad\quad\quad$ $u = \text{postorder}[u_{post}]$;
$\quad\quad\quad$ $u_{post} = \text{post\_parent}[u_{post}]$;
$\quad\quad\quad$ `// Output or count nonzero` $l_{iu}$
$\quad\quad$ **end**
$\quad$ **end**
**end**

# Performance Evaluation: Setup

- Code available in the Ginkgo HPC library 📚 https://github.com/ginkgo-project/ginkgo

- Comparison against symbolic part of CHOLMOD

- Compiled using gcc 11.3.0, CUDA 11.8, ROCM 5.1.1 with –O3 flags

- Inputs: (almost) all pattern-symmetric matrices from SuiteSparse with between $10^4$ and $10^7$ rows/columns

- Benchmarks with input ordering (*natural*), AMD and Nested Dissection (*nd*) ordering

- Removed all inputs that overflow 32 bit indices in row pointers

- Validated as correct for small sample of matrices

|                | natural | AMD   | ND    |
|----------------|---------|-------|-------|
| #matrices      | 458     | 579   | 601   |
| median fill-in | 131x    | 8.5x  | 7.7x  |

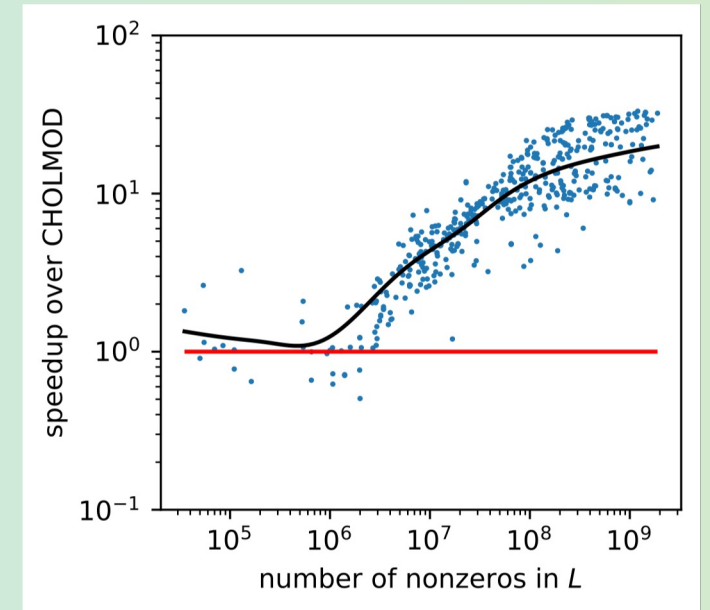|                 | NVIDIA             | AMD         |
|-----------------|--------------------|-------------|
|                 | Intel              | AMD         |
| CPU             | Xeon Platinum 8368 | EPYC 7543   |
| Sockets         | 2                  | 2           |
| Cores/Socket    | 38                 | 32          |
| L3 Cache/Socket | 57 MB              | 256 MB      |
|                 |                    |             |
| GPU             | NVIDIA A100        | AMD MI210   |
| VRAM            | 40 GB              | 64 GB       |
| Memory BW       | 1555 GB/s          | 1600 GB/s   |
| FP32 FLOPS      | 19.5 TFLOPS/s      | 22.6 TFLOP/s |

