

Argonne Training Program on

EXTREME-SCALE COMPUTING



August 3 – August 15, 2014

Adaptive Linear Solvers and Eigensolvers

Jack Dongarra

University of Tennessee
Oak Ridge National Laboratory
University of Manchester

Dense Linear Algebra

Common Operations

$$Ax = b; \quad \min_x \|Ax - b\|; \quad Ax = \lambda x$$

- .. A major source of large dense linear systems is problems involving the solution of boundary integral equations.
 - The price one pays for replacing three dimensions with two is that what started as a sparse problem in $O(n^3)$ variables is replaced by a dense problem in $O(n^2)$.
- .. Dense systems of linear equations are found in numerous other applications, including:
 - airplane wing design;
 - radar cross-section studies;
 - flow around ships and other off-shore constructions;
 - diffusion of solid bodies in a liquid;
 - noise reduction; and
 - diffusion of light through small particles.

Existing Math Software - Dense LA

DIRECT SOLVERS	License	Support	Type		Language			Mode		
			Real	Complex	F77	C	C++	Shared	GPU	Dist
Eigen	Mozilla	yes	X	X			X	X		
Elemental	BSD	yes	X	X			X			M
FLAME	LGPL	yes	X	X	X	X		X		
FLENS	BSD	yes	X	X			X	X		
LAPACK	BSD	yes	X	X	X	X		X		
LAPACK95	BSD	yes	X	X	F95			X		
MAGMA	BSD	yes	X	X	X	X		X	C/O/X	
NAPACK	BSD	yes	X		X			X		
PLAPACK	?	no	X	X	X	X				M
PLASMA	BSD	yes	X	X	X	X		X		
PRISM	?	no	X		X			X		M
rejtrix	by-nc-sa	yes	X				X	X		
ScaLAPACK	BSD	yes	X	X	X	X				M/P
Trilinos/Pliris	BSD	yes	X	X		X	X			M
ViennaCL	MIT	yes	X				X	X	C/O/X	

<http://www.netlib.org/utk/people/JackDongarra/la-sw.html>

•• LINPACK, EISPACK, LAPACK, ScaLAPACK

➤ PLASMA, MAGMA



DLA Solvers

- We are interested in developing Dense Linear Algebra Solvers
- Retool LAPACK and ScaLAPACK for multicore and hybrid architectures

What do you mean by performance?

◆ What is a x flop/s?

- **x flop/s is a rate of execution, some number of floating point operations per second.**
 - » Whenever this term is used it will refer to 64 bit floating point operations and the operations will be either addition or multiplication.

◆ What is the theoretical peak performance?

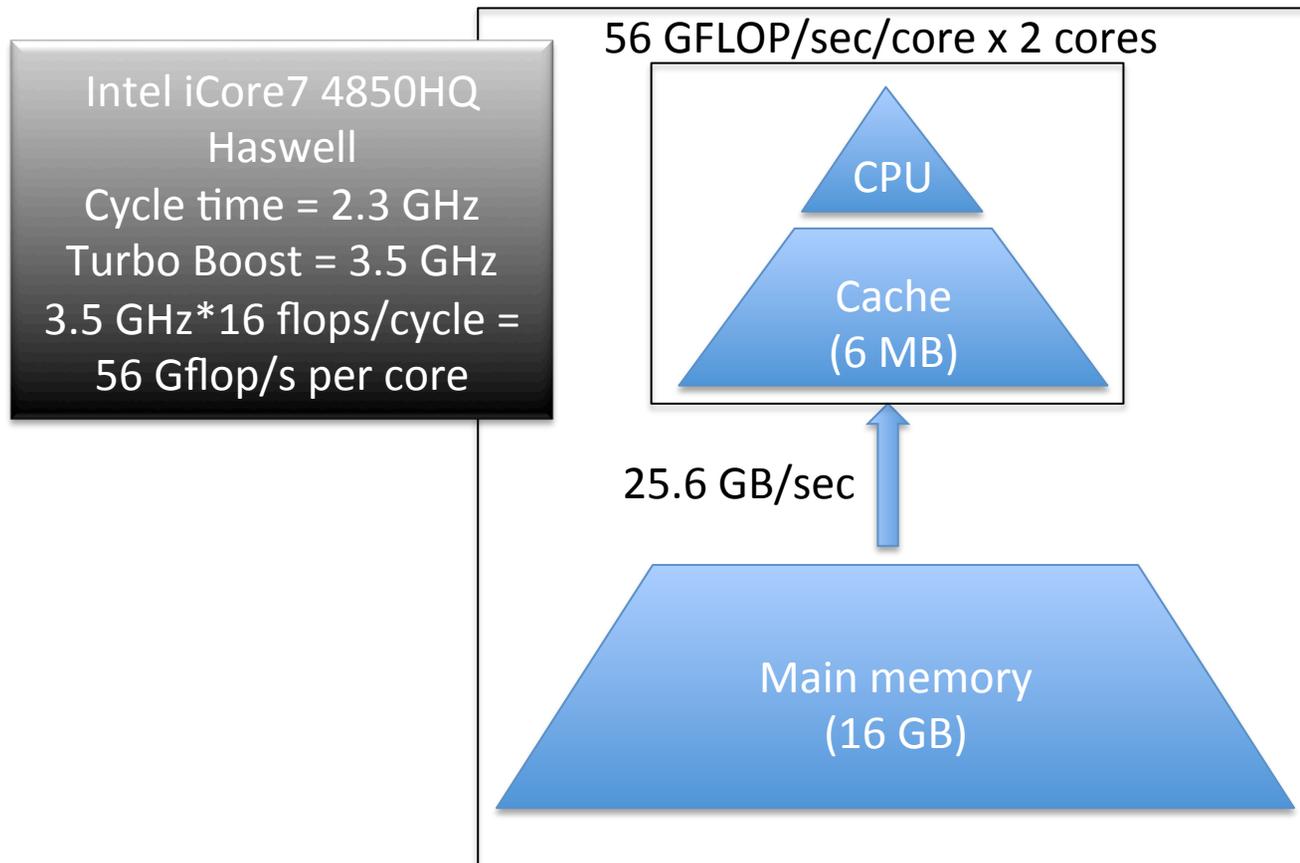
- **The theoretical peak is based not on an actual performance from a benchmark run, but on a paper computation to determine the theoretical peak rate of execution of floating point operations for the machine.**
- **The theoretical peak performance is determined by counting the number of floating-point additions and multiplications (in full precision) that can be completed during a period of time, usually the cycle time of the machine.**
- **For example, an Intel Xeon 5570 quad core at 2.93 GHz can complete 4 floating point operations per cycle or a theoretical peak performance of 11.72 GFlop/s per core or 46.88 Gflop/s for the socket.**

Computing Theoretical Peak

- **To get the number of operation per cycle per core, you need to look for special instruction set:**
 - **Most of the computer nowadays have FMA (Fused multiple add):**
(i.e. $x \leftarrow x + y * z$ in one cycle)
 - **Intel Pentium and AMD Opteron have SSE2**
 - 2 flops/cycle DP & 4 flops/cycle SP
 - **Intel Xeon Nehalem & Westmire have AVX**
 - 4 flops/cycle DP & 8 flops/cycle SP
 - **Intel Xeon Sandy Bridge & Ivy Bridge have AVX**
 - 8 flops/cycle DP & 16 flops/cycle SP
 - **Intel Xeon Haswell has AVX**
 - 16 flops/cycle DP & 32 flops/cycle SP
 - **IBM PowerPC has AltiVec**
 - 8 flops/cycle SP
 - 4 flops/cycle DP

Memory transfer

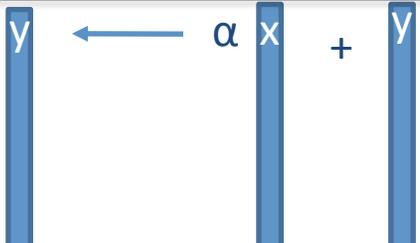
- One level of memory model on my laptop:



(Omitting latency here.)

The model IS simplified (see next slide) but it provides an upper bound on performance as well. I.e., we will never go faster than what the model predicts. (And, of course, we can go slower ...)

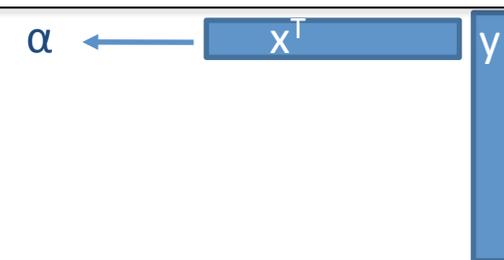
FMA: fused multiply-add

AXPY:  $\alpha x + y$

```
for ( j = 0; j < n; j++)  
    y[i] += a * x[i];
```

(without increment)

n MUL
n ADD
2n FLOP
n FMA

DOT:  $\alpha \leftarrow x^T y$

```
alpha = 0e+00;  
for ( j = 0; j < n; j++)  
    alpha += x[i] * y[i];
```

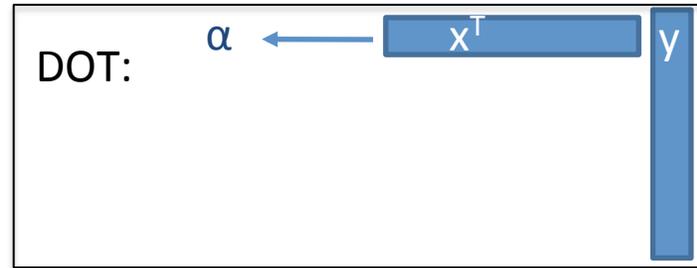
(without increment)

n MUL
n ADD
2n FLOP
n FMA

Note: It is reasonable to expect the one loop codes shown here to perform as well as their Level 1 BLAS counterpart (on multicore with an OpenMP pragma for example).

The true gain these days with using the BLAS is (1) Level 3 BLAS, and (2) portability.

- Take two double precision vectors x and y of size $n=375,000$.



- Data size:
 - (375,000 double) * (8 Bytes / double) = 3 MBytes per vector
 - (Two vectors fit in cache (6 MBytes). OK.)
- Time to move the vectors in cache:
 - (6 MBytes) / (25.6 GBytes/sec) = **0.23 ms**
- Time to perform computation of AXPY:
 - (2n flop) / (56 Gflop/sec) = **0.01 ms**

$$\begin{aligned} \text{total_time} &\geq \text{time_comm} + \text{time_comp} \\ &= 0.23\text{ms} + 0.01\text{ms} = 0.24\text{ms} \end{aligned}$$

1. The equation assumes that communication and computation do not overlap. This assumption is hardly true. For a AXPY (without INCX), a nice streaming/pipeline occurs. Data locality is perfect.
2. We will assume overlapping communication and computation. We are looking for a strict lower bound, and there is no need to second guessing and make assumptions.

$$\begin{aligned} \text{total_time} &\geq \max (\text{time_comm} , \text{time_comp}) \\ &= \max(0.23\text{ms} , 0.01\text{ms}) = 0.23\text{ms} \end{aligned}$$

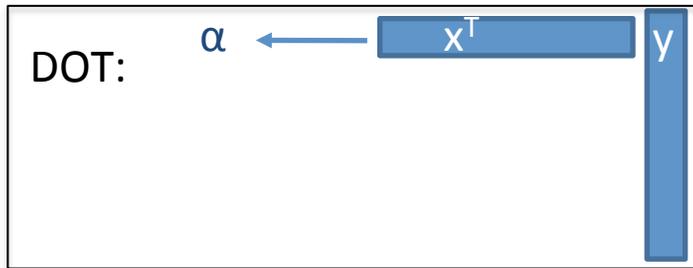
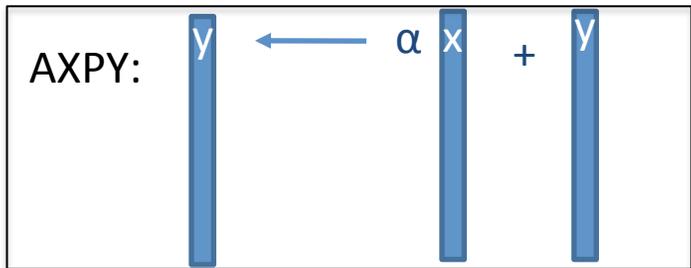
performance for DOT ≤ 3.2 Gflop/s

We say that the performance operation is communication bounded.

