# SYMPACK: A 2D TASK-BASED FACTORIZATION ALGORITHM FOR SPARSE SYMMETRIC MATRICES

Mathias Jacquelin
mjacquelin@lbl.gov

Esmond Ng

February 27 2019

Scalable Solvers Group
Computational Research Department
Lawrence Berkeley National Laboratory

Motivations:

- Sparse matrices arise in many applications:
  - Optimization problems
  - Discretized PDEs
  - Electronic structure theory
  - . . .

- Some sparse direct methods require:
  - Sparse factorizations
  - Computing some inverse elements

- Matrix $A$ is symmetric in many cases
- Symmetric storage: only lower triangular part of $A$ is stored
  - Lower memory consumption
  - Fewer floating point operations

- Many ways to schedule computations, partition data

- Challenging problem: irregular computation load
- **Crucial to remove synchronization points**

- Only lower triangular part of $A$ is stored
- Basic algorithm:

---
**Algorithm 1:** Basic Cholesky algorithm

---
for *column $j = 1$ to $n$* do

    $\ell_{j,j} = \sqrt{A_{j,j}}$
    for *row $i = j + 1$ to $n$* do
        $\ell_{i,j} = A_{i,j}/\ell_{j,j}$
    end

    for *column $k = j + 1$ to $n$* do
        for *row $i = k$ to $n$* do
            $A_{i,k} = A_{i,k} - \ell_{i,j} \cdot \ell_{k,j}$
        end

    end

end

---

- Only lower triangular part of $A$ is stored
- Basic algorithm:

---

**Algorithm 1:** Basic Cholesky algorithm

---

for *column* $j = 1$ *to* $n$ do

$\quad\ell_{j,j} = \sqrt{A_{j,j}}$
$\quad$for *row* $i = j + 1$ *to* $n$ do
$\quad\quad\ell_{i,j} = A_{i,j}/\ell_{j,j}$
$\quad$end $\qquad\qquad\qquad$ Factor column $j$

$\quad$for *column* $k = j + 1$ *to* $n$ do
$\quad\quad$for *row* $i = k$ *to* $n$ do
$\quad\quad\quad A_{i,k} = A_{i,k} - \ell_{i,j} \cdot \ell_{k,j}$
$\quad\quad$end

$\quad$end

end

---

- Only lower triangular part of $A$ is stored
- Basic algorithm:

---

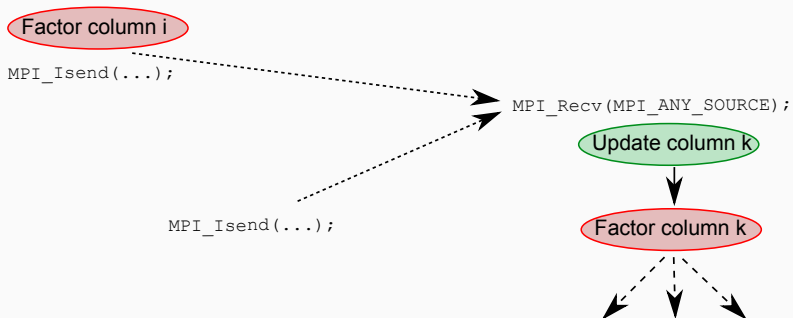**Algorithm 1:** Basic Cholesky algorithm

---

**for** *column $j = 1$ to $n$* **do**

$\quad \ell_{j,j} = \sqrt{A_{j,j}}$
$\quad$**for** *row $i = j + 1$ to $n$* **do**
$\quad\quad \ell_{i,j} = A_{i,j}/\ell_{j,j}$
$\quad$**end**

$\qquad\qquad\qquad\qquad\qquad$ Factor column $j$

$\quad$**for** *column $k = j + 1$ to $n$* **do**
$\quad\quad$**for** *row $i = k$ to $n$* **do**
$\quad\quad\quad A_{i,k} = A_{i,k} - \ell_{i,j} \cdot \ell_{k,j}$
$\quad\quad$**end**

$\quad$**end**

$\qquad\qquad\qquad\qquad\qquad$ Update next columns

**end**

---

- Only lower triangular part of $A$ is stored
- Basic algorithm:

**Algorithm 1:** Basic Cholesky algorithm

for *column $j = 1$ to $n$* do

$\quad \ell_{j,j} = \sqrt{A_{j,j}}$

$\quad$ for *row $i = j + 1$ to $n$* do

$\quad\quad \ell_{i,j} = A_{i,j}/\ell_{j,j}$

$\quad$ end

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ **Factor** column $j$

$\quad$ for *column $k = j + 1$ to $n$* do

$\quad\quad\quad\quad\quad\quad\quad\quad$ **Update** next columns
$\quad\quad\quad\quad\quad\quad\quad$ and **Aggregate** updates

$\quad\quad$ for *row $i = k$ to $n$* do

$\quad\quad\quad A_{i,k} = A_{i,k} - \ell_{i,j} \cdot \ell_{k,j}$

$\quad\quad$ end

$\quad$ end

end

- Only lower triangular part of $A$ is stored
- Basic algorithm:

