

# SCHEDULING SPARSE SYMMETRIC FAN-BOTH CHOLESKY FACTORIZATION

---

Mathias Jacquelin

`mjacquelin@lbl.gov`

Esmond Ng, Kathy Yelick and Yili Zheng

`engng|kayelick|yzheng@lbl.gov`

May 18 2016

Scalable Solvers Group

Computational Research Department

Lawrence Berkeley National Laboratory

Background and motivation

Fan-In, Fan-Out and Fan-Both factorizations

Parallel distributed memory implementation, a.k.a. **symPACK**

Numerical experiments

### Motivations:

- Sparse matrices arise in many applications:
  - Optimization problems
  - Discretized PDEs
  - ...
- Some sparse matrices are symmetric

### Motivations:

- Sparse matrices arise in many applications:
  - Optimization problems
  - Discretized PDEs
  - ...
- Some sparse matrices are symmetric

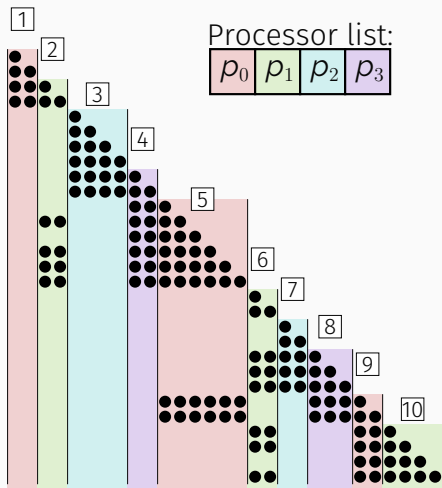
### Challenges for current and future platforms:

- Higher relative communication costs
- Lower amount of memory per core

### Objective:

- Compute sparse  $A = LL^T$  factorization
- $A$  is sparse symmetric matrix
- $A$  is positive definite
- Need to exploit symmetry
- $L$  is a lower triangular matrix

# SPARSE MATRICES AND ELIMINATION TREE

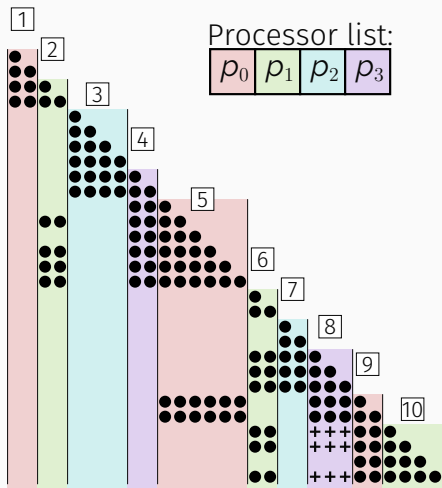


· Fill in,  $\Omega(A) \subseteq \Omega(L)$

$$A = LL^T$$

$\Omega(A)$  is the sparsity pattern of  $A$

# SPARSE MATRICES AND ELIMINATION TREE

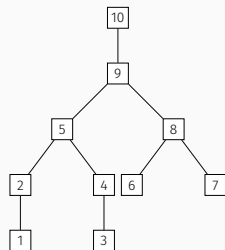
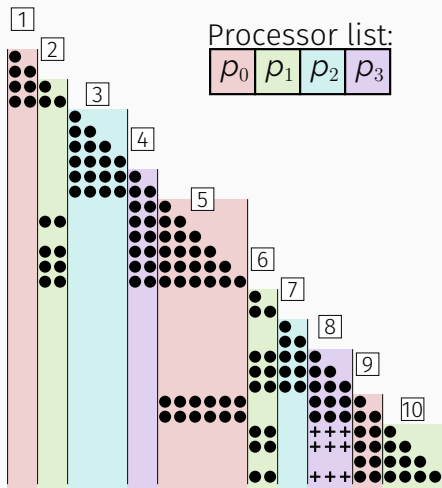


· Fill in,  $\Omega(A) \subseteq \Omega(L)$

$$A = LL^T$$

$\Omega(A)$  is the sparsity pattern of  $A$

# SPARSE MATRICES AND ELIMINATION TREE



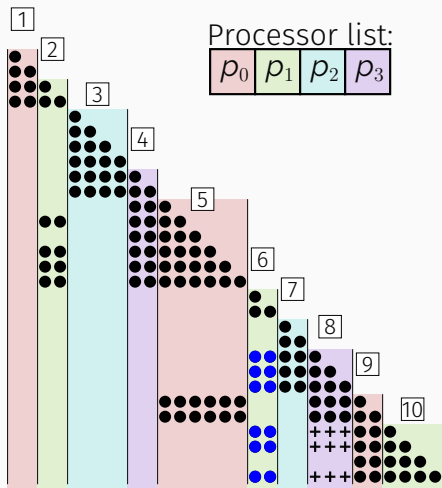
- Elim. tree represents column dependences
- Fill in,  $\Omega(A) \subseteq \Omega(L)$