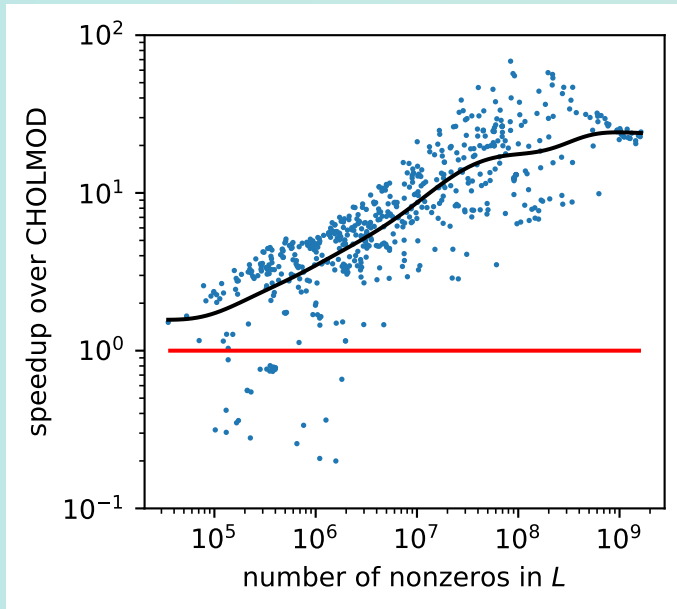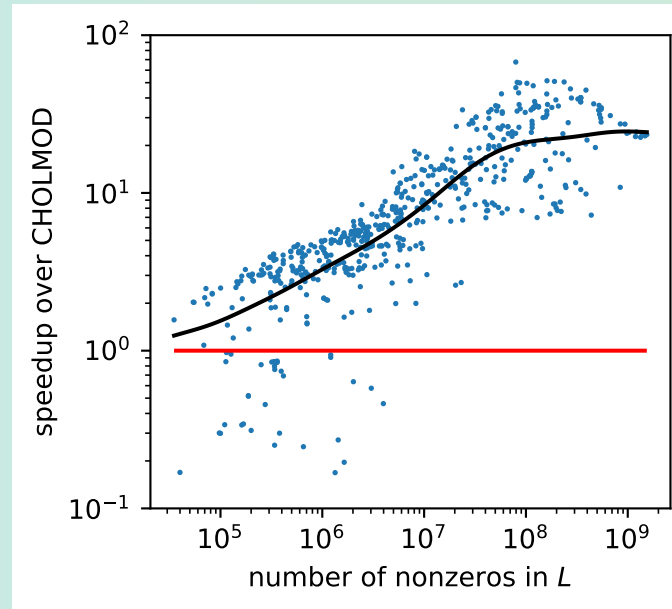# Performance Evaluation: Results Intel CPU (Speedup)



