

Task-parallel Sparse Incomplete Cholesky Factorization using Kokkos Portable APIs.

Kyungjoo Kim Sivasankaran Rajamanickam
H. Carter Edwards Stephen L. Olivier George Stelle

Collaborators : E. Boman J. Booth A. Bradley C. Dohrmann S. Hammond

Center for Computing Research, Sandia National Labs

Trilinos User Group Meeting, October 26-29, 2015

Overview

Kokkos Portable Task Parallel Programming Model

Task-parallel (In)complete Cholesky Factorization

Numerical Examples

Conclusion

Deep hierarchical features of current hardware

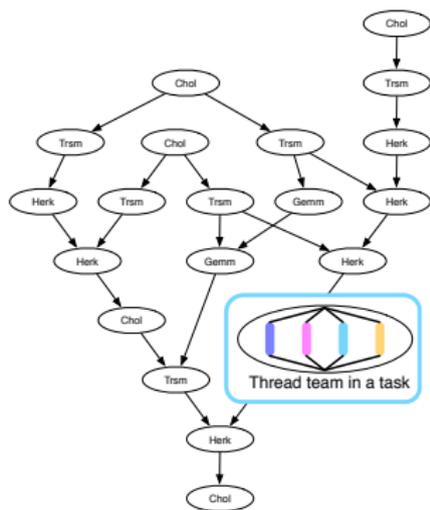
- Multi-socket, multi-processor, multi(or many)-core, multiple hardware threads.
- Multiple NUMA regions, multiple levels of caches, segmented and shared cache.

Need to expose more fine-grained parallelism

- Task parallelism is suitable for irregular problems:
e.g., producer/consumer, recursive algorithms.
 - Kokkos addresses high-level abstractions for data parallelism.
- Nested data-parallelism within a task provides better locality exploiting hardware threads.

Research highlights

- Abstractions harnessing **multiple tasking backends** to heterogeneous devices.
- Dependence driven **asynchronous** task execution with **data-parallel thread team**.
- Wait-free **respawn** task mechanism,
e.g., a task on GPUs cannot wait on dependence.
- Mini-apps (*e.g.*, sparse factorization) to support and evaluate development.



- **TaskPolicy** coordinates how and where tasks are executed
e.g., create, add_dependence, spawn (or respawn), wait;
- **Future** is a handle for tasks and allows dependence among them.

```
1 void SimpleTask() {
2     typedef Kokkos::Threads exec_space; //Serial, Threads, Qthread
3
4     Kokkos::TaskPolicy<exec_space> policy;
5     Kokkos::Future<int> f = policy.create( Functor<exec_space>() );
6
7     policy.spawn( f );
8
9     Kokkos::wait( f );
10 }
```

- **Functor** includes a user-defined function body and associated data sets.

```
11 class Functor<exec_space> {
12 public:
13     Kokkos::View<exec_space> data;
14
15     void apply( int &r_val ) {
16         r_val = doSomething( data );
17     }
18 };
```

- **DAG of tasks** is implicitly formed to guide asynchronous task execution.

```
1 void SimpleDAG() {
2     typedef Kokkos::Threads exec_space;
3
4     Kokkos::TaskPolicy<exec_space> policy;
5     Kokkos::View<exec_space> x, y; // data sets for task
6
7     auto /* future */ fx = policy.create( Functor<exec_space>( x ) );
8     auto /* future */ fy = policy.create( Functor<exec_space>( y ) );
9
10    // dependence of tasks is expressed before spawning
11    policy.add_dependence( fx, fy ); // fx is scheduled after fy
12
13    policy.spawn( fx ); // wait for now
14    policy.spawn( fy ); // may immediately execute
15
16    Kokkos::wait( policy ); // wait for all tasks to complete
17 }
```