$$A = LL^T$$

$\Omega(A)$  is the sparsity pattern of  $A$

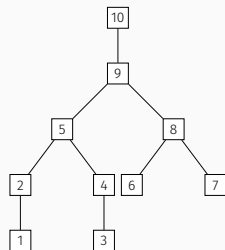


# SPARSE MATRICES AND ELIMINATION TREE



$$A = LL^T$$

$\Omega(A)$  is the sparsity pattern of  $A$



- Elim. tree represents column dependences
- Fill in,  $\Omega(A) \subseteq \Omega(L)$
- Supernode, same structure below diagonal block

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

---

**Algorithm 1:** Basic Cholesky algorithm

---

```
for column  $j = 1$  to  $n$  do
     $\ell_{jj} = \sqrt{A_{jj}}$ 
    for row  $i = j + 1$  to  $n$  do
        |  $\ell_{ij} = A_{ij} / \ell_{jj}$ 
    end

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

$$\ell_{jj} = \sqrt{A_{jj}}$$

for row  $i = j + 1$  to  $n$  do

$$\ell_{ij} = A_{ij} / \ell_{jj}$$

end

} Factor column  $j$

for column  $k = j + 1$  to  $n$  do

for row  $i = k$  to  $n$  do

$$A_{i,k} = A_{i,k} - \ell_{ij} \cdot \ell_{kj}$$

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

$\ell_{jj} = \sqrt{A_{jj}}$   
    for row  $i = j + 1$  to  $n$  do  
        |  $\ell_{ij} = A_{ij} / \ell_{jj}$   
    end  
    } Factor column  $j$

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

$\ell_{jj} = \sqrt{A_{jj}}$   
    for row  $i = j + 1$  to  $n$  do  
        |  $\ell_{ij} = A_{ij} / \ell_{jj}$   
    end  
    } Factor column  $j$

    for column  $k = j + 1$  to  $n$  do  
        for row  $i = k$  to  $n$  do  
            |  $A_{i,k} = A_{i,k} - \ell_{ij} \cdot \ell_{kj}$   
        end  
    end  
    } Update next columns  
    and Aggregate updates

end

---

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

---

## Algorithm 1: Basic Cholesky algorithm

---

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

$\ell_{jj} = \sqrt{A_{jj}}$ for row $i = j + 1$ to $n$ do   $\ell_{ij} = A_{ij} / \ell_{jj}$ end	}	Factor column $j$
---	---	-------------------

for column $k = j + 1$ to $n$ do for row $i = k$ to $n$ do   $A_{i,k} = A_{i,k} - \ell_{ij} \cdot \ell_{kj}$ end end	}	Update next columns and Aggregate updates for row $i = k$ to $n$ do   $tmp_i = tmp_i + \ell_{ij} \cdot \ell_{kj}$ end $A_{*,k} = A_{*,k} - tmp_*$
--	---	--

end

---

- Three families [Ashcraft'95]:
  - Fan-In: “fanning-in updates”
    - Reduce **aggregate vectors** (updates)
    - Factorize column