Peak is 56 Gflop/s

Level 1, 2 and 3 BLAS

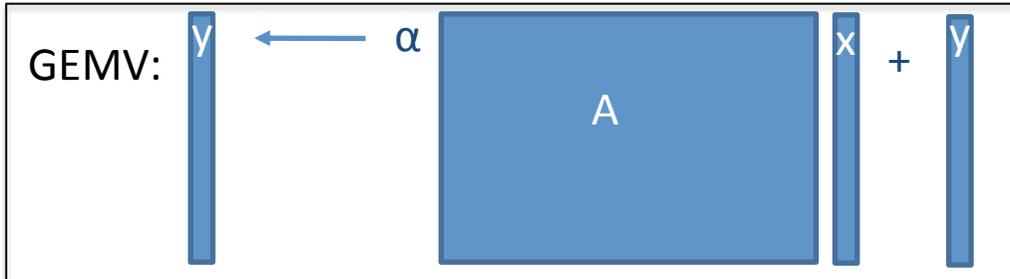
Level 1 BLAS Matrix-Vector operations



2n FLOP
2n memory reference
AXPY: 2n READ, n WRITE
DOT: 2n READ

RATIO: 1

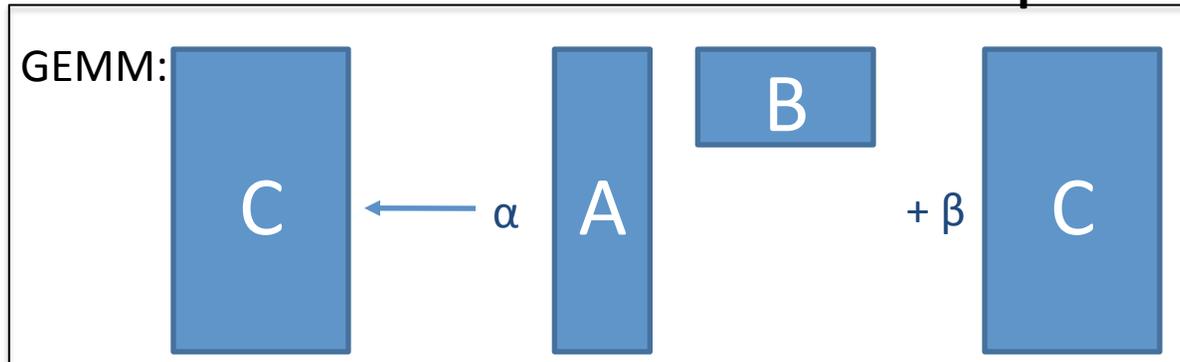
Level 2 BLAS Matrix-Vector operations



2n² FLOP
n² memory references

RATIO: 2

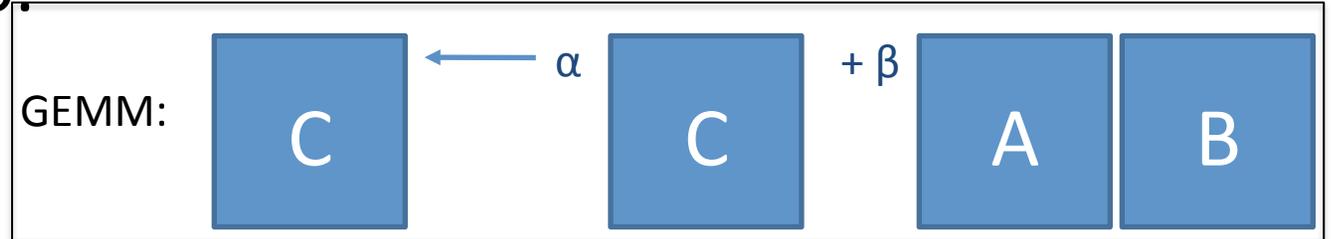
Level 3 BLAS Matrix-Matrix operations



2n³ FLOP
3n² memory references
3n² READ, n² WRITE

RATIO: 2/3 n

- Take three double precision matrices A, B, and C of size $n=500$.



- Data size:
 - $(500^2 \text{ double}) * (8 \text{ Bytes / double}) = 2 \text{ MBytes per matrix}$
 - (Three matrices fit in cache (6 MBytes). OK.)
- Time to move the matrices in cache:
 - $(6 \text{ MBytes}) / (25.6 \text{ GBytes/sec}) = \mathbf{0.23 \text{ ms}}$
- Time to perform computation of GEMM:
 - $(2n^3 \text{ flop}) / (56 \text{ Gflop/sec}) = \mathbf{4.46 \text{ ms}}$

$$\begin{aligned} \text{total_time} &\geq \max (\text{time_comm} , \text{time_comp}) \\ &= \max(0.23\text{ms} , 4.46\text{ms}) = 4.46\text{ms} \end{aligned}$$

For this example, communication time is less than 6% of the computation time.

There is a lots of data reuse in a GEMM – $2/3n$ ops per data element. We have good temporal locality.

If we assume $\text{total_time} \approx \text{time_comm} + \text{time_comp}$, we get
performance for GEMM ≈ 53.3 Gflop/sec

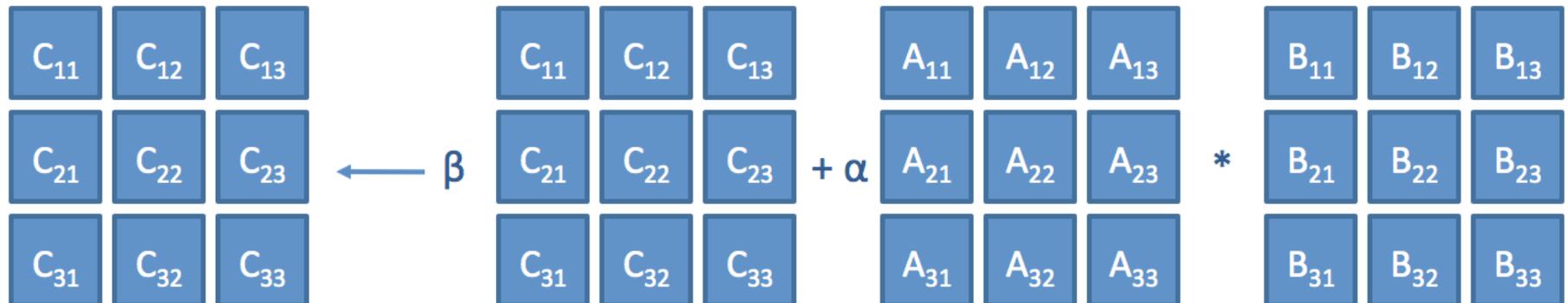
(Out of 56 Gflop/sec possible, so that would be 95% peak performance efficiency.)

Issues

- Reuse based on matrices that fit into cache.
- What if you have matrices bigger than cache?

Issues

- Reuse based on matrices that fit into cache.
- What if you have matrices bigger than cache?
- Break matrices into blocks or tiles that will fit.



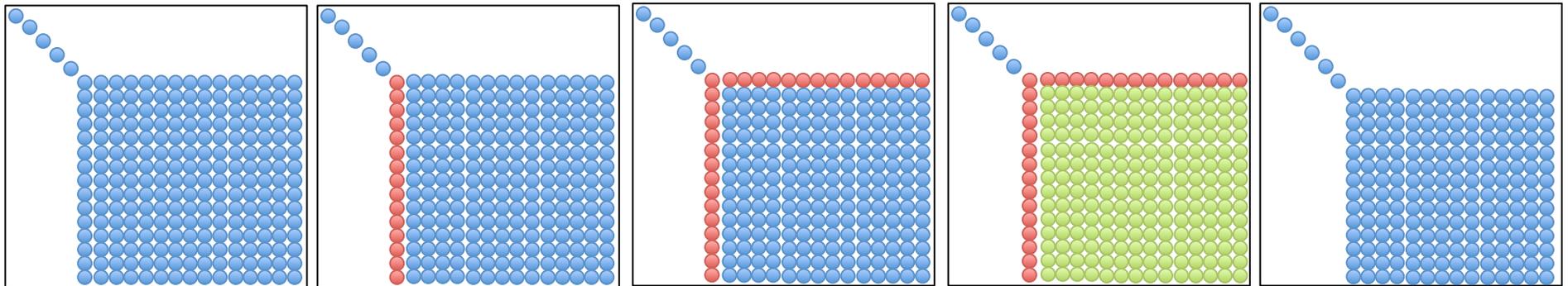
By the way

Performance for your laptop

- If you are interested in running the Linpack Benchmark on your system see:
<https://software.intel.com/en-us/node/157667?wapkw=mkl+linpack>
- Also Intel has a power meter, see:
<https://software.intel.com/en-us/articles/intel-power-gadget-20>

The Standard LU Factorization LINPACK

1970's HPC of the Day: Vector Architecture



Factor column
with Level 1
BLAS

Divide by
Pivot
row

Schur
complement
update
(Rank 1 update)