■ Nested data parallelism with a team of threads

```
1  class Functor<exec_space> {
2  public:
3      Kokkos::View<exec_space> data;
4
5      // member is mpi-like thread communicator interface
6      // i.e., member.{team_rank,team_size,team_barrier,team_reduce}
7      void apply( const policy_type::member_type &member, int &r_val ) {
8          Kokkos::parallel_for( TeamThreadRange( member, data.size() ),
9                               [&](const int i ) {
10                                 // different indexing may be required for different
11                                 // execution space e.g., GPU interleaved data layout
12                                 int id = Index<exec_space>( member, i );
13                                 doSomething( data(id) );
14                             } );
15     }
16 };
17
18 void SimpleTaskData() {
19     typedef Kokkos::Threads exec_space;
20
21     Kokkos::TaskPolicy<exec_space> policy;
22     auto /* future */ f = policy.create_team( Functor<exec_space>() );
23
24     policy.spawn( f );
25
26     Kokkos::wait( f );
27 }
```

Task-parallel (In)complete Cholesky Factorization

Standard procedure

1. Fill-reduced (or band-reduced) ordering: Scotch.
2. Symbolic factorization: Hysom and Pothen[1].
3. **Numeric factorization.**

Design considerations for task-parallelism

- Structure-based (in)complete factorization; fills are statically determined.
 - A set of **self-contained data** within a task.
 - **Cache-friendly** numeric kernels.
 - Separation of concerns (concurrency is separated from parallelism):
 - Algorithm decomposes factorization into subproblems and provides dependence among them.
 - Runtime schedules tasks to parallel compute units.
- **Objective: portable performance on most of heterogeneous architectures.**

1 D.Hysom and A.Pothen, Level-based incomplete LU factorization: Graph model and algorithms, 2002.

Algorithms-by-blocks

- Originally developed for distributed parallel *out-of-core* matrix computations.
- Converts basic computing units from *scalar* to *blocks*.
- Used for thread-level task parallelism[2,3]: **asynchronous tasking** and **efficient level 3 BLAS**.

Algorithm: $A := \text{CHOL_UNB}(A)$

Partition $A \rightarrow \left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$
where A_{TL} is 0×0
while $\text{length}(A_{TL}) < \text{length}(A)$ **do**

Repartition

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$$

where α_{11} is a scalar

$$\alpha_{11} := \sqrt{\alpha_{11}}$$

$$a_{12}^T := a_{12}^T / \alpha_{11}$$

$$A_{22} := A_{22} - a_{12} a_{12}^T$$

Continue with

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$$

endwhile

- 2 G.Quintana-Ortí *et al.*, Programming Matrix Algorithms-by-blocks for Thread-level Parallelism, 2009.
- 3 A.Buttari *et al.*, A Class of Parallel Tiled Linear Algebra Algorithms for Multicore Architectures, 2009.

Algorithms-by-blocks

- Originally developed for distributed parallel *out-of-core* matrix computations.
- Converts basic computing units from *scalar* to *blocks*.
- Used for thread-level task parallelism[2,3]: **asynchronous tasking** and **efficient level 3 BLAS**.

<p>Algorithm: $A := \text{CHOL_BLK}(A)$</p> <p>Partition $A \rightarrow \begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix}$ where A_{TL} is 0×0 while $\text{length}(A_{TL}) < \text{length}(A)$ do Determine block size b Repartition $\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \rightarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$ where A_{11} is $b \times b$</p> <hr/> <p>$A_{11} := \text{CHOL_UNB}(A_{11})$ $A_{12} := \text{TRIU}(A_{11})^{-1} A_{12}$ $A_{22} := A_{22} - A_{12}^T A_{12}$</p> <hr/> <p>Continue with $\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \leftarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$</p> <p>endwhile</p>

k -iteration

$$\begin{pmatrix} \bar{A}_{00} & \dots & \bar{A}_{0K} & \dots & \bar{A}_{0,N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \bar{A}_{K0} & \dots & \bar{A}_{KK} & \dots & \bar{A}_{K,N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \bar{A}_{N-1,0} & \dots & \bar{A}_{N-1,K} & \dots & \bar{A}_{N-1,N-1} \end{pmatrix}$$