    - Compute all updates from *that column* locally

- Three families [Ashcraft'95]:
  - Fan-In: “fanning-in updates”
    - Reduce **aggregate vectors** (updates)
    - Factorize column
    - Compute all updates from *that column* locally
  - Fan-Out: “fanning-out factors”
    - Factorize column
    - **Distribute the Cholesky factor**
    - Compute and apply all updates to *my column*.



- Three families [Ashcraft'95]:
  - Fan-In: “fanning-in updates”
    - Reduce **aggregate vectors** (updates)
    - Factorize column
    - Compute all updates from *that column* locally
  - Fan-Out: “fanning-out factors”
    - Factorize column
    - **Distribute the Cholesky factor**
    - Compute and apply all updates to *my column*.

Family determined by type of data exchanged

- Three families [Ashcraft'95]:
  - Fan-In: “fanning-in updates”
    - Reduce **aggregate vectors** (updates)
    - Factorize column
    - Compute all updates from *that column* locally
  - Fan-Out: “fanning-out factors”
    - Factorize column
    - **Distribute the Cholesky factor**
    - Compute and apply all updates to *my column*.

Family determined by type of data exchanged

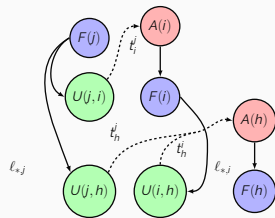
Fan-In, Fan-Out  $\subset$  Fan-Both

- Three families [Ashcraft'95]: **Fan-In**, **Fan-Out**  $\subset$  **Fan-Both**
- Task based algorithm:
  - $A(i)$ : accumulation of **aggregate vectors** (updates) to column  $i$ 

Reduces the aggregate vectors  $t_i^*$
  - $F(j)$ : factorization of col.  $j$ 

Produces cholesky **factor**  $\ell_{*,j}$
  - $U(j,i)$ : update of col.  $i$  with col.  $j$ 

Put the update in an (temporary) **aggregate vector**  $t_i^j$

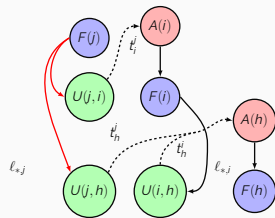


- Three families [Ashcraft'95]: **Fan-In**, **Fan-Out**  $\subset$  **Fan-Both**
- Task based algorithm:
  - $A(i)$ : accumulation of **aggregate vectors** (updates) to column  $i$ 

Reduces the aggregate vectors  $t_i^*$
  - $F(j)$ : factorization of col.  $j$ 

Produces cholesky **factor**  $\ell_{*,j}$
  - $U(j,i)$ : update of col.  $i$  with col.  $j$ 

Put the update in an (temporary) **aggregate vector**  $t_i^j$

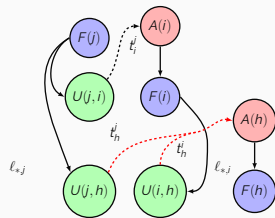


- Three families [Ashcraft'95]: **Fan-In**, **Fan-Out**  $\subset$  **Fan-Both**
- Task based algorithm:
  - $A(i)$ : accumulation of **aggregate vectors** (updates) to column  $i$ 