Nested Dissection

AMD

Natural
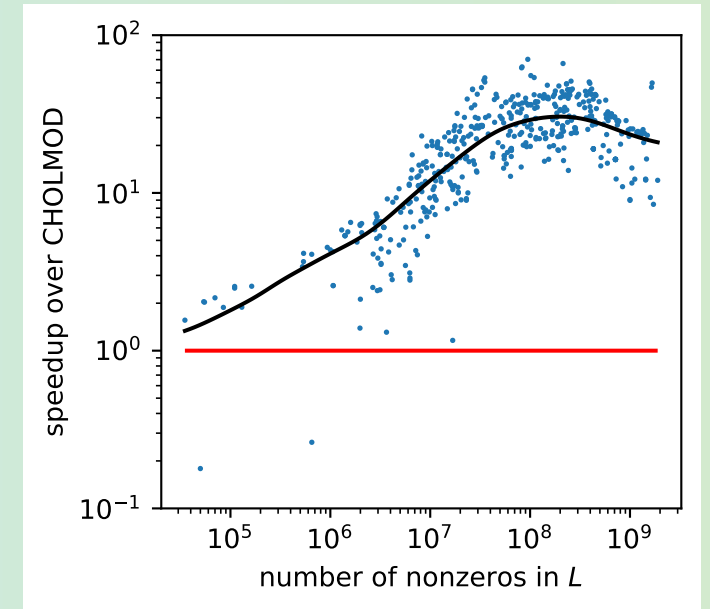
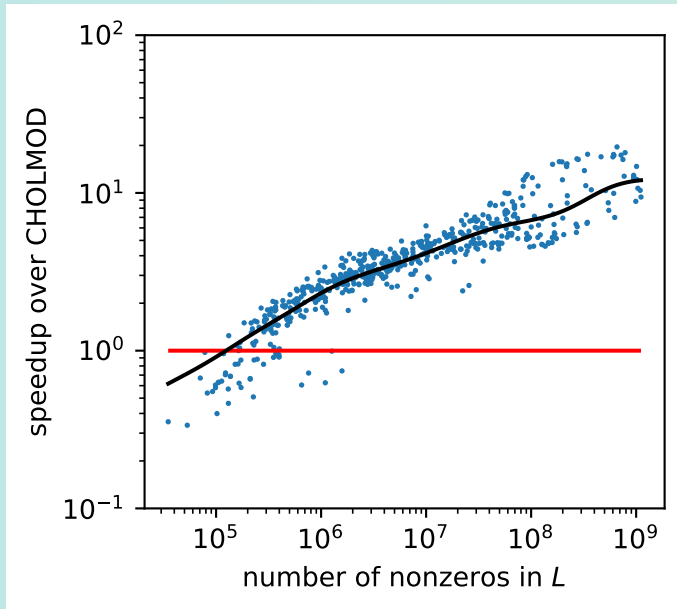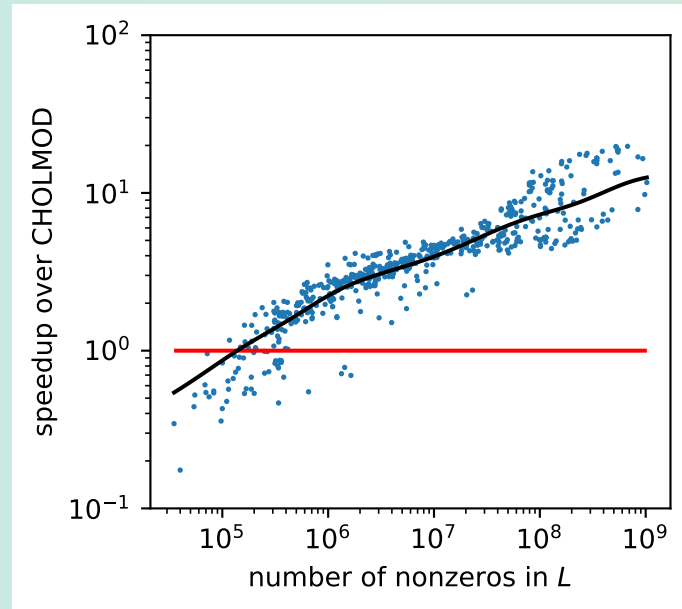# Performance Evaluation: Results AMD MI210 (Speedup)
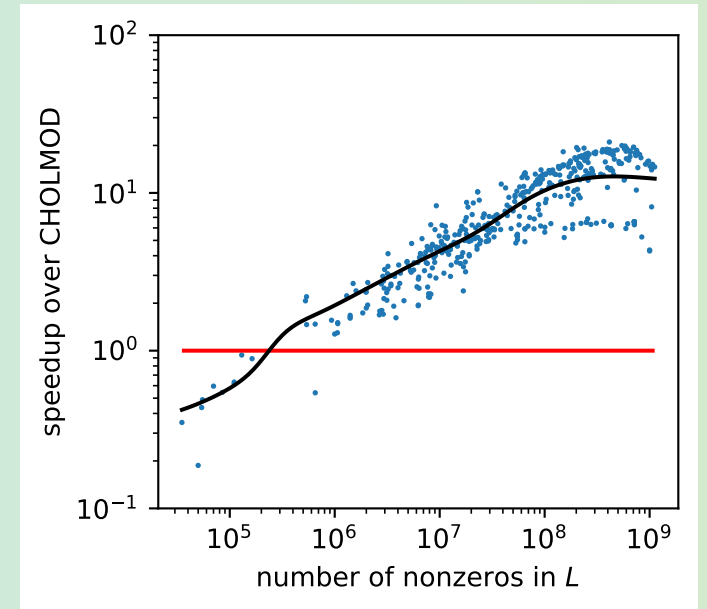