Next Step

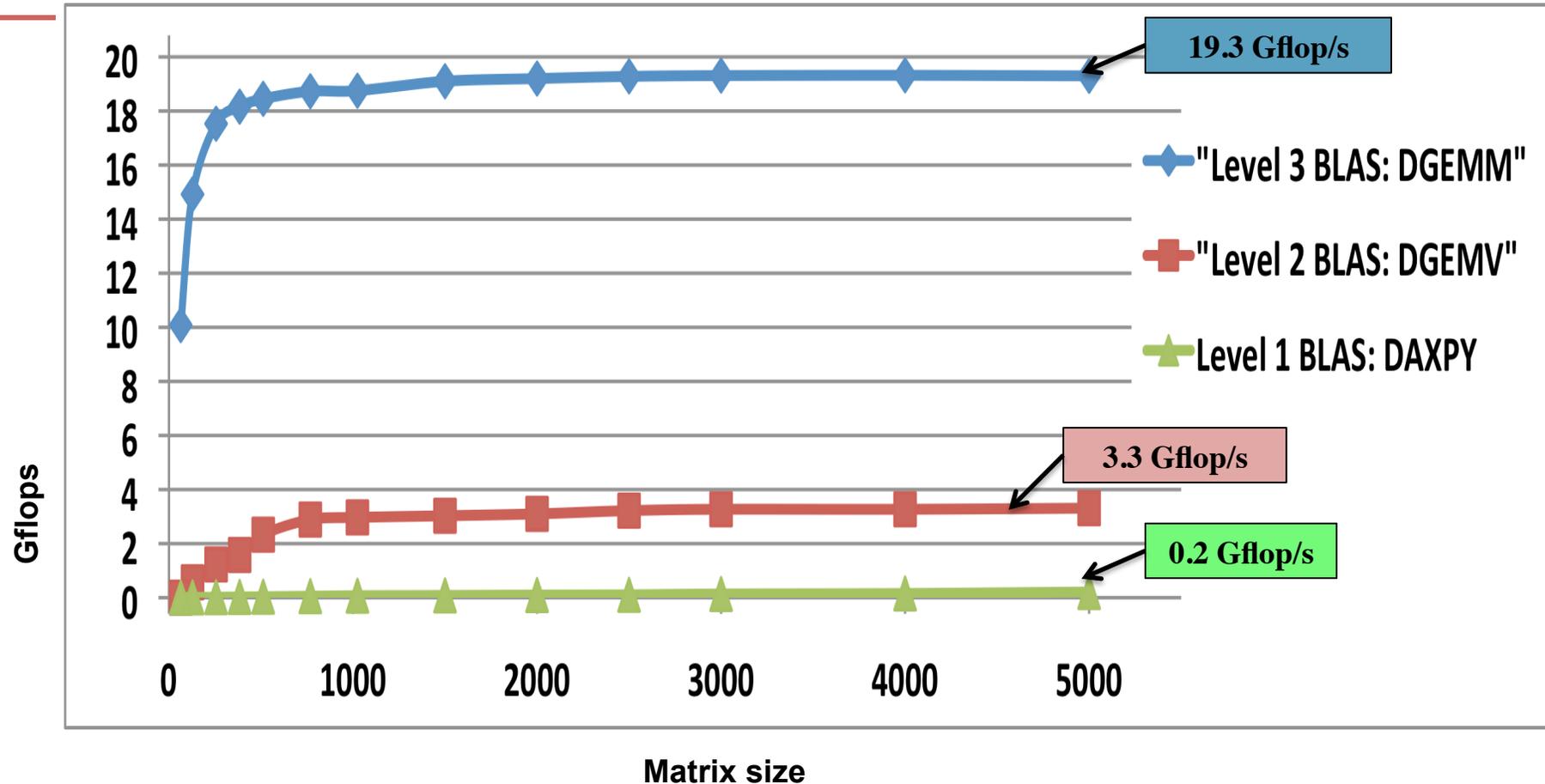
Main points

- Factorization column (zero) mostly sequential due to memory bottleneck
- Level 1 BLAS
- Divide pivot row has little parallelism
- Rank -1 Schur complement update is the only easy parallelize task
- Partial pivoting complicates things even further
- Bulk synchronous parallelism (fork-join)
 - Load imbalance
 - Non-trivial Amdahl fraction in the panel
 - Potential workaround (look-ahead) has complicated implementation



Level 1, 2 and 3 BLAS

1 core Intel Xeon E5-2670 (Sandy Bridge); 2.6 GHz; Peak = 20.8 Gflop/s



1 core Intel Xeon E5-2670 (Sandy Bridge), 2.6 GHz.

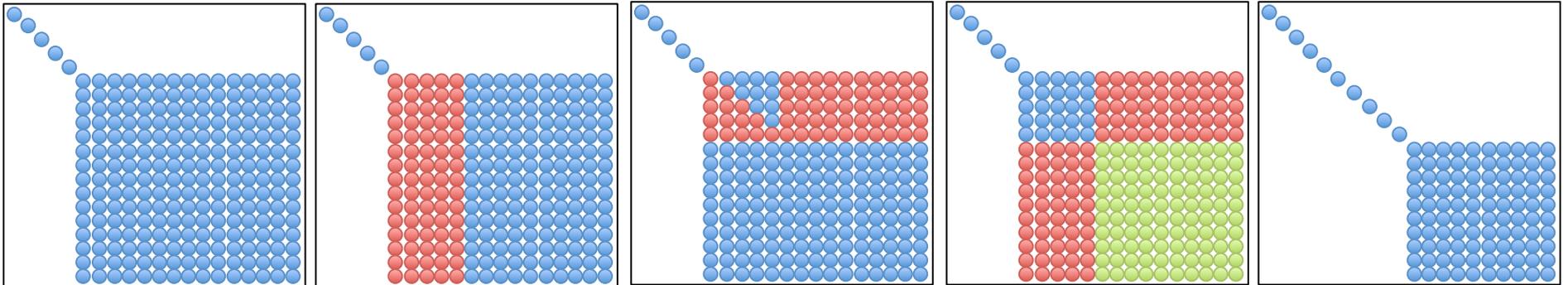
24 MB shared L3 cache, and each core has a private 256 KB L2 and 64 KB L1.

The theoretical peak per core DP is $8 \text{ flop/cycle} * 2.6 \text{ GHz} = 20.8 \text{ Gflop/s}$ per core.

Compiled with gcc 4.4.6 and using MKL_composer_xe_2013.3.163

The Standard LU Factorization LAPACK

1980's HPC of the Day: Cache Based SMP



Factor panel
with Level 1,2
BLAS

Triangular
update

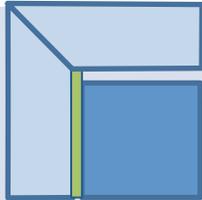
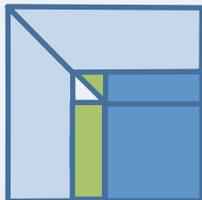
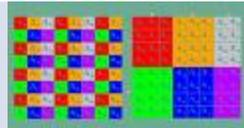
Schur
complement
update

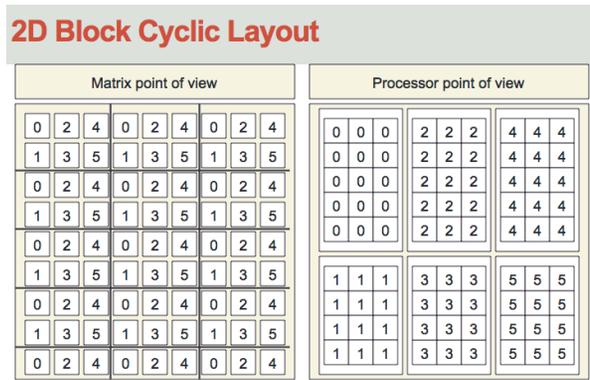
Next Step

Main points

- Panel factorization mostly sequential due to memory bottleneck
- Triangular solve has little parallelism
- Schur complement update is the only easy parallelize task
- Partial pivoting complicates things even further
- Bulk synchronous parallelism (fork-join)
 - Load imbalance
 - Non-trivial Amdahl fraction in the panel
 - Potential workaround (look-ahead) has complicated implementation

Last Generations of DLA Software

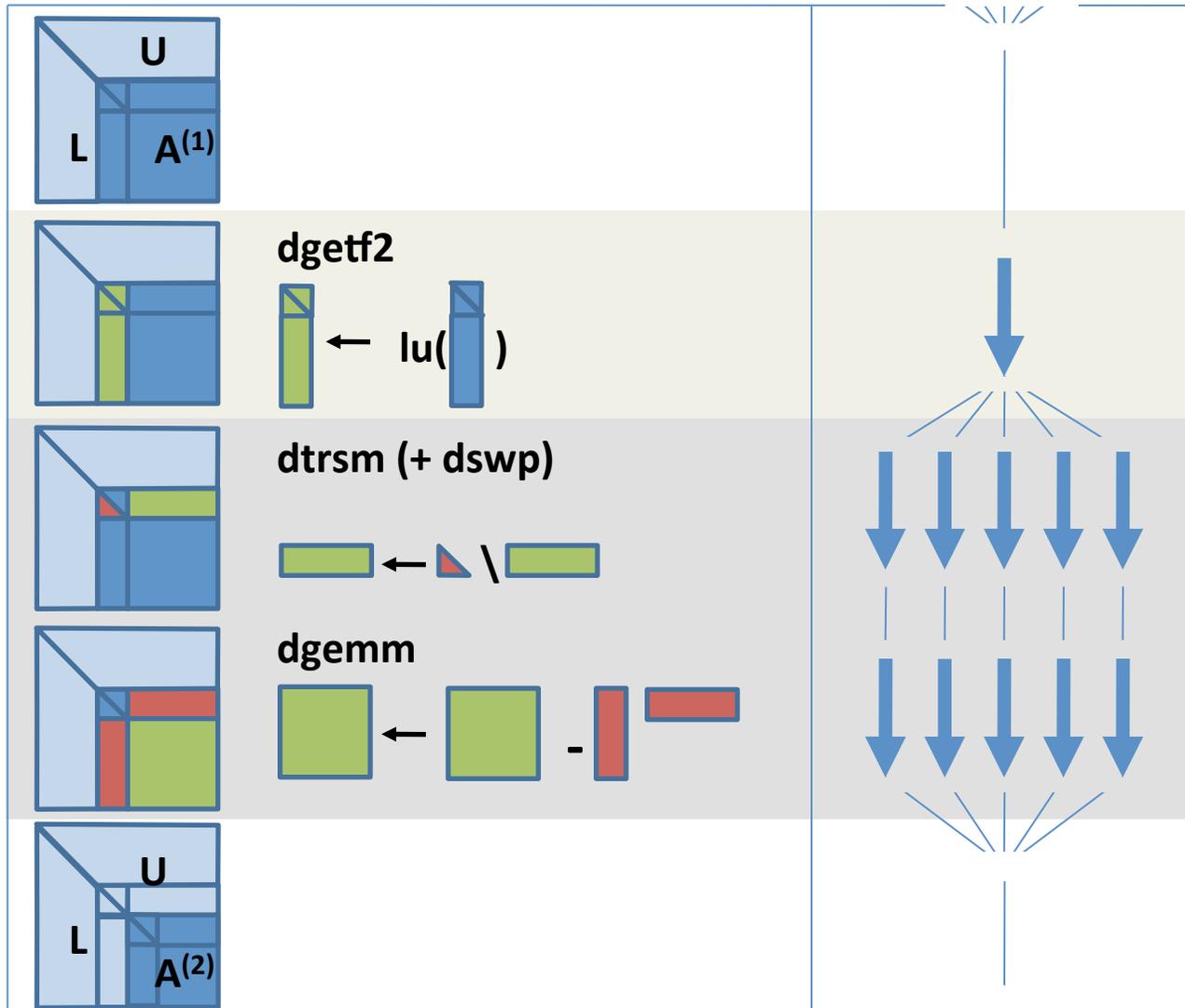
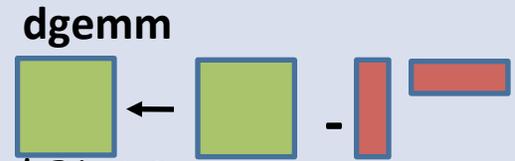
Software/Algorithms follow hardware evolution in time		
LINPACK (70's) (Vector operations)		Rely on - Level-1 BLAS operations
LAPACK (80's) (Blocking, cache friendly)		Rely on - Level-3 BLAS operations
ScaLAPACK (90's) (Distributed Memory)		Rely on - PBLAS Mess Passing



Parallelization of LU and QR.