Reduces the aggregate vectors  $t_i^*$
  - $F(j)$ : factorization of col.  $j$ 

Produces cholesky **factor**  $\ell_{*,j}$
  - $U(j,i)$ : update of col.  $i$  with col.  $j$ 

Put the update in an (temporary) **aggregate vector**  $t_i^j$



# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size

1D Cyclic distribution

1	2	3	4
---	---	---	---



Virtual 2D mapping  $\mathcal{M}$

1	1	3	3
2	2	4	4
1	1	3	3
2	2	4	4

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal

1D Cyclic distribution

1	2	3	4
---	---	---	---



Virtual 2D mapping  $\mathcal{M}$

1	1	3	3
2	2	4	4
1	1	3	3
2	2	4	4

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings

1D Cyclic distribution

1	2	3	4
---	---	---	---



Virtual 2D mapping  $\mathcal{M}$

1	1	3	3
2	2	4	4
1	1	3	3
2	2	4	4



- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings

1D Cyclic distribution

1	2	3	4
---	---	---	---

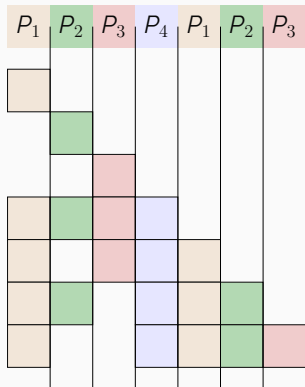


Virtual 2D mapping  $\mathcal{M}$

1	2	1	2
2	2	1	2
1	1	3	4
2	2	4	4

# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings



# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings

$P_1$	$P_2$	$P_3$	$P_4$	$P_1$	$P_2$	$P_3$
$P_1$	$P_1$	$P_3$	$P_3$	$P_1$	$P_1$	$P_3$
$P_1$	$P_2$	$P_4$	$P_4$	$P_2$	$P_2$	$P_4$
$P_3$	$P_4$	$P_3$	$P_3$	$P_1$	$P_1$	$P_3$
$P_3$	$P_4$	$P_3$	$P_4$	$P_2$	$P_2$	$P_4$
$P_1$	$P_2$	$P_1$	$P_2$	$P_1$	$P_1$	$P_3$
$P_1$	$P_2$	$P_1$	$P_2$	$P_1$	$P_2$	$P_4$
$P_3$	$P_4$	$P_3$	$P_4$	$P_3$	$P_4$	$P_3$

# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$

	$P_1$	$P_2$	$P_3$	$P_4$	$P_1$	$P_2$	$P_3$
$P_1$	$P_1$	$P_1$	$P_3$	$P_3$	$P_1$	$P_1$	$P_3$
	$P_1$	$P_2$	$P_4$	$P_4$	$P_2$	$P_2$	$P_4$
	$P_3$	$P_4$	$P_3$	$P_3$	$P_1$	$P_1$	$P_3$
	$P_3$	$P_4$	$P_3$	$P_4$	$P_2$	$P_2$	$P_4$
	$P_1$	$P_2$	$P_1$	$P_2$	$P_1$	$P_1$	$P_3$
	$P_1$	$P_2$	$P_1$	$P_2$	$P_1$	$P_2$	$P_4$
	$P_3$	$P_4$	$P_3$	$P_4$	$P_3$	$P_4$	$P_3$

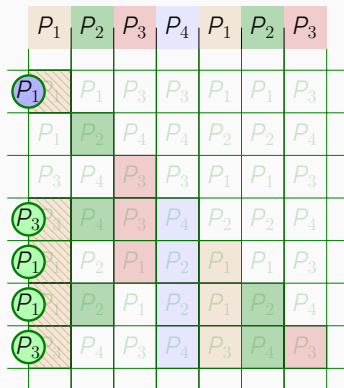
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i,i)$

	$P_1$	$P_2$	$P_3$	$P_4$	$P_1$	$P_2$	$P_3$
$P_1$	$P_1$	$P_1$	$P_3$	$P_3$	$P_1$	$P_1$	$P_3$
	$P_1$	$P_2$	$P_4$	$P_4$	$P_2$	$P_2$	$P_4$
	$P_3$	$P_4$	$P_3$	$P_3$	$P_1$	$P_1$	$P_3$
	$P_3$	$P_4$	$P_3$	$P_4$	$P_2$	$P_2$	$P_4$
	$P_1$	$P_2$	$P_1$	$P_2$	$P_1$	$P_1$	$P_3$
	$P_1$	$P_2$	$P_1$	$P_2$	$P_1$	$P_2$	$P_4$
	$P_3$	$P_4$	$P_3$	$P_4$	$P_3$	$P_4$	$P_3$