$$A_{11} := \text{CHOL}(\bar{A}_{kk})$$

$$\begin{aligned} A_{12} &:= \bar{A}_{kk}^{-1} (\bar{A}_{k,k+1} \dots \bar{A}_{k,n-1}) \\ &= (\bar{A}_{kk}^{-1} \bar{A}_{k,k+1} \dots \bar{A}_{kk}^{-1} \bar{A}_{k,n-1}) \end{aligned}$$

$$\begin{aligned} A_{22} &:= \begin{pmatrix} \bar{A}_{k+1,k+1} & \dots & \bar{A}_{k+1,n-1} \\ \vdots & \ddots & \vdots \\ \bar{A}_{n-1,k+1} & \dots & \bar{A}_{n-1,n-1} \end{pmatrix} - \begin{pmatrix} \bar{A}_{k+1,k} \\ \vdots \\ \bar{A}_{n-1,k} \end{pmatrix} (\bar{A}_{k,k+1} \dots \bar{A}_{k,n-1}) \\ &= \begin{pmatrix} \bar{A}_{k+1,k+1} - \bar{A}_{k+1,k} \bar{A}_{k,k+1} & \dots & \bar{A}_{k+1,n-1} - \bar{A}_{k+1,k} \bar{A}_{k,n-1} \\ \vdots & \ddots & \vdots \\ \bar{A}_{n-1,k+1} - \bar{A}_{n-1,k} \bar{A}_{k,k+1} & \dots & \bar{A}_{n-1,n-1} - \bar{A}_{n-1,k} \bar{A}_{k,n-1} \end{pmatrix} \end{aligned}$$

- G. Quintana-Ortí *et al.*, Programming Matrix Algorithms-by-blocks for Thread-level Parallelism, 2009.
- A. Buttari *et al.*, A Class of Parallel Tiled Linear Algebra Algorithms for Multicore Architectures, 2009.

Algorithms-by-blocks

- Originally developed for distributed parallel *out-of-core* matrix computations.
- Converts basic computing units from *scalar* to *blocks*.
- Used for thread-level task parallelism[2,3]: **asynchronous tasking** and **efficient level 3 BLAS**.

Algorithm: $A := \text{CHOL_BLK}(A)$

Partition $A \rightarrow \begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix}$

where A_{TL} is 0×0

while $\text{length}(A_{TL}) < \text{length}(A)$ **do**

Determine block size b

Repartition

$$\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \rightarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$$

where A_{11} is $b \times b$

$A_{11} := \text{CHOL_UNB}(A_{11})$

$A_{12} := \text{TRIU}(A_{11})^{-1} A_{12}$

$A_{22} := A_{22} - A_{12}^T A_{12}$

Continue with

$$\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \leftarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$$

endwhile

k -iteration

$$\begin{pmatrix} \bar{A}_{00} & \dots & \bar{A}_{0K} & \dots & \bar{A}_{0,N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \bar{A}_{K0} & \dots & \bar{A}_{KK} & \dots & \bar{A}_{K,N-1} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \bar{A}_{N-1,0} & \dots & \bar{A}_{N-1,K} & \dots & \bar{A}_{N-1,N-1} \end{pmatrix}$$

$A_{11} := \text{CHOL}(\bar{A}_{kk})$

$A_{12} := \bar{A}_{kk}^{-1} (\bar{A}_{k,k+1} \dots \bar{A}_{k,n-1})$ ← Tri-Solve

$$= (\bar{A}_{kk}^{-1} \bar{A}_{k,k+1} \dots \bar{A}_{kk}^{-1} \bar{A}_{k,n-1})$$

$$A_{22} := \begin{pmatrix} \bar{A}_{k+1,k+1} & \dots & \bar{A}_{k+1,n-1} \\ \text{sym} & \ddots & \vdots \\ \bar{A}_{n-1,k+1} & \dots & \bar{A}_{n-1,n-1} \end{pmatrix} - \begin{pmatrix} \bar{A}_{k+1,k} \\ \vdots \\ \bar{A}_{n-1,k} \end{pmatrix} (\bar{A}_{k,k+1} \dots \bar{A}_{k,n-1})$$