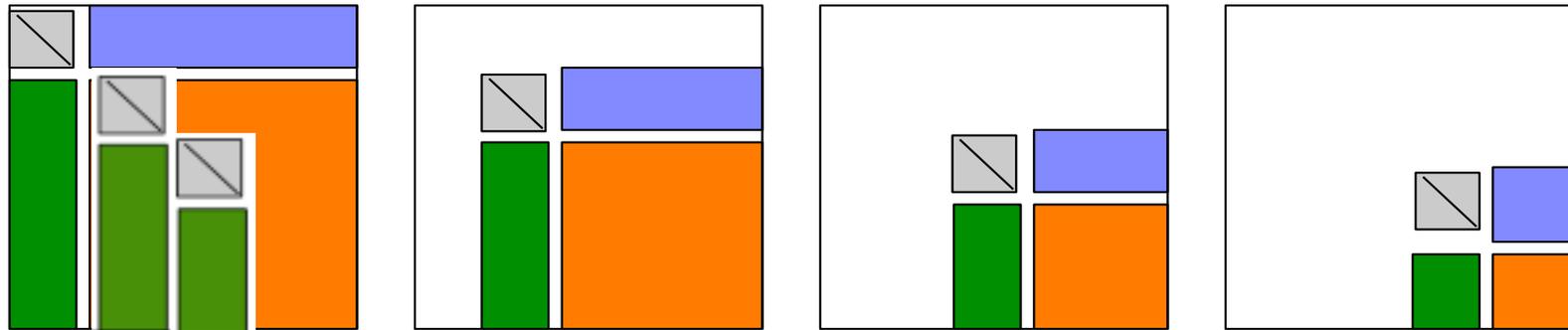
Parallelize the update:

- Easy and done in any reasonable software.
- This is the $2/3n^3$ term in the FLOPs count.
- Can be done efficiently with LAPACK+multithreaded BLAS

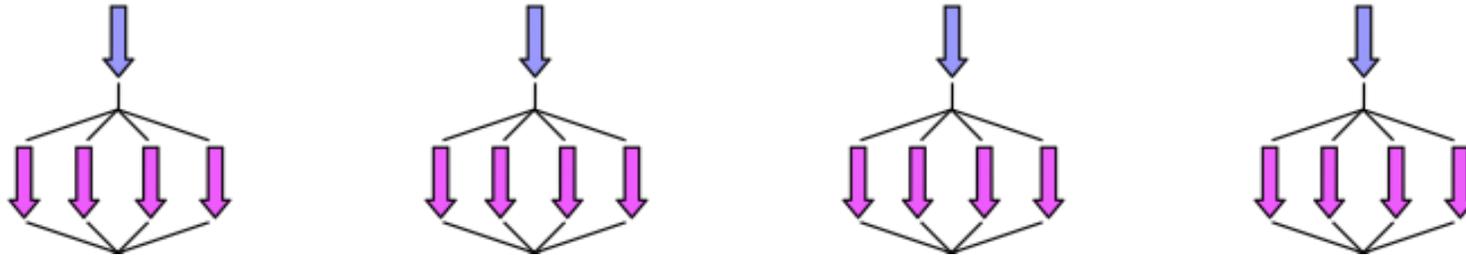


Fork - Join parallelism
 Bulk Sync Processing

Synchronization (in LAPACK LU)



Step 1 → Step 2 → Step 3 → Step 4 ...



DGETF2
(Factor a panel)



DLSWP
(Backward swap)



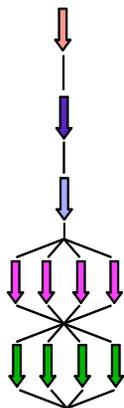
DLSWP
(Forward swap)



DTRSM
(Triangular solve)



DGEMM
(Matrix multiply)



LAPACK

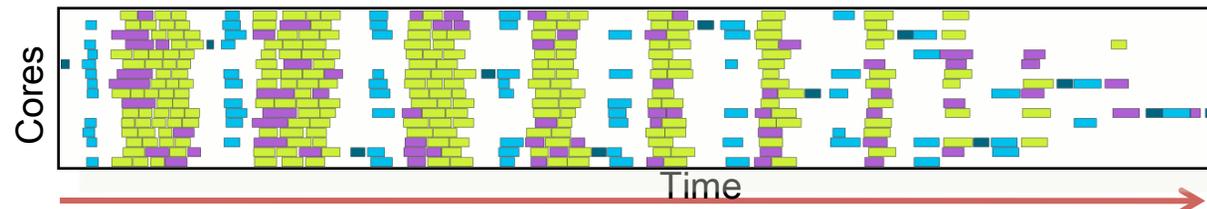
LAPACK

LAPACK

BLAS

BLAS

- fork join
- bulk synchronous processing

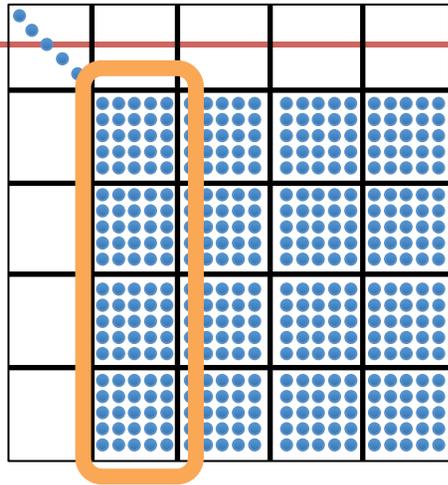




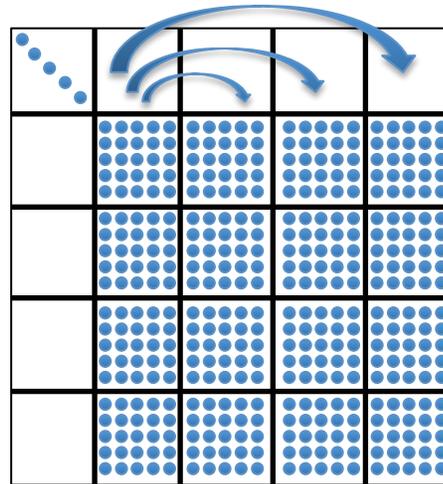
PLASMA LU Factorization

Dataflow Driven

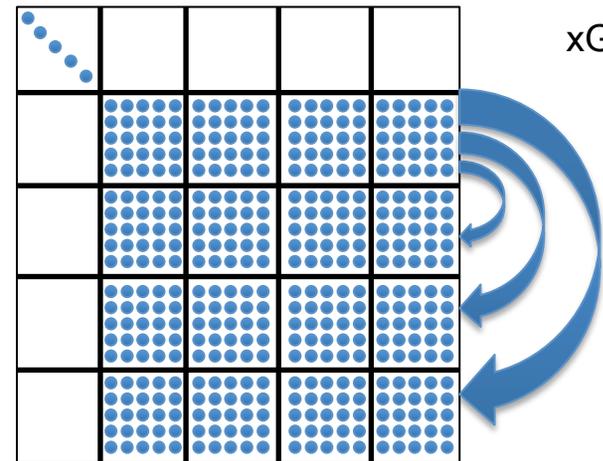
Numerical program generates tasks and run time system executes tasks respecting data dependences.



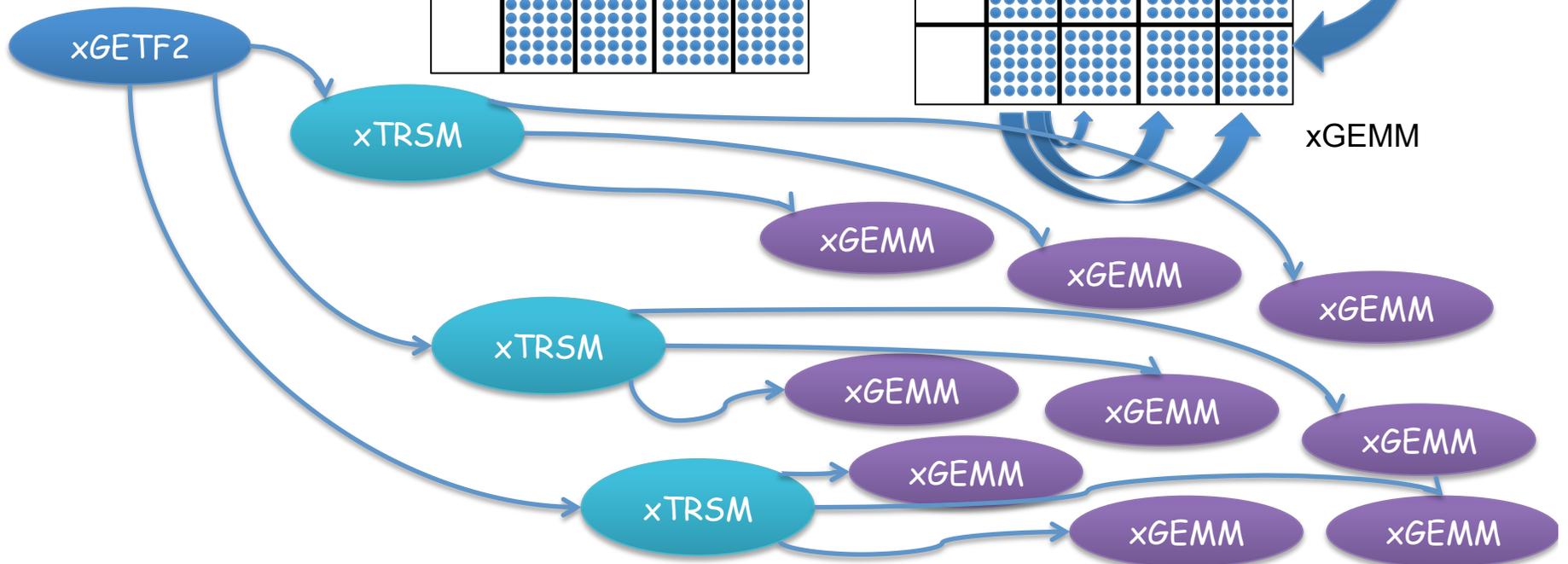
xTRSM



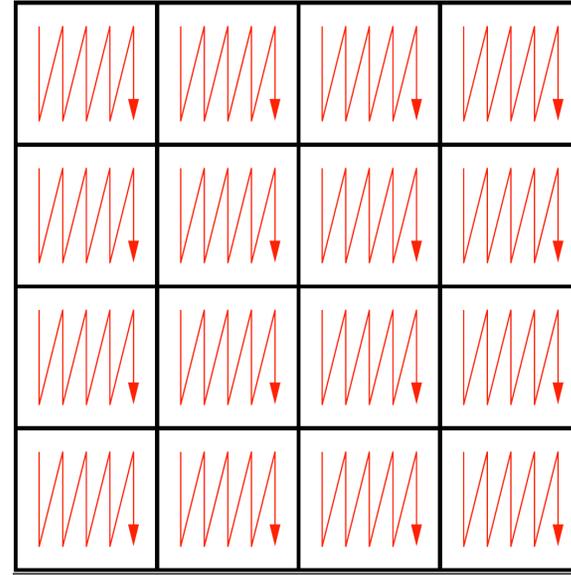
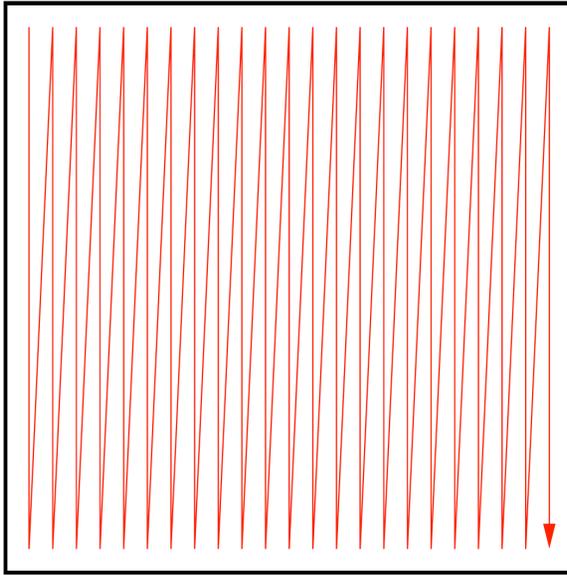
xGEMM