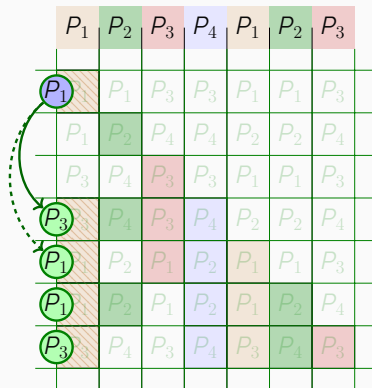
# fan-both MAPPINGS

- How do we map tasks?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$



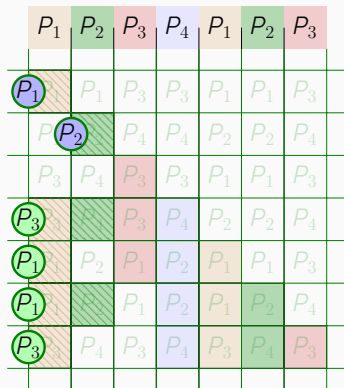
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$



# fan-both MAPPINGS

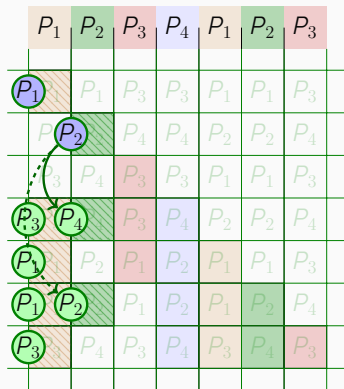
- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$





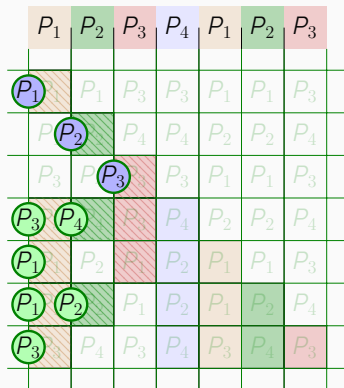
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$



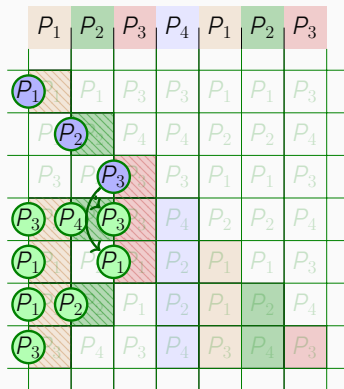
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$



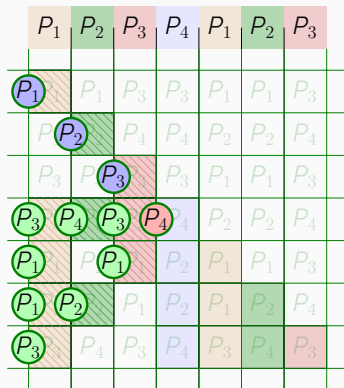
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$



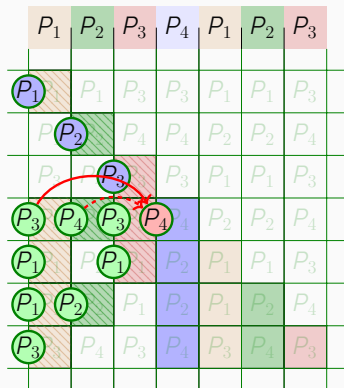
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$
  - $A(j)$  on  $\mathcal{M}(j, j)$



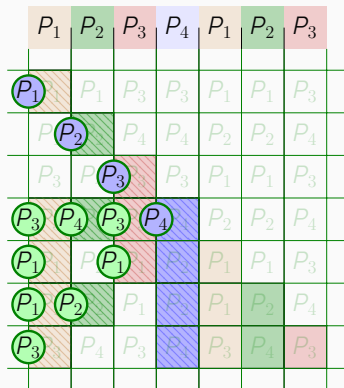
# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$
  - $A(j)$  on  $\mathcal{M}(j, j)$



# fan-both MAPPINGS

- How do we map tasks ?  
(independently of data)
- Use of 2D computation mapping grid  $\mathcal{M}$ 
  - Mapping grid “extends” to matrix size
  - Better if  $P$  processors on diagonal
  - Many possible mappings
  - $F(i)$  on proc.  $\mathcal{M}(i, i)$
  - $U(j, i)$  on  $\mathcal{M}(j, i)$
  - $A(j)$  on  $\mathcal{M}(j, j)$



0	1	2	3	0	1
0	1	2	3	0	1
0	1	2	3	0	1
0	1	2	3	0	1
0	1	2	3	0	1
0	1	2	3	0	1

**Fan-In**

$$\mathcal{M}_{i,j} = \text{mod}(i, P)$$

0	0	0	0	0	0
1	1	1	1	1	1
2	2	2	2	2	2
3	3	3	3	3	3
0	0	0	0	0	0
1	1	1	1	1	1

**Fan-Out**

$$\mathcal{M}_{i,j} = \text{mod}(j, P)$$

0	0	2	2	0	0
1	1	3	3	1	1
0	0	2	2	0	0
1	1	3	3	1	1
0	0	2	2	0	0
1	1	3	3	1	1

**Fan-Both**

$$\mathcal{M}_{i,j} = \frac{\text{mod}(\min(i,j), P) + \text{mod}(\max(i,j), P)}{P}$$

Three different computation maps, corresponding to  
Fan-In, Fan-Out and Fan-Both

- Remove synchronization points
  - Asynchronous point to point send
  - Group communication:  
(MPI) Collectives probably not the way to go
    - Requires too many communicators
    - Efficient non blocking collectives needed
    - Collective nature



- Remove synchronization points
  - Asynchronous point to point send
  - Group communication:  
(MPI) Collectives probably not the way to go
    - Requires too many communicators
    - Efficient non blocking collectives needed
    - Collective nature