$$= \begin{pmatrix} \bar{A}_{k+1,k+1} - \bar{A}_{k+1,k} \bar{A}_{k,k+1} & \dots & \bar{A}_{k+1,n-1} - \bar{A}_{k+1,k} \bar{A}_{k,n-1} \\ \text{sym} & \ddots & \vdots \\ \bar{A}_{n-1,k+1} - \bar{A}_{n-1,k} \bar{A}_{k,k+1} & \dots & \bar{A}_{n-1,n-1} - \bar{A}_{n-1,k} \bar{A}_{k,n-1} \end{pmatrix}$$

← General Matrix-Matrix mult.

← Hermitian Rank K-update

- G. Quintana-Ortí *et al.*, Programming Matrix Algorithms-by-blocks for Thread-level Parallelism, 2009.
- A. Buttari *et al.*, A Class of Parallel Tiled Linear Algebra Algorithms for Multicore Architectures, 2009.

Algorithm: $A := \text{CHOL_BY_BLOCKS}(A)$

$$\text{Partition } A \rightarrow \begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix}$$

where A_{TL} is 0×0

while $\text{length}(A_{TL}) < \text{length}(A)$ do

Repartition

$$\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \rightarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$$

 where A_{11} is 1×1

 genTaskChol (A_{11})

 genTaskTrsm (A_{11} , A_{12})

 genTaskHerk (A_{12} , A_{22})

Continue with

$$\begin{pmatrix} A_{TL} & A_{TR} \\ A_{BL} & A_{BR} \end{pmatrix} \leftarrow \begin{pmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \\ A_{20} & A_{21} & A_{22} \end{pmatrix}$$

endwhile

```
function genTaskChol :
  Future f = create( Chol, A11(0,0) )
  add_dependency( f, A11(0,0).getFuture() )
  A11(0,0).setFuture( f )
  spawn( f )
```

```
function genTaskTrsm :
  for j in A12.nnz()
    Future f = create( Trsm, A11(0,0), A12(0,j) )
    add_dependency( f, A11(0,0).getFuture() )
    add_dependency( f, A12(0,j).getFuture() )
    A12(0,j).setFuture( f )
    spawn( f )
```

```
function genTaskHerk :
  for i in A12.nnz()
    for j in A12.nnz()
      if exist( A22(i,j) )
        Future f = create( i==j ? Herk : Gemv,
                          A12(0,i), A12(0,j),
                          A22(i,j) )
        add_dependency( f, A12(0,i).getFuture() )
        add_dependency( f, A12(0,j).getFuture() )
        add_dependency( f, A22(i,j).getFuture() )
        A21(i,j).setFuture( f )
        spawn( f )
```

- Degree of concurrency still depends on **nested dissection ordering**.
- Parallelism is **not strictly tied** with the nested dissection tree.
- Each block records future and dependence is **explicitly** described from the algorithm.

	0	1	2	3	4	5	6	7	8	9	10	11
0	0	x			x							
1		1	x			x						
2			2	x			x					
3				3				x				
4					4	x			x			
5						5	x			x		
6							6	x			x	
7								7				x
8									8	x		
9										9	x	
10											10	x
11												11

Natural ordering from a 2D mesh

Hierarchical (recursive) definition of matrices

- `CrsMatrixBase` contains sparse data arrays
i.e., row pointers, column indices, value array.
- `MatrixView` defines a rectangular region
i.e., view offsets and dimensions;
 - light-weight object with meta data only.

```
typedef CrsMatrixBase<value_type := scalar> ScalarMatrix;  
typedef MatrixView<base_object := ScalarMatrix> ViewOnScalarMatrix;  
  
typedef CrsMatrixBase<value_type := ViewOnScalarMatrix> BlockMatrix;  
typedef MatrixView<base_object := BlockMatrix> ViewOnBlockMatrix;
```