**Algorithm 1:** Basic Cholesky algorithm

for *column $j = 1$ to $n$* do

$\quad \ell_{j,j} = \sqrt{A_{j,j}}$
$\quad$ for *row $i = j + 1$ to $n$* do
$\quad\quad \ell_{i,j} = A_{i,j}/\ell_{j,j}$
$\quad$ end
$\qquad\qquad\qquad\qquad\qquad$ **Factor** column $j$

$\quad$ for *column $k = j + 1$ to $n$* do
$\quad\quad$ for *row $i = k$ to $n$* do
$\quad\quad\quad A_{i,k} = A_{i,k} - \ell_{i,j} \cdot \ell_{k,j}$
$\quad\quad$ end

$\quad$ end

$\qquad\qquad\qquad\qquad\qquad$ **Update** next columns
$\qquad\qquad\qquad\qquad\qquad$ and **Aggregate** updates
$\qquad\qquad\qquad\qquad\qquad$ for *row $i = k$ to $n$* do
$\qquad\qquad\qquad\qquad\qquad\quad tmp_i = tmp_i + \ell_{i,j} \cdot \ell_{k,j}$
$\qquad\qquad\qquad\qquad\qquad$ end
$\qquad\qquad\qquad\qquad\qquad A_{*,k} = A_{*,k} - tmp_*$

end

Asynchronous comm. becomes blocking when out of buffer

Asynchronous comm. becomes blocking when out of buffer

Deadlock issues

Asynchronous comm. becomes blocking when out of buffer

Deadlock issues

- Deadlock prevention is difficult:
  - Order in operations/messages

Asynchronous comm. becomes blocking when out of buffer

<span style="color:red">Deadlock issues</span>

- Deadlock prevention is difficult:
  - Order in operations/messages
    <span style="color:red">Potential over-synchronization</span>

Asynchronous comm. becomes blocking when out of buffer

## Deadlock issues

· Deadlock prevention is difficult:
  · Order in operations/messages
      Potential over-synchronization

· "Pull" strategy (one sided communications)
  · Signal data when available
  · Receiver gets data when ready

Run times on boneS10 for three variants of **symPACK**

- **symPACK- Push**
- **symPACK- Pull**

Time (s)

Processor count

n=914,898    nnz(A)=20,896,803    nnz(L)=318,019,434

- Per-task dependency counts
- Update dependencies as messages are flowing in
- Maintain a list of tasks ready for execution

Run times on boneS10 for three variants of **symPACK**

n=914,898     nnz(A)=20,896,803     nnz(L)=318,019,434

- 2D block cyclic used in many solvers
- Works well in practice for sparse matrices as well
- However, nothing is explicitly balanced

- · 2D block cyclic used in many solvers
- · Works well in practice for sparse matrices as well
- · However, nothing is explicitly balanced

- · Can we do better?

- 2D block cyclic used in many solvers
- Works well in practice for sparse matrices as well
- However, nothing is explicitly balanced

- Can we do better?
- Can we store mapping information?
  - Cell: block "delimited" by supernode partition
  - Block: set of contiguous rows in a given supernode
  - A cell can hold multiple blocks

· Subtree-to-subcube mapping

· Subtree-to-subcube mapping

· Subtree-to-subcube mapping

- Subtree-to-subcube mapping
- Non-empty cells distributed within groups

· Subtree-to-subcube mapping
· Non-empty cells distributed within groups

· Subtree-to-subcube mapping
· Non-empty cells distributed within groups

· Subtree-to-subcube mapping
· Non-empty cells distributed within groups

- A **Future** is a synchronization object for asynchronous operations:
  - When the operation is complete, future becomes *ready*
  - A *callback* can be attached to a future
- A **Promise** can be thought as a counter:
  - Associated with a future
  - Future is ready when the count reaches 0.

incoming dependence
available remotely

incoming dependence
transferred

- A **Future** is a synchronization object for asynchronous operations:
  - When the operation is complete, future becomes *ready*
  - A *callback* can be attached to a future
- A **Promise** can be thought as a counter:
  - Associated with a future
  - Future is ready when the count reaches 0.
- Each task has two **Promise**s (a counter)

incoming dependence available remotely

incoming dependence transferred

data available tasks

ready tasks

T$_k$

① 

n  m

data

incoming dependence
available remotely

incoming dependence
transferred

data available tasks

ready tasks

Run times for Flan_1565

n=1,564,794     nnz(L)=1,574,541,576

Run times for audikw_1

Legend:
- symPACK_2D
- symPACK_1D
- pastix_5_2_3

Time (s) vs. Processes (32, 64, 128, 256, 512, 1024)

n=943,695      nnz(L)=1,261,342,196

Run times for serena

n=1,391,349    nnz(L)=2,821,178,652

Run times for phosphorene

n=512,000     nnz(L)=1,697,433,600

- Aggregate updates using a tree pattern
- 1D data distribution at leaves
- Use tasks to implement 3D type of layout at higher levels (multiple tasks on the same cell)

- Accelerator / GPU support
  - Upcoming UPC++ with seamless local/remote host/device memory accesses
  - Batched BLAS

- Acknowledgments:
  - DOE SciDAC FASTMath, CompCat, ComPASS4
  - ECP Pagoda