  - Asynchronous tree-based group communications
  - Non-collectives = full asynchronicity

- Remove synchronization points
  - Asynchronous point to point send
  - Group communication:  
(MPI) Collectives probably not the way to go
    - Requires too many communicators
    - Efficient non blocking collectives needed
    - Collective nature
    - Asynchronous tree-based group communications
    - Non-collectives = full asynchronicity
- Minimize memory operations
  - Row-major layout

- Remove synchronization points
  - Asynchronous point to point send
  - Group communication:  
(MPI) Collectives probably not the way to go
    - Requires too many communicators
    - Efficient non blocking collectives needed
    - Collective nature
    - Asynchronous tree-based group communications
    - Non-collectives = full asynchronicity
- Minimize memory operations
  - Row-major layout
  - Avoid making extra copies when sending data

- All operations described by task  $T_{src \rightarrow tgt}$
- Message  $Msg_{src \rightarrow tgt}$
- “Push” strategy natural with MPI

- All operations described by task  $T_{src \rightarrow tgt}$
- Message  $Msg_{src \rightarrow tgt}$
- “Push” strategy natural with MPI

**Asynchronous comm. becomes blocking when out of buffer**

- All operations described by task  $T_{src \rightarrow tgt}$
- Message  $Msg_{src \rightarrow tgt}$
- “Push” strategy natural with MPI

Asynchronous comm. becomes blocking when out of buffer

Deadlock issues

- All operations described by task  $T_{src \rightarrow tgt}$
- Message  $Msg_{src \rightarrow tgt}$
- “Push” strategy natural with MPI

**Asynchronous comm. becomes blocking when out of buffer**

## Deadlock issues

- Deadlock prevention is difficult:
  - Total order in operations/messages  
(Also observed by Amestoy et al.)
  - Order by non decreasing  $tgt$ , then  $src$ :  
⇒ Use of priority queue for tasks/messages

- All operations described by task  $T_{src \rightarrow tgt}$
- Message  $Msg_{src \rightarrow tgt}$
- “Push” strategy natural with MPI

Asynchronous comm. becomes blocking when out of buffer

## Deadlock issues

- Deadlock prevention is difficult:
  - Total order in operations/messages  
(Also observed by Amestoy et al.)
  - Order by non decreasing  $tgt$ , then  $src$ :
    - ⇒ Use of priority queue for tasks/messages