Nested Dissection

AMD

Natural

# Performance Evaluation: Results NVIDIA A100 (Speedup)
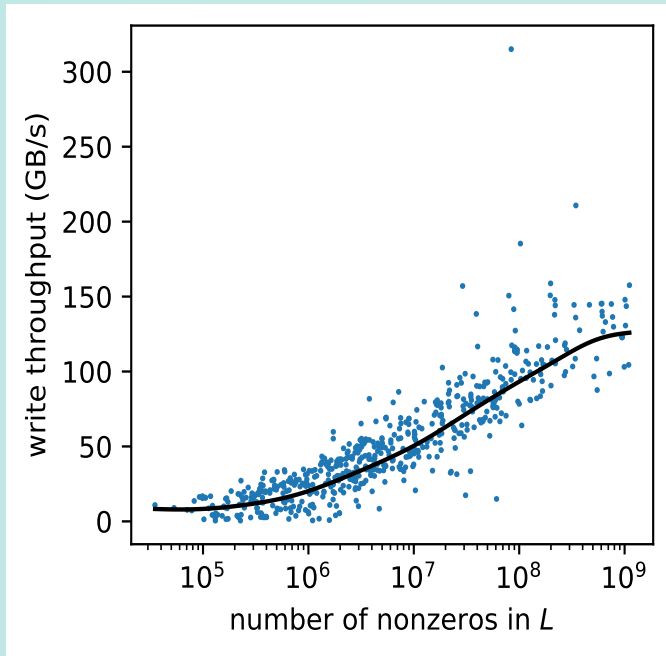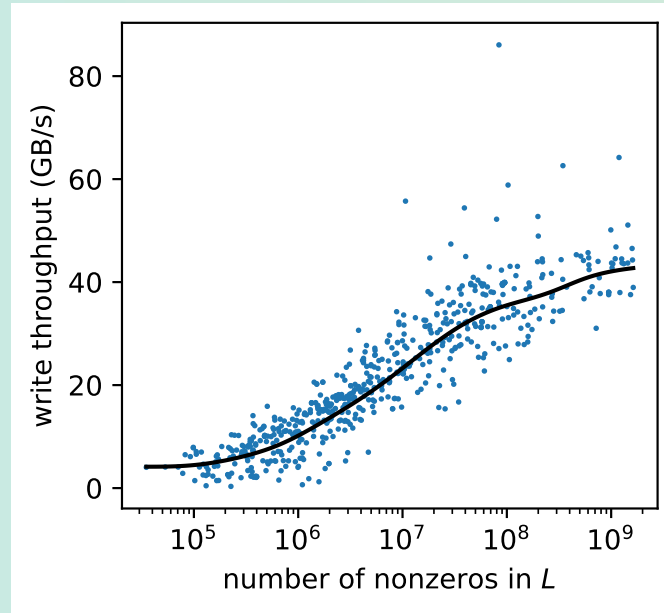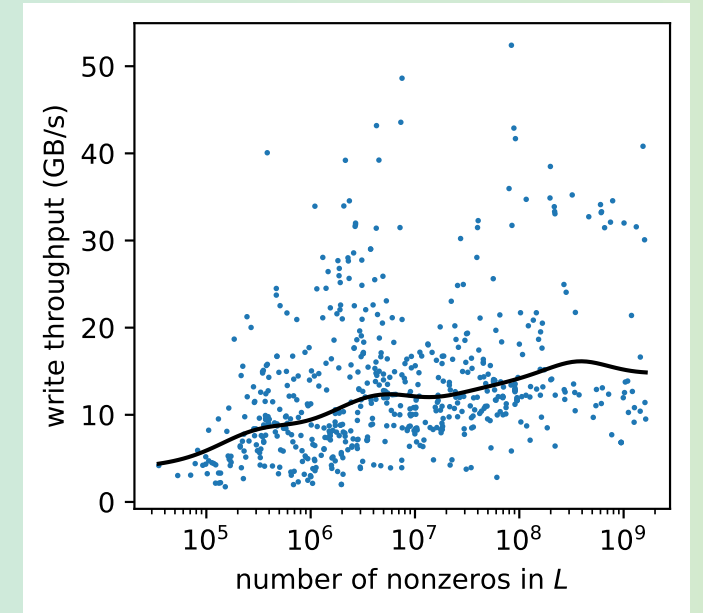