xGEMM



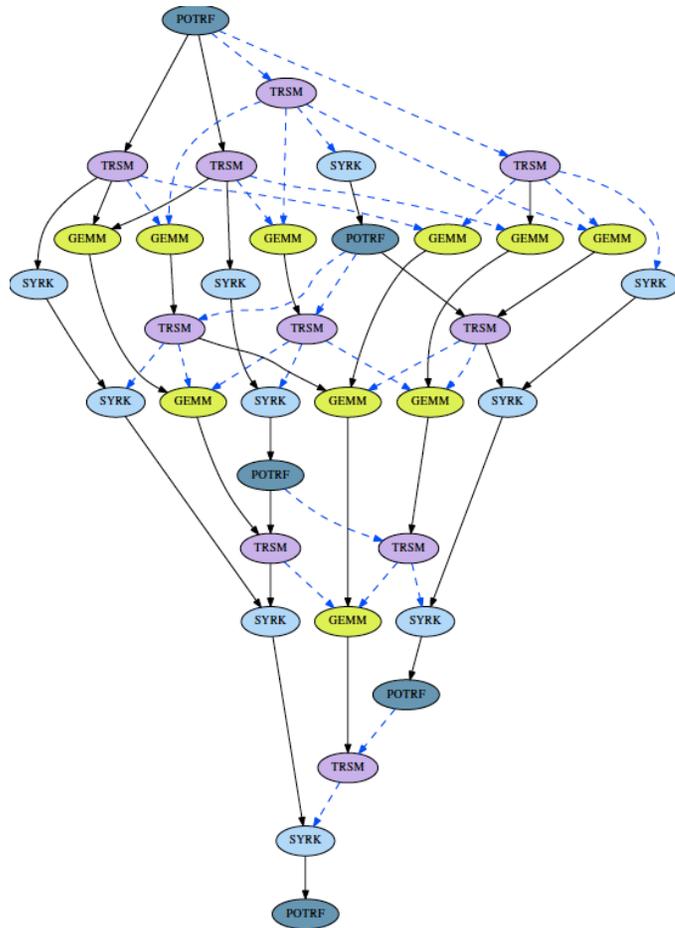
Data Layout is Critical



- Tile data layout where each data tile is contiguous in memory
- Decomposed into several fine-grained tasks, which better fit the memory of the small core caches

PLASMA LU: Tile Algorithm and Nested Parallelism

- Operates on one, two, or three matrix tiles at a time using a single core
 - This is called a kernel; executed independently of other kernels
 - Mostly Level 3 BLAS are used
- Data flows between kernels as prescribed by the programmer
- Coordination is done transparently via runtime scheduler (QUARK)
 - Parallelism level adjusted at runtime
 - Look-ahead adjusted at runtime
- Uses single-threaded BLAS with all the optimization benefits
- Panel is done on multiple cores
 - Recursive formulation of LU for better BLAS use
 - Level 1 BLAS are faster because they work on combined cache size



A runtime environment for the dynamic execution of precedence-constraint tasks (DAGs) in a multicore machine

- Translation
- If you have a serial program that consists of computational kernels (tasks) that are related by data dependencies, QUARK can help you execute that program (relatively efficiently and easily) in parallel on a multicore machine

```
FOR k = 0..TILES-1  
  A[k][k] ← DPOTRF(A[k][k])  
  FOR m = k+1..TILES-1  
    A[m][k] ← DTRSM(A[k][k], A[m][k])  
  FOR m = k+1..TILES-1  
    A[m][m] ← DSYRK(A[m][k], A[m][m])  
    FOR n = k+1..m-1  
      A[m][n] ← DGEMM(A[m][k], A[n][k], A[m][n])
```

definition – pseudocode

The Purpose of a QUARK Runtime

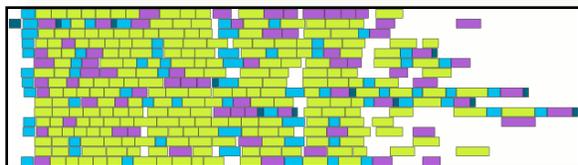
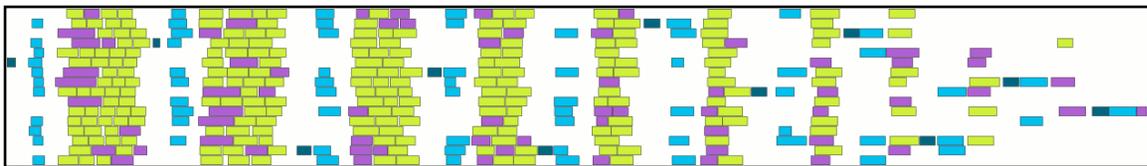
Objectives

- High utilization of each core
- Scaling to large number of cores
- Synchronization reducing algorithms

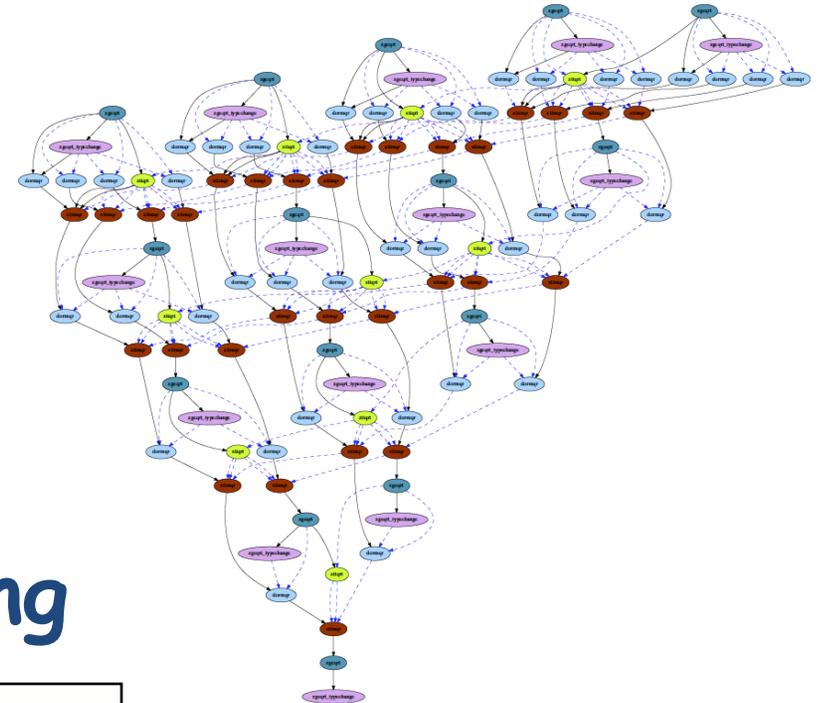
Methodology

- Dynamic DAG scheduling (QUARK)
- Explicit parallelism
- Implicit communication
- Fine granularity / block data layout

Arbitrary DAG with dynamic scheduling



DAG scheduled parallelism

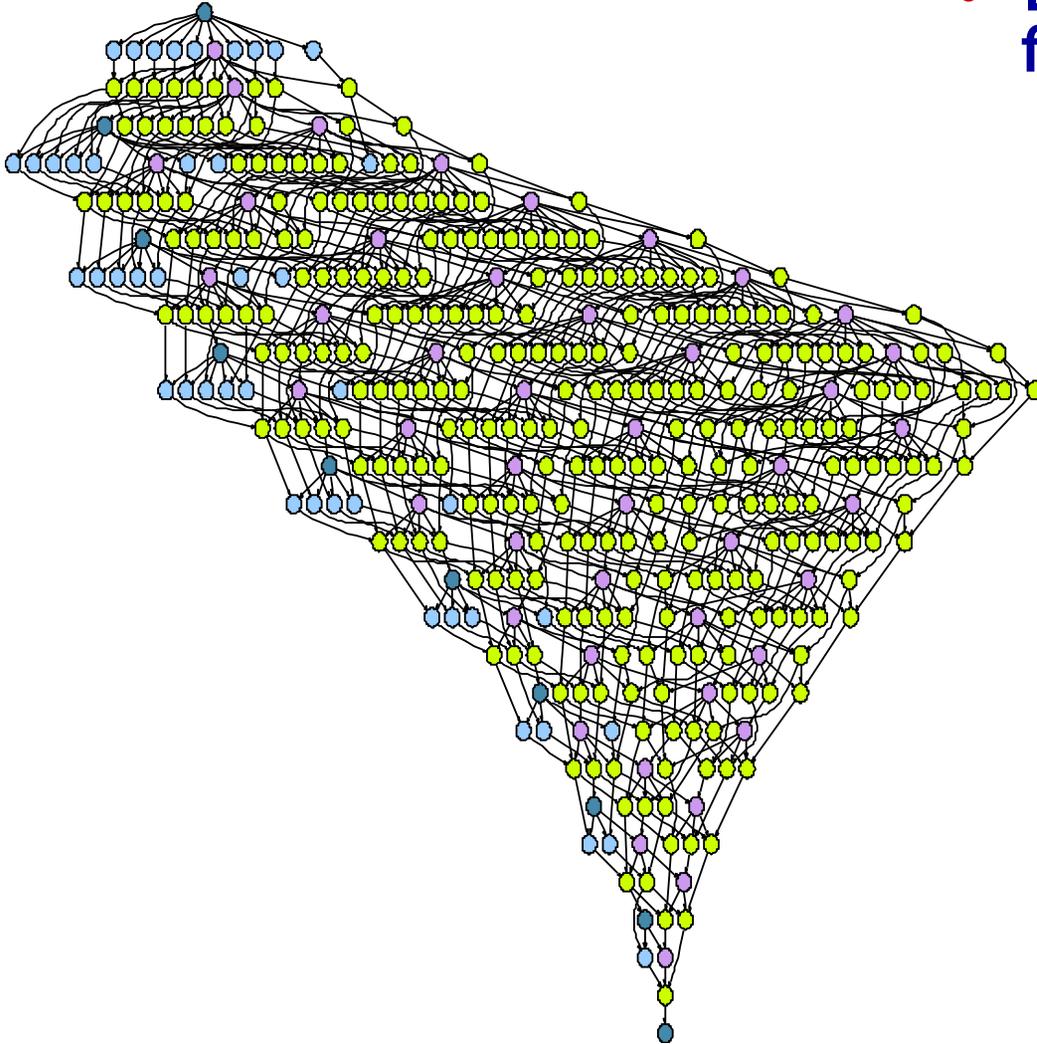


Fork-join parallelism
Notice the synchronization penalty in the presence of heterogeneity.