Potential over-synchronization



- All operations described by task  $T_{src \rightarrow tgt}$
- Message  $Msg_{src \rightarrow tgt}$
- “Push” strategy natural with MPI

**Asynchronous comm. becomes blocking when out of buffer**

## Deadlock issues

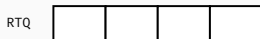
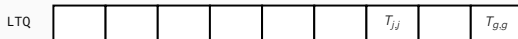
- Deadlock prevention is difficult:
  - Total order in operations/messages  
(Also observed by Amestoy et al.)
  - Order by non decreasing  $tgt$ , then  $src$ :  
⇒ Use of priority queue for tasks/messages

## Potential over-synchronization

- “Pull” strategy (one sided communications)
  - Signal data when available
  - Receiver gets data when ready

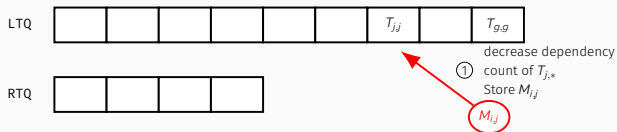
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ



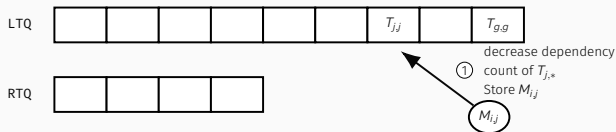
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count



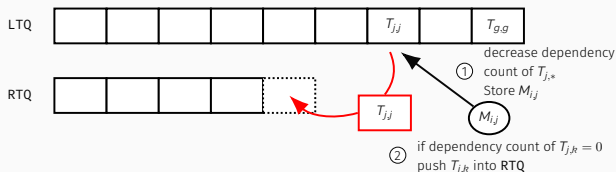
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count



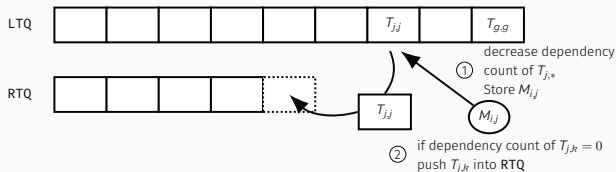
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ



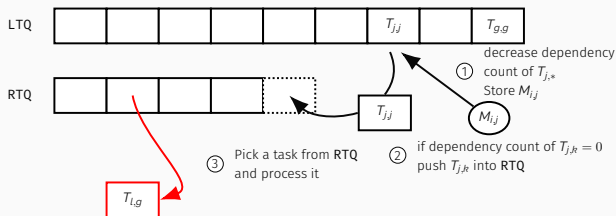
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ



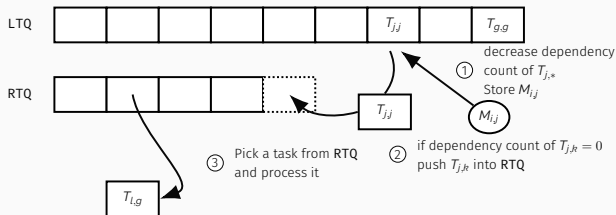
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ



# TASK SCHEDULING IN sympack

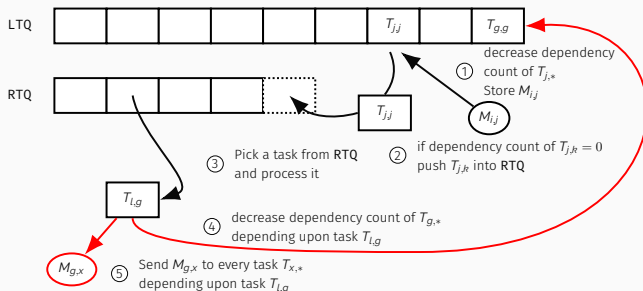
- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ





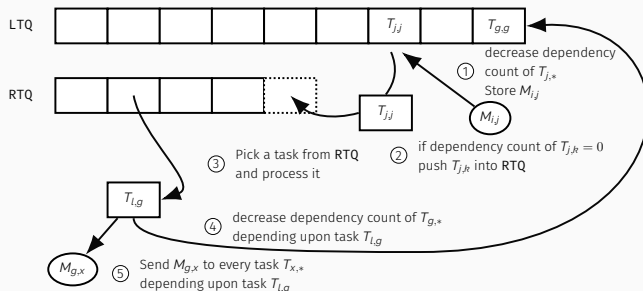
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ



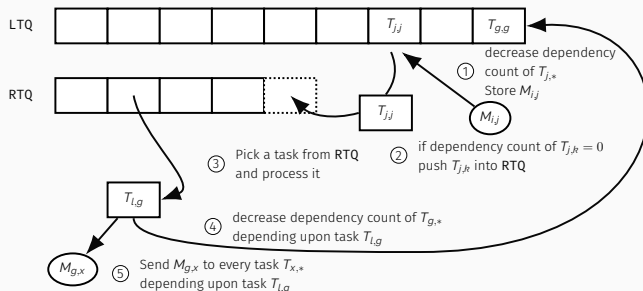
# TASK SCHEDULING IN sympack

- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ



# TASK SCHEDULING IN sympack

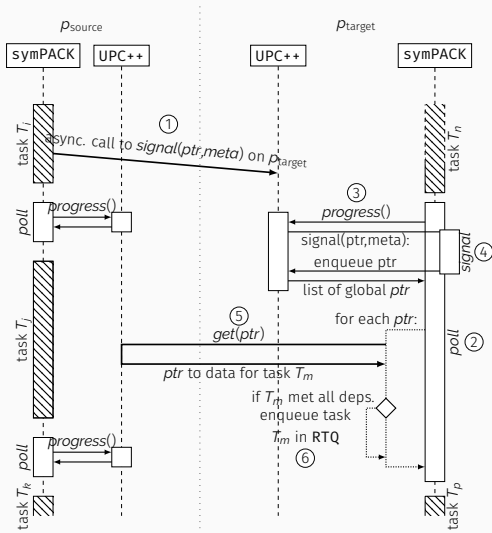
- Tasks  $T_{src \rightarrow tgt}$
- Tasks currently mapped statically
- Processor manages local task queue LTQ
  - Dependency count
  - Ready tasks are placed in RTQ



Scheduling policy ? FIFO, close to diagonal, etc.

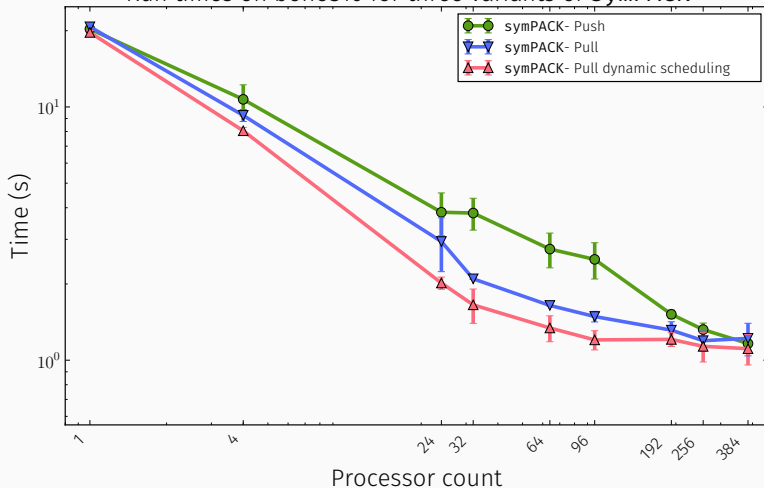
# NOTIFICATION AND COMMUNICATIONS IN symPACK

- UPC++ and GASNet for communications
- global pointer to remote memory
- one-sided communications
- asynchronous remote functions calls



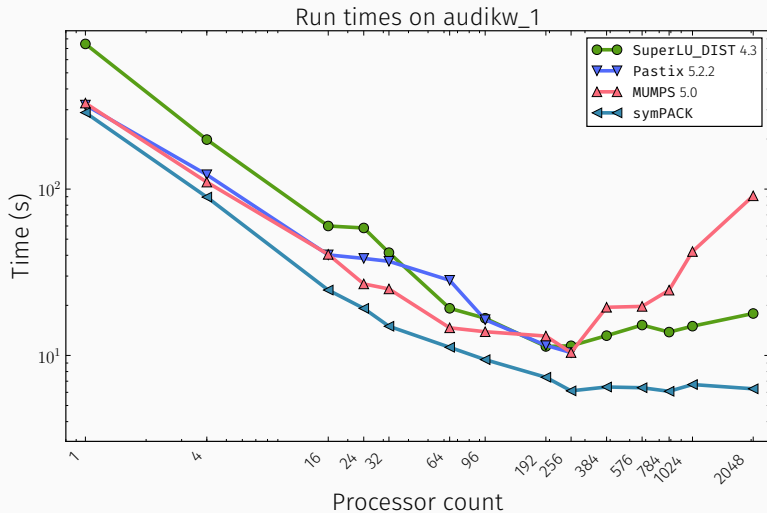
# IMPACT OF COMMUNICATION STRATEGY AND SCHEDULING

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