Nested Dissection

AMD

Natural

# Performance Evaluation: Throughput
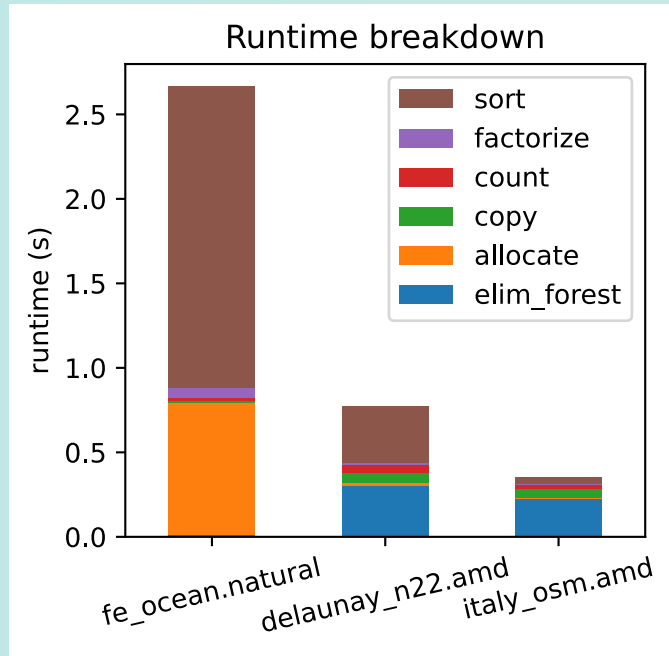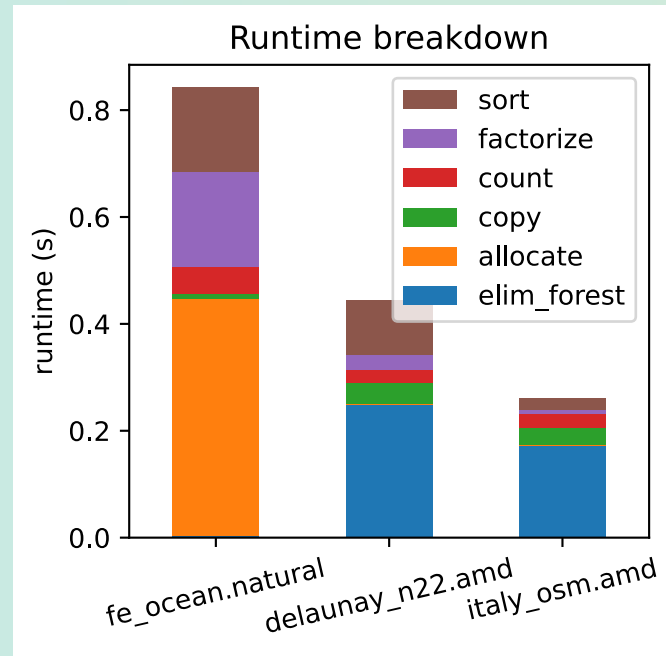


A100

MI210

Intel CPU

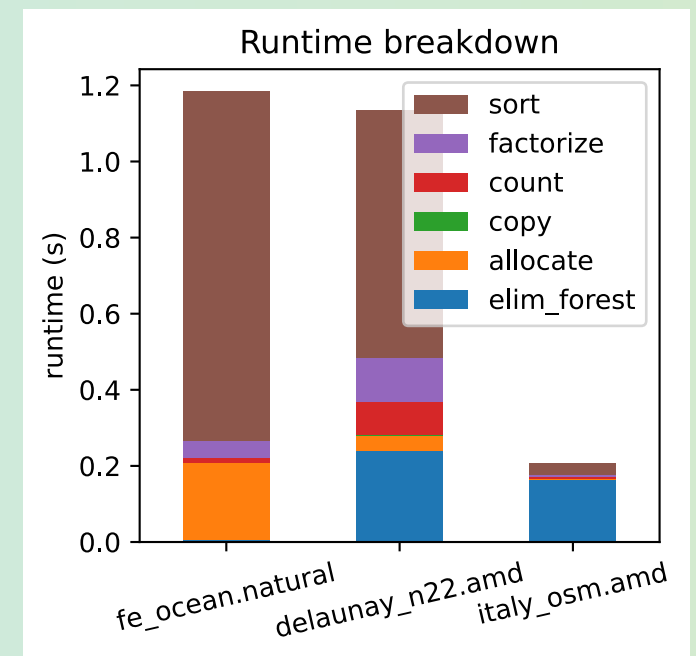Sparsity Pattern generation only on Nested Dissection ordering