PLASMA Local Scheduling

Dynamic Scheduling: Sliding Window

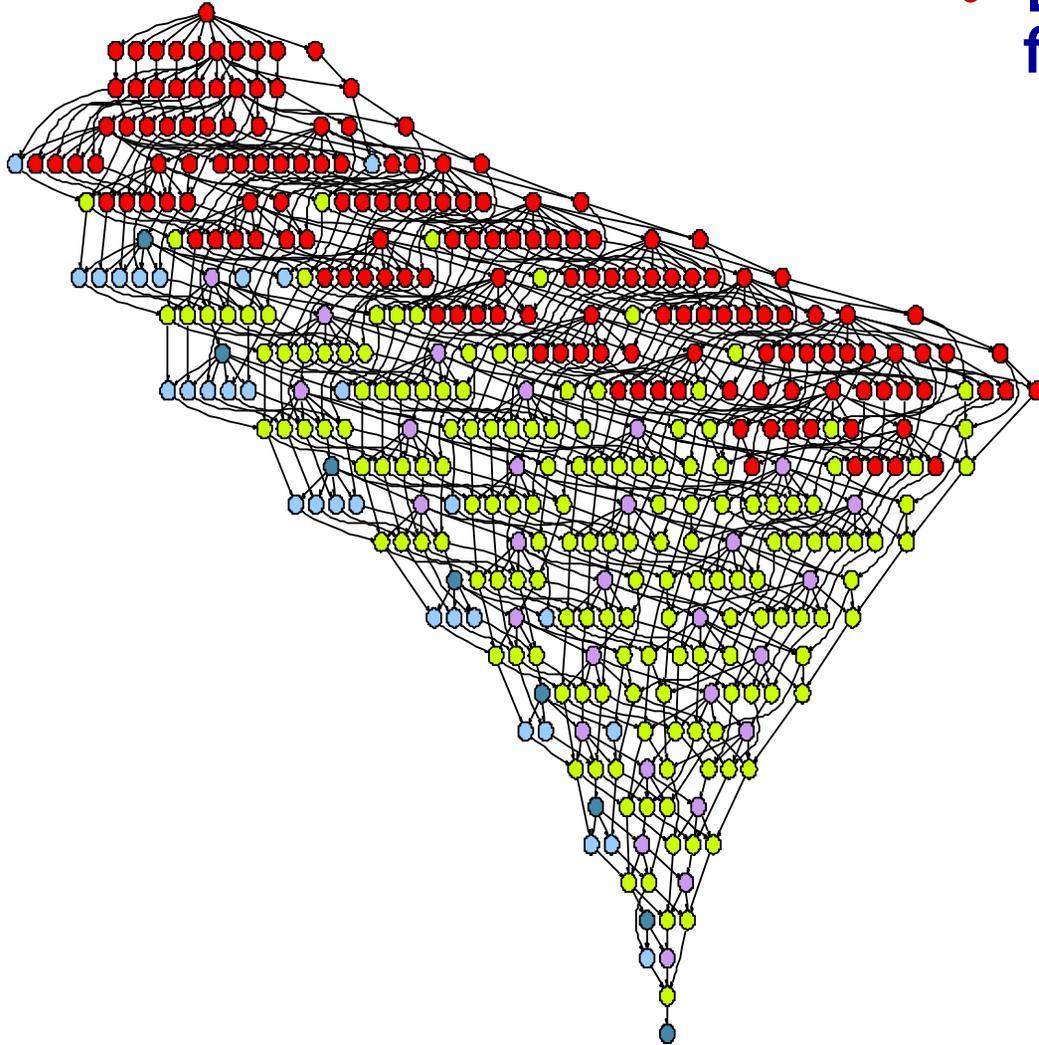


- **DAGs get very big, very fast**
 - **So windows of active tasks are used; this means no global critical path**
 - **Matrix of $NB \times NB$ tiles; NB^3 operation**
 - **$NB=100$ gives 1 million tasks**



PLASMA Local Scheduling

Dynamic Scheduling: Sliding Window

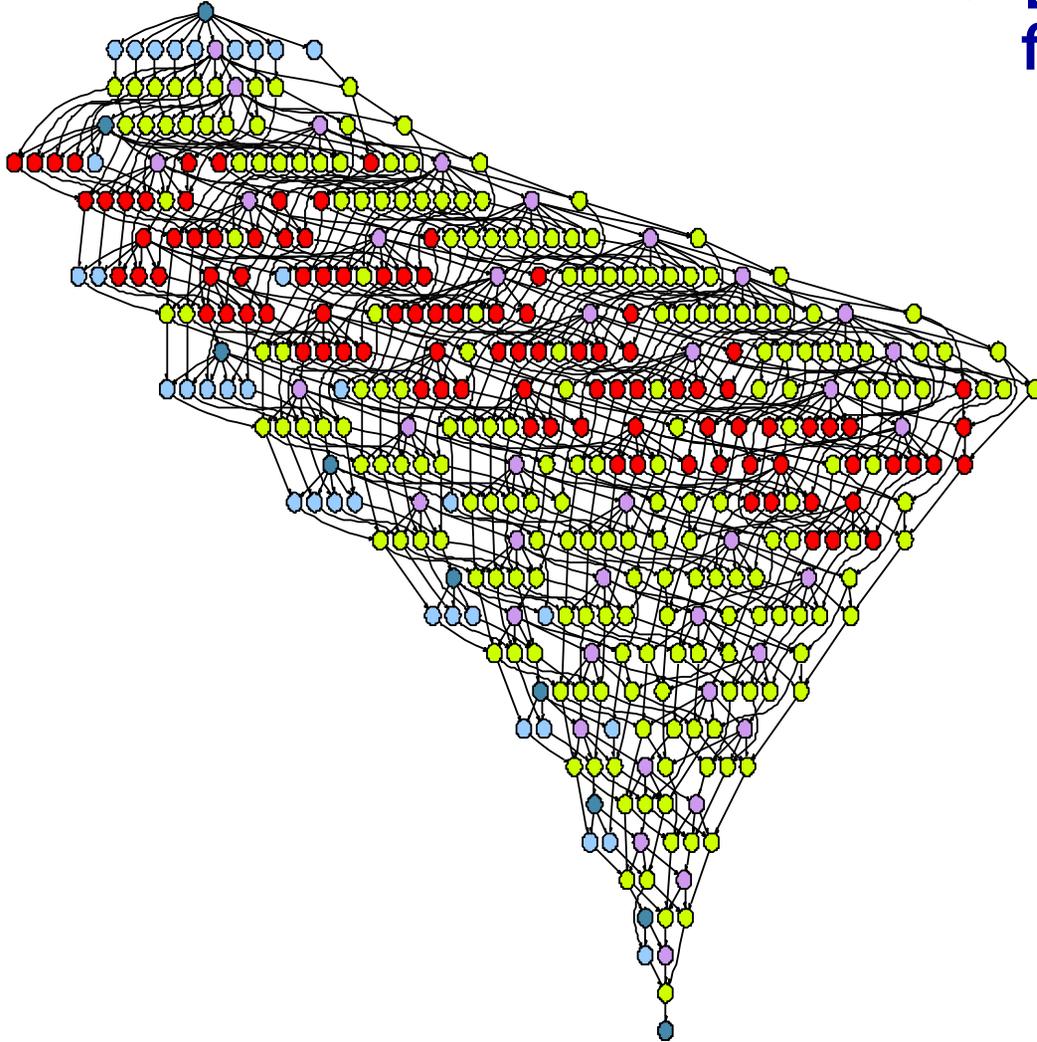


- **DAGs get very big, very fast**
 - So windows of active tasks are used; this means no global critical path
 - Matrix of $NB \times NB$ tiles; NB^3 operation
 - $NB=100$ gives 1 million tasks



PLASMA Local Scheduling

Dynamic Scheduling: Sliding Window

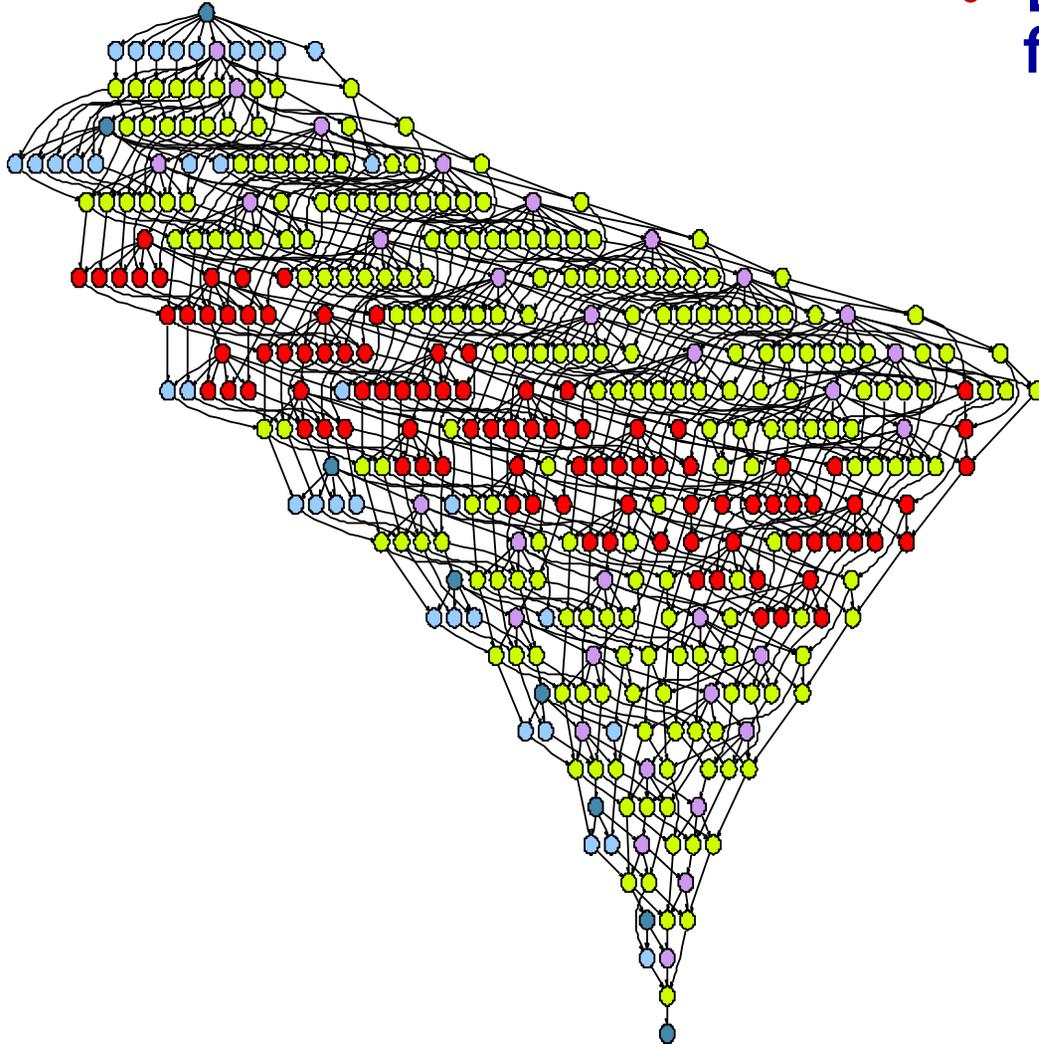


- **DAGs get very big, very fast**
 - So windows of active tasks are used; this means no global critical path
 - Matrix of $NB \times NB$ tiles; NB^3 operation
 - $NB=100$ gives 1 million tasks