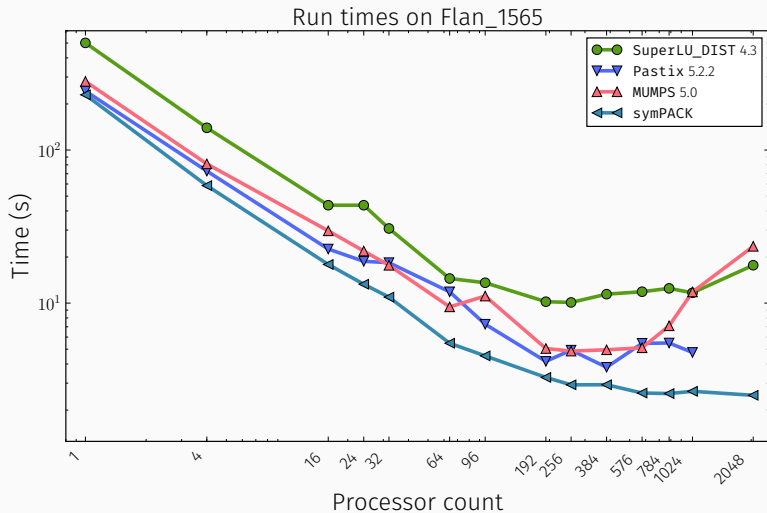
$n=914,898$     $\text{nnz}(A)=20,896,803$     $\text{nnz}(L)=318,019,434$

# STRONG SCALING VS. STATE-OF-THE-ART



$n=943,695$   $\text{nnz}(A)=39,297,771$   $\text{nnz}(L)=1,221,674,796$

# STRONG SCALING VS. STATE-OF-THE-ART



$n=1,564,794$   $\text{nnz}(A)=57,865,083$   $\text{nnz}(L)=1,574,541,576$

## SPEEDUP VS. STATE-OF-THE-ART VS. SUMMARY

Problem	Speedup vs. sym.			Speedup vs. best		
	min	max	<b>avg.</b>	min	max	<b>avg.</b>
G3_circuit	0.24	5.70	<b>1.07</b>	0.24	5.70	<b>1.07</b>
Flan_1565	1.06	9.40	<b>2.11</b>	1.06	7.07	<b>1.94</b>
af_shell7	0.89	10.61	<b>3.61</b>	0.89	7.77	<b>3.21</b>
audikw_1	1.11	14.46	<b>3.14</b>	1.11	2.84	<b>1.77</b>
boneS10	–	–	–	0.86	4.73	<b>1.75</b>
bone010	1.06	16.83	<b>3.34</b>	1.06	2.03	<b>1.47</b>



- Reduces communication cost in theory [Ashcraft'95]
- Increases parallelism during updates

- Reduces communication cost in theory [Ashcraft'95]
- Increases parallelism during updates
- Avoiding deadlocks is challenging (Similar to observation by Larkar et al.)
- New symmetric solver **symPACK**
  - implements **Fan-Both**
  - Task based Cholesky requires fine / dynamic scheduling
  - **One sided approach using UPC++**
  - Asynchronous task execution model
  - dynamic scheduling

- 2D wrap mapping performance
- Conflict with load balancing (proportional mapping) ?
- Tree-based group communications
- Hybrid parallelism (OpenMP)
- Data distribution (2D, block based ?)
- Scheduling strategies
- New task mapping policies
- Parallel ordering becomes a bottleneck

- 2D wrap mapping performance
- Conflict with load balancing (proportional mapping) ?
- Tree-based group communications
- Hybrid parallelism (OpenMP)
- Data distribution (2D, block based ?)
- Scheduling strategies
- New task mapping policies
- Parallel ordering becomes a bottleneck

Async. model important for scalability and to tolerate variability

- 2D wrap mapping performance
- Conflict with load balancing (proportional mapping) ?
- Tree-based group communications
- Hybrid parallelism (OpenMP)
- Data distribution (2D, block based ?)
- Scheduling strategies
- New task mapping policies
- Parallel ordering becomes a bottleneck

Async. model important for scalability and to tolerate variability

[www.sympack.org](http://www.sympack.org)