2D Partitioned-block Matrix

	0	1	2	3	4	5	6	7	8	9	10	11
0	0	X								X		
1		1	X	X								
2			2							X	X	X
3				3								X
4					4	X				X	X	
5						5			X			
6							6	X			X	X
7								7	X			
8									8		X	
9										9		
10											10	
11												11

Nested dissection ordering by Scotch

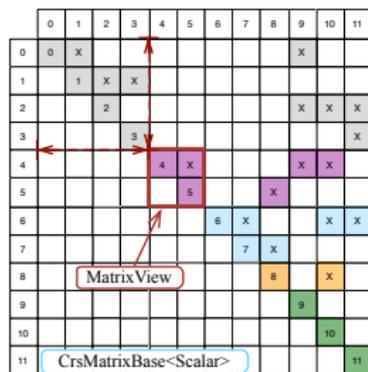
Hierarchical (recursive) definition of matrices

- `CrsMatrixBase` contains sparse data arrays *i.e.*, row pointers, column indices, value array.
- `MatrixView` defines a rectangular region *i.e.*, view offsets and dimensions;
 - light-weight object with meta data only.

```
typedef CrsMatrixBase<value_type := scalar> ScalarMatrix;
typedef MatrixView<base_object := ScalarMatrix> ViewOnScalarMatrix;

typedef CrsMatrixBase<value_type := ViewOnScalarMatrix> BlockMatrix;
typedef MatrixView<base_object := BlockMatrix> ViewOnBlockMatrix;
```

2D Partitioned-block Matrix



Matrix view

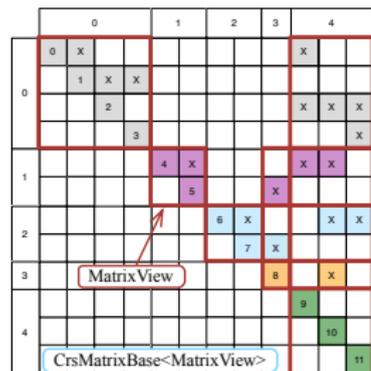
Hierarchical (recursive) definition of matrices

- `CrsMatrixBase` contains sparse data arrays *i.e.*, row pointers, column indices, value array.
- `MatrixView` defines a rectangular region *i.e.*, view offsets and dimensions;
 - light-weight object with meta data only.

```
typedef CrsMatrixBase<value_type := scalar> ScalarMatrix;  
typedef MatrixView<base_object := ScalarMatrix> ViewOnScalarMatrix;
```

```
typedef CrsMatrixBase<value_type := ViewOnScalarMatrix> BlockMatrix;  
typedef MatrixView<base_object := BlockMatrix> ViewOnBlockMatrix;
```

2D Partitioned-block Matrix



Matrix of blocks

Hierarchical (recursive) definition of matrices

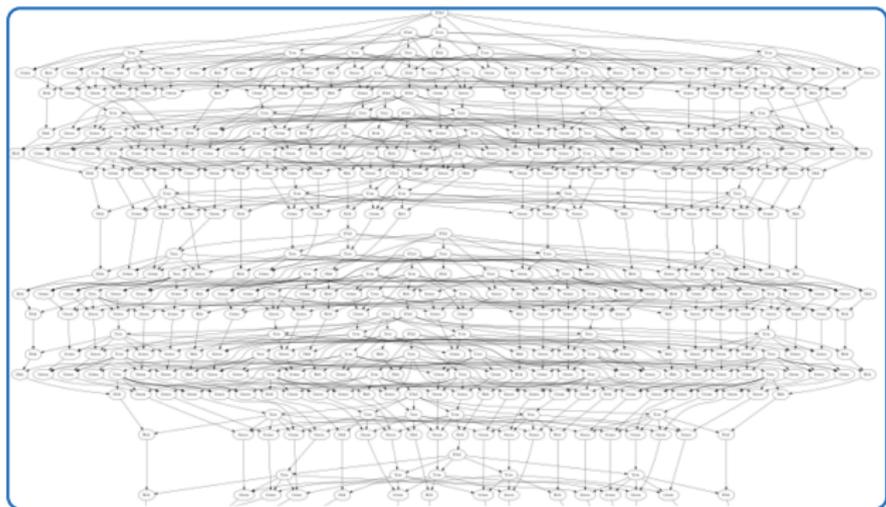
- `CrsMatrixBase` contains sparse data arrays *i.e.*, row pointers, column indices, value array.
- `MatrixView` defines a rectangular region *i.e.*, view offsets and dimensions;
 - light-weight object with meta data only.

```
typedef CrsMatrixBase<value_type := scalar> ScalarMatrix;
typedef MatrixView<base_object := ScalarMatrix> ViewOnScalarMatrix;

typedef CrsMatrixBase<value_type := ViewOnScalarMatrix> BlockMatrix;
typedef MatrixView<base_object := BlockMatrix> ViewOnBlockMatrix;
```