PLASMA Local Scheduling

Dynamic Scheduling: Sliding Window



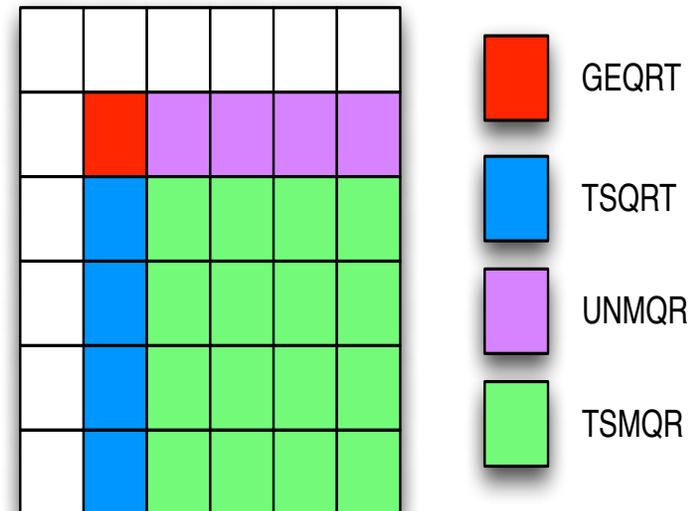
- **DAGs get very big, very fast**
 - So windows of active tasks are used; this means no global critical path
 - Matrix of $NB \times NB$ tiles; NB^3 operation
 - $NB=100$ gives 1 million tasks

Example: QR Factorization

```

FOR k = 0 .. SIZE - 1
  A[k][k], T[k][k] <- GEQRT( A[k][k] )
  FOR m = k+1 .. SIZE - 1
    A[k][k]|Up, A[m][k], T[m][k] <-
      TSQRT( A[k][k]|Up, A[m][k], T[m][k] )
    FOR n = k+1 .. SIZE - 1
      A[k][n] <- UNMQR( A[k][k]|Low, T[k][k], A[k][n] )
      FOR m = k+1 .. SIZE - 1
        A[k][n], A[m][n] <-
          TSMQR( A[m][k], T[m][k], A[k][n], A[m][n] )

```

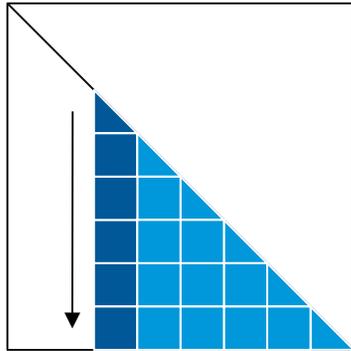




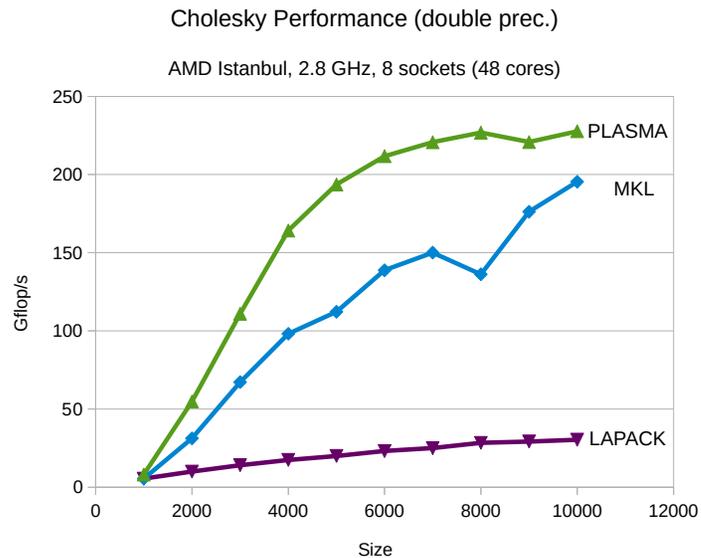
Input Format - Quark (PLASMA)

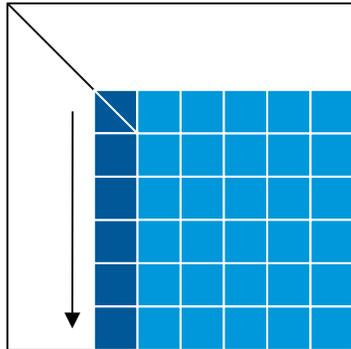
```
for (k = 0; k < A.mt; k++) {
  Insert_Task( zgeqrt, A[k][k], INOUT,
              T[k][k], OUTPUT);
  for (m = k+1; m < A.mt; m++) {
    Insert_Task( ztsqrt, A[k][k], INOUT | REGION_D|REGION_U,
                A[m][k], INOUT | LOCALITY,
                T[m][k], OUTPUT);
  }
  for (n = k+1; n < A.nt; n++) {
    Insert_Task( zunmqr, A[k][k], INPUT | REGION_L,
                T[k][k], INPUT,
                A[k][m], INOUT);
    for (m = k+1; m < A.mt; m++) {
      Insert_Task( ztismqr, A[k][n], INOUT,
                  A[m][n], INOUT | LOCALITY,
                  A[m][k], INPUT,
                  T[m][k], INPUT);
    }
  }
}
```

- Sequential C code
- Annotated through QUARK-specific syntax
 - **Insert_Task**
 - **INOUT, OUTPUT, INPUT**
 - **REGION_L, REGION_U, REGION_D, ...**
 - **LOCALITY**
- Executes thru the QUARK RT to run on multicore SMPs



- **Algorithm**
 - equivalent to LAPACK
- **Numerics**
 - same as LAPACK
- **Performance**
 - comparable to vendor on few cores
 - much better than vendor on many cores



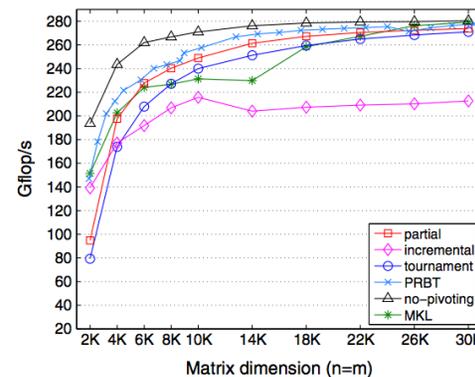


- **Algorithm**
 - equivalent to LAPACK
 - same pivot vector
 - same L and U factors
 - same forward substitution procedure

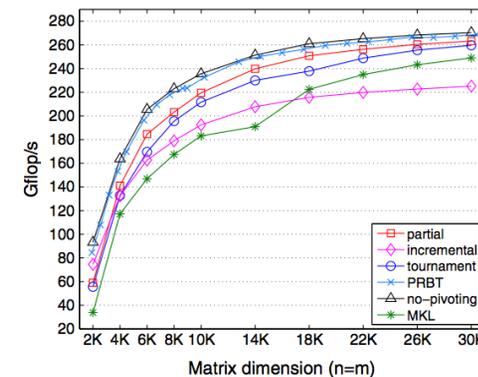
- **Numerics**
 - same as LAPACK

- **Performance**
 - comparable to vendor on few cores
 - much better than vendor on many cores

16 Sandy Bridge cores



Factorization alone, using 16 cores



Factorization and solve with iterative refinement, using 16 cores

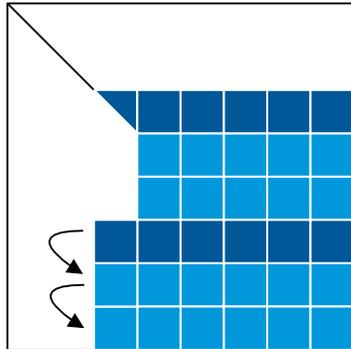


ICL

Algorithms

incremental QR Factorization

```
PLASMA_[scdz]geqrt[_Tile][_Async]()
```



- **Algorithm**
 - the same R factor as LAPACK (absolute values)
 - different set of Householder reflectors
 - different Q matrix
 - different Q generation / application procedure
- **Numerics**
 - same as LAPACK
- **Performance**
 - comparable to vendor on few cores
 - much better than vendor on many cores



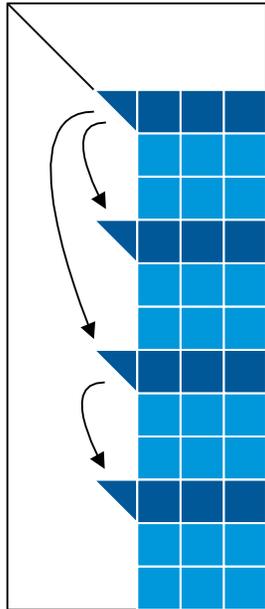
ICL

Algorithms

incremental QR Factorization (Communication Avoiding)

```
PLASMA_[scdz]geqrt[_Tile][_Async]()
```