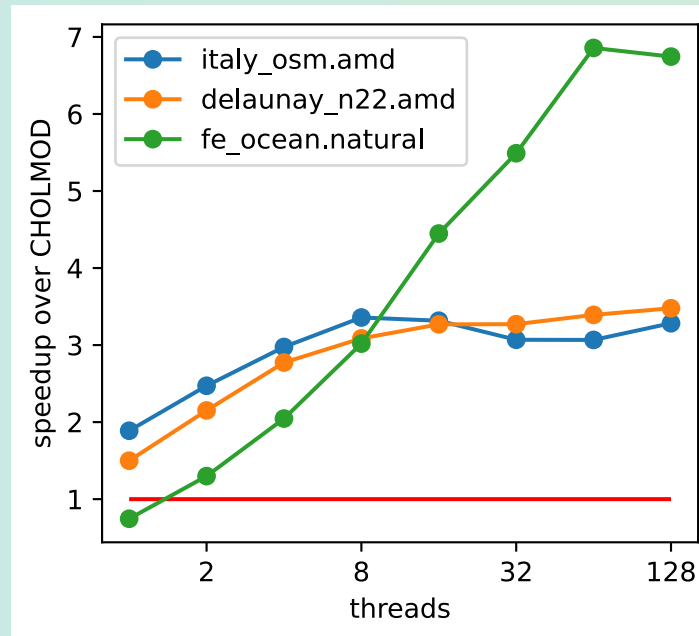# Performance Evaluation: Breakdown



A100

MI210

Intel CPU

# Performance Evaluation: Scaling on Intel CPU

# Future Work in Progress: Near-Symmetric LU

- If we symmetrize via $A_{symm} = A + A^T$, then the fill-in of $A$ is a subset of the factors $L_{symm} + L_{symm}^T$ of $A_{symm}$
  $\rightarrow$ Follow Symbolic Cholesky on $A_{symm}$ with "numerical factorization" with $A$ stored inside 0/1 binary $A_{symm}$
  $\rightarrow$ Result tells us which fill-in entries in $L_{symm} + L_{symm}{}^T$ are actually present in factors $L + L^T$ of $A$

- First results (A100 vs. sequential baseline):

  - 3.5x speedup for AMD-ordered matrix

  - 18x speedup for input-ordered matrix

# Future Work

- Performance optimization of fill-in kernels (faster LCA lookup, parallel path traversal)

- Reduced data movement cost by computing skeleton graph of $A$

- Fully on-GPU elimination tree computation

- Fully on-GPU symbolic factorization

- Tuning numerical factorizations

# References

- Joseph W.H. Liu. 1990. The Role of Elimination Trees in Sparse Factorization, SIAM J. Matrix Anal. Appl, https://doi.org/10.1137/0611010

- Timothy A. Davis. 2006. Direct methods for sparse linear systems, SIAM

- Yanqing Chen, Timothy A. Davis, William W. Hager, and Sivasankaran Rajamanickam. 2008. Algorithm 887: CHOLMOD, Supernodal Sparse Cholesky Factorization and Update/Downdate. ACM TOMS, https://doi.org/10.1145/1391989.1391995

- Timothy A. Davis and Yifan Hu. 2011. The University of Florida sparse matrix collection. ACM TOMS, https://doi.org/10.1145/2049662.2049663

- Hartwig Anzt, Terry Cojean, Goran Flegar, Fritz Göbel, Thomas Grützmacher, Pratik Nayak, Tobias Ribizel, Yuhsiang Mike Tsai, and Enrique S. Quintana-Ortí. 2022. Ginkgo: A Modern Linear Operator Algebra Framework for High Performance Computing. ACM TOMS, https://doi.org/10.1145/3480935

# Backup Slides