Example from Real Problem: pwtk

- Entire task DAG is constructed to demonstrate the degree of concurrency.
- Explicit DAG is never formed and not used in task scheduling in both Pthreads and Qthreads task polices.



Test problems from UFL sparse collection

Matrix ID	# of rows(n)	# of nonzeros (nnz)	nnz/n	Description
ecology2	999,999	4,995,991	4.99	Circuit theory
pwtk	217,918	11,524,432	52.88	Stiffness matrix

Machine specifications

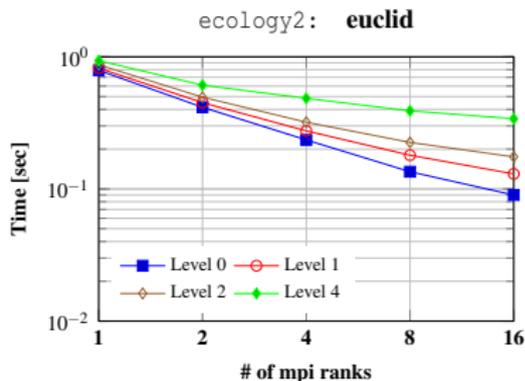
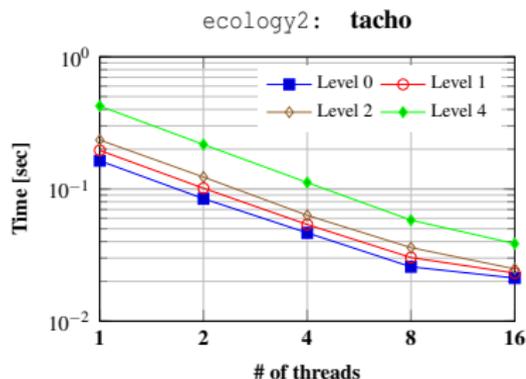
Processors	Intel Xeon E5-2670	Intel Xeon Phi	IBM Power8
# of cores	2x8	1x57	4x5
Clock speed	2.6 GHz	1.1 GHz	3.4 GHz
L2 per core	256 KB	512 KB	512 KB
L3	20 MB shared	-	8 MB per core
Compiler	Intel 15.2.164		GNU 4.9.2

Kokkos

- Pthreads backend with task only interface (team size = 1).

Comparison with Euclid (Hypre)