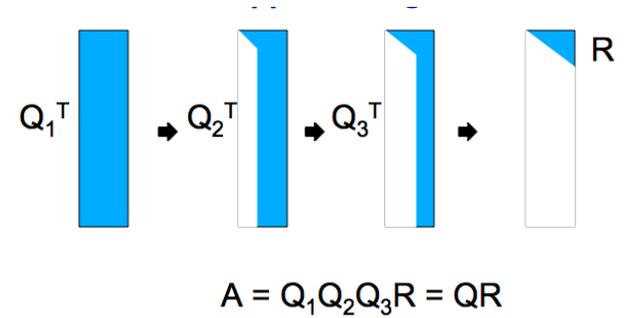
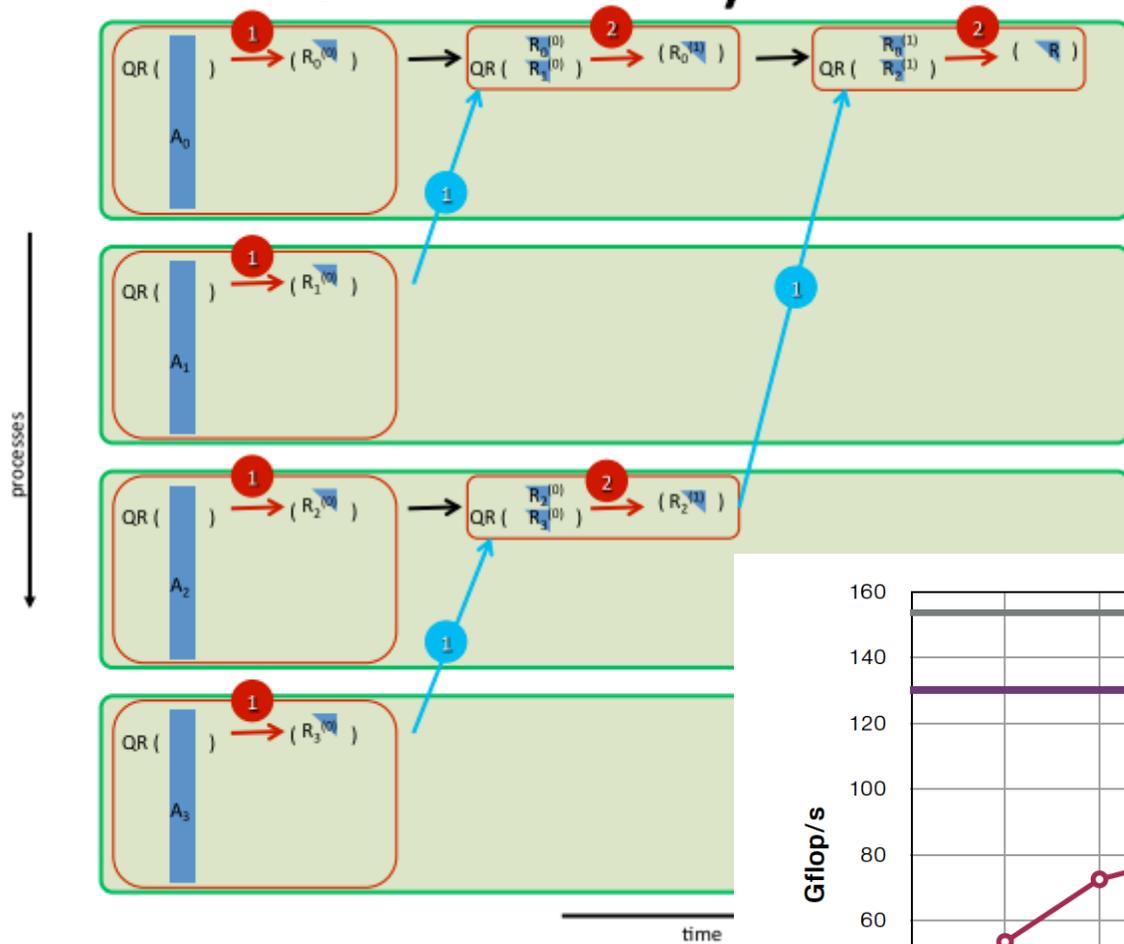
```
PLASMA_Set(  
    PLASMA_HOUSEHOLDER_MODE,  
    PLASMA_TREE_HOUSEHOLDER);
```



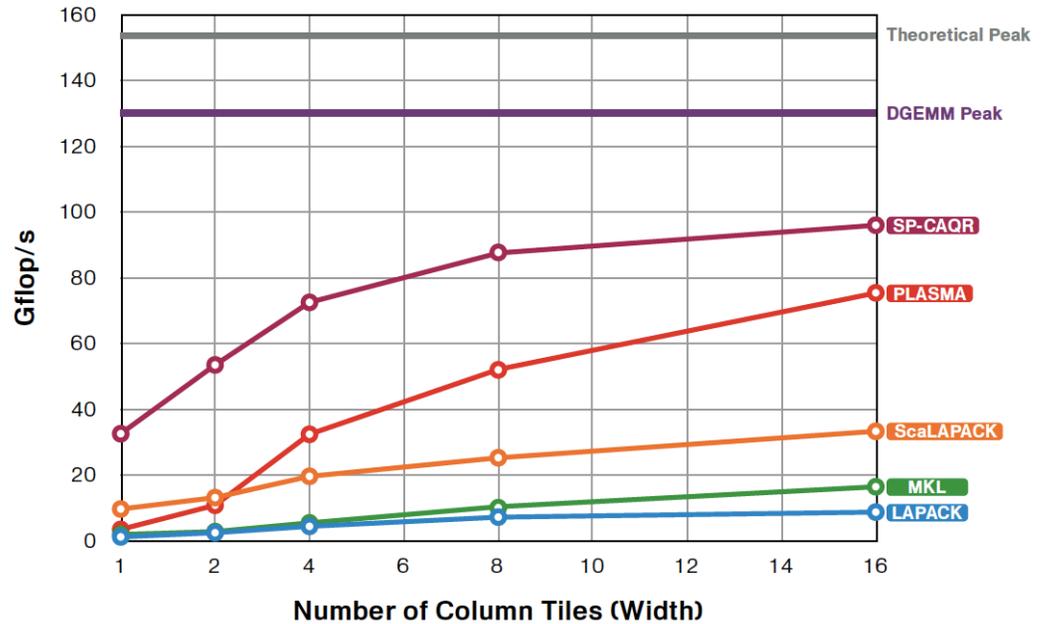
- **Algorithm**
 - the same R factor as LAPACK (absolute values)
 - different set of Householder reflectors
 - different Q matrix
 - different Q generation / application procedure
- **Numerics**
 - same as LAPACK
- **Performance**
 - absolutely superior for tall matrices

Communication Avoiding QR

Example



Quad-socket, quad-core machine Intel Xeon EMT64 E7340 at 2.39 GHz.
 Theoretical peak is 153.2 Gflop/s with 16 cores.
 Matrix size 51200 by 3200

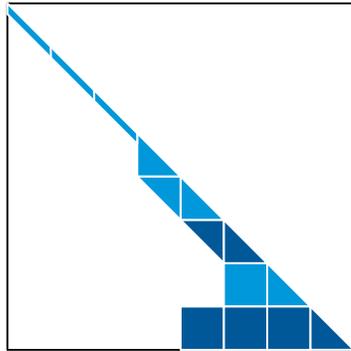




Algorithms

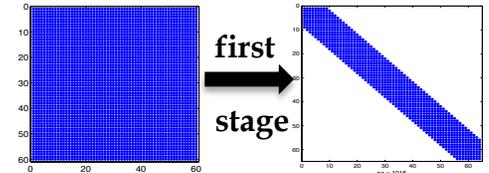
three-stage symmetric EVP

PLASMA_[scdz]syev[_Tile][_Async]()



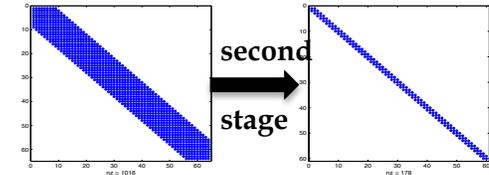
- **Algorithm**

- two-stage tridiagonal reduction + QR Algorithm
- fast eigenvalues, slower eigenvectors (possibility to calculate a subset)



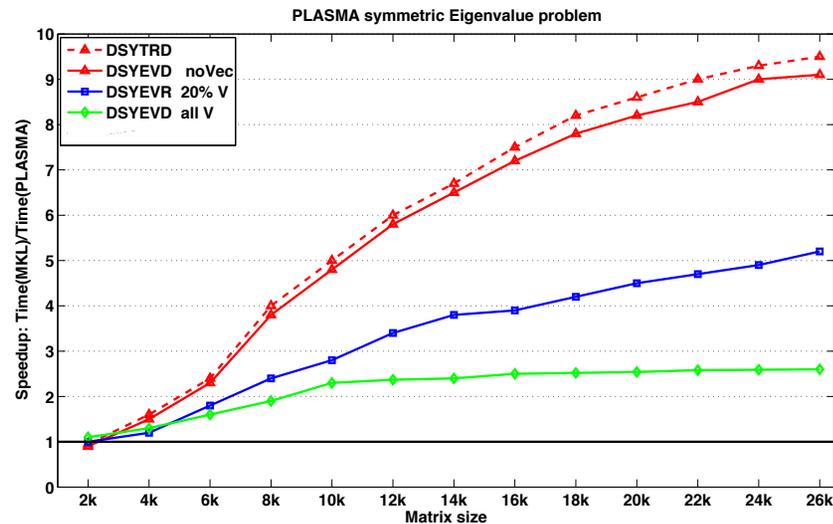
- **Numerics**

- same as LAPACK

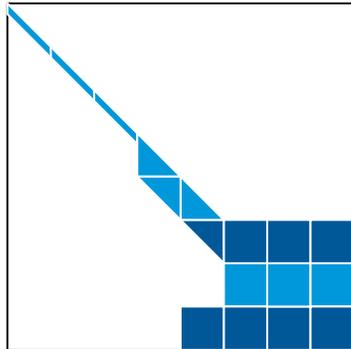


- **Performance**

- comparable to MKL for very small problems
- absolutely superior for larger problems



16 cores of Intel Sandy Bridge



- **Algorithm**

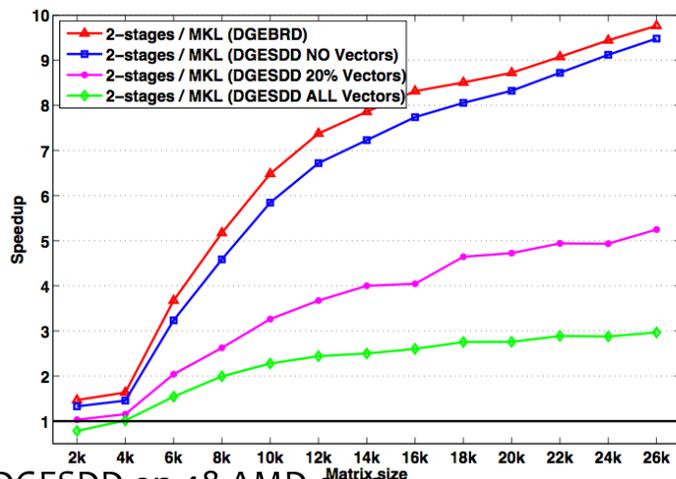
- two-stage bidiagonal reduction + QR iteration
- fast singular values, slower singular vectors (possibility of calculating a subset)

- **Numerics**

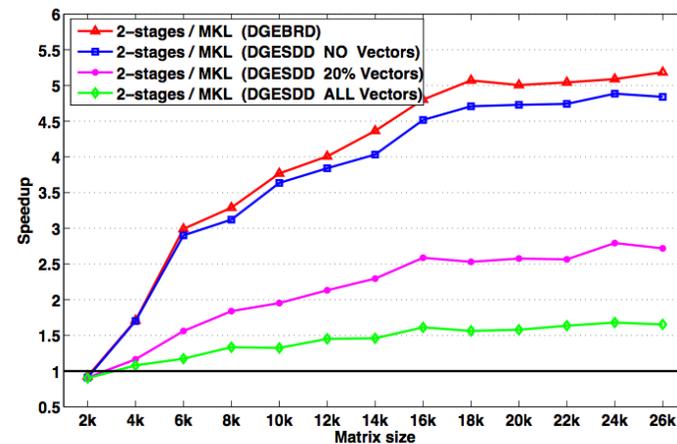
- same as LAPACK

- **Performance**

- comparable with MKL for very small problems
- absolutely superior for larger problems



DGESDD on 48 AMD cores

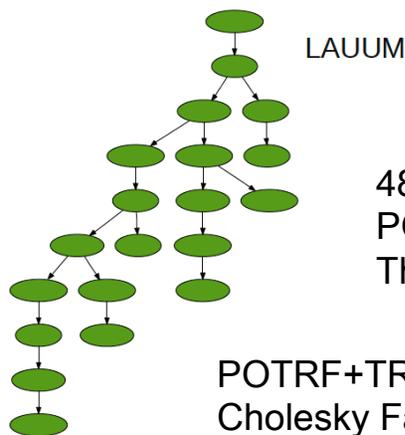