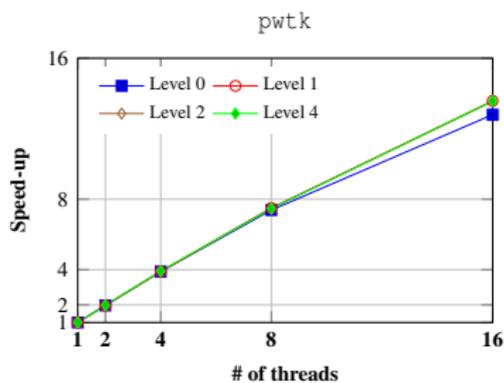
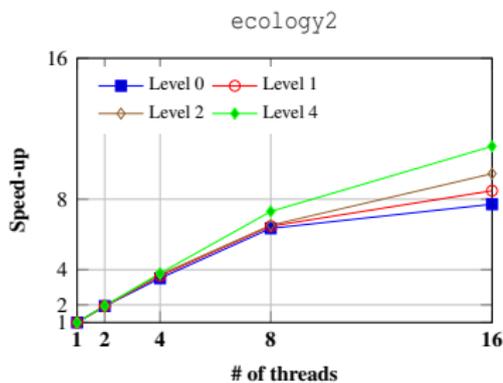
- Euclid performs **MPI-parallel** incomplete LU.
- Parallelism is extracted from **1D rowwise** partition of a matrix.
- **Reverse Cuthill McKee (RCM)** ordering is used to reduce the bandwidth of the matrix.
- Bandwidth of matrices increases with an increasing level of fills.



Intel Sandy Bridge

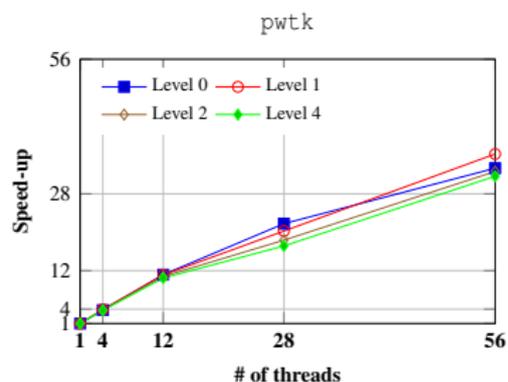
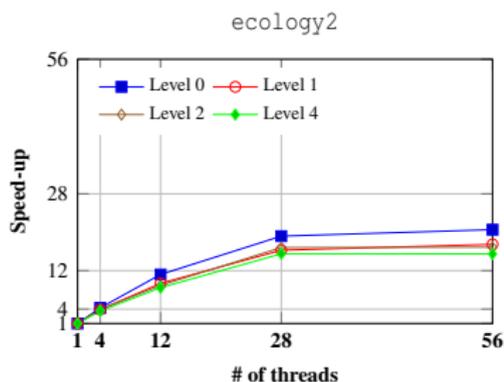
Strong Scaling : Intel Sandy Bridge

- Speed-up = $\frac{\text{Time for single-threaded Cholesky-by-blocks}}{\text{Time for parallel Cholesky-by-blocks}}$.
- Performance depends on matrix sparsity: `ecology2` is **sparser** and `pwtk` is **denser**.
- With increasing fill-level, factorization is **more compute-intensive**.
- Tasking overhead is **constant** and amortized during asynchronous parallel execution.

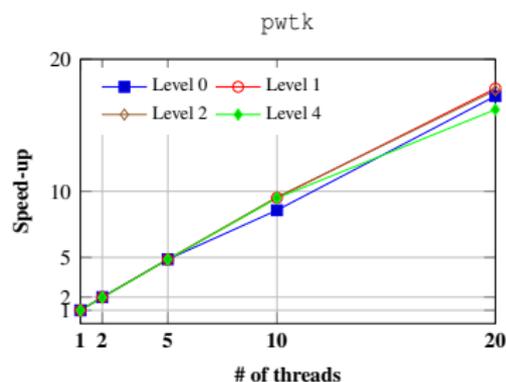
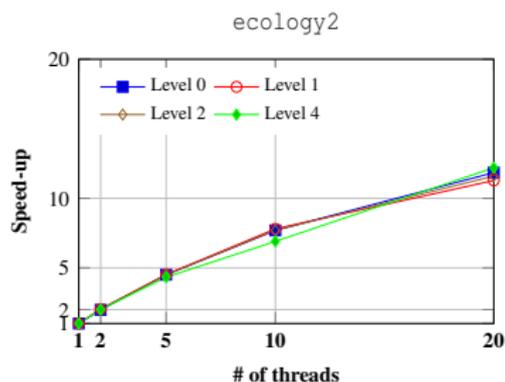


Strong Scaling : Intel Xeon Phi

- Speed-up = $\frac{\text{Time for single-threaded Cholesky-by-blocks}}{\text{Time for parallel Cholesky-by-blocks}}$.
- Performance depends on matrix sparsity: `ecology2` is **sparser** and `pwtk` is **denser**.
- With increasing fill-level, factorization is **more compute-intensive**.
- Tasking overhead is **constant** and amortized during asynchronous parallel execution.

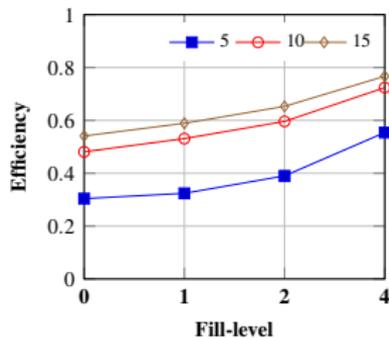
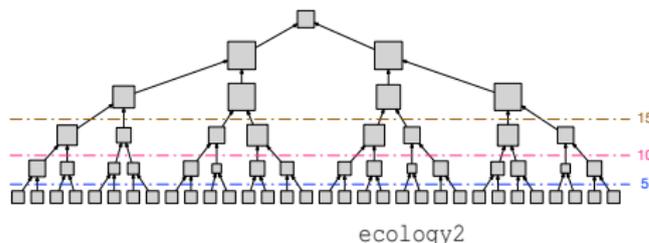


- Speed-up = $\frac{\text{Time for single-threaded Cholesky-by-blocks}}{\text{Time for parallel Cholesky-by-blocks}}$.
- Performance depends on matrix sparsity: `ecology2` is **sparser** and `pwtk` is **denser**.
- With increasing fill-level, factorization is **more compute-intensive**.
- Tasking overhead is **constant** and amortized during asynchronous parallel execution.

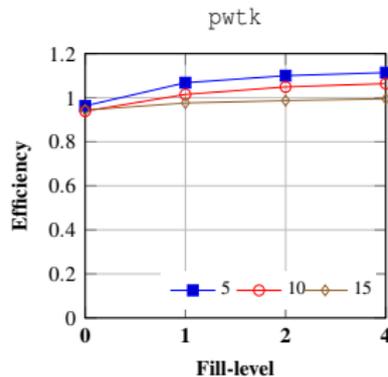


Tasking Overhead: Intel Xeon Phi

- Efficiency = $\frac{\text{Time for sequential Cholesky}}{\text{Time for single-threaded Cholesky-by-blocks}}$
- With increasing fill-level, cache-friendly sparse matrix operations on blocks exploit better data locality.
- Task granularity is **problem-specific** on irregular problems.



Treecut	# blocks	
	ecology2	pwtk
15	83	17
10	3,071	905
5	105,864	33,182



Kokkos hybrid task-data programming model

- Presented abstractions for task-data parallelism.
- Developed dependence driven task model.
- Harnessed two tasking backends: Pthreads and Qthreads.
- <https://github.com/kokkos/kokkos>

Task-data parallel sparse matrix factorization

- Presented sparse algorithm-by-blocks for task parallel Cholesky (in)complete factorization.
- As mini-app, provided support and feedback to design task-data interface.
- Demonstrated portable performance on multicore and manycore architectures.
- [Trilinos/shylu/tacho](#)

Kokkos tasking API and Cholesky miniapp are in the experimental phase

Kokkos hybrid task-data programming model

- Asynchronous tasking on GPUs.

Task-data parallel sparse matrix factorization

- Performance optimization for task-data hybrid parallelism: *e.g.*, algorithm design and thread team overhead.
- Supernodal direct factorization: Cholesky and LDL.
- Leverage to domain decomposition FE solver in collaboration with Clark Dohrmann (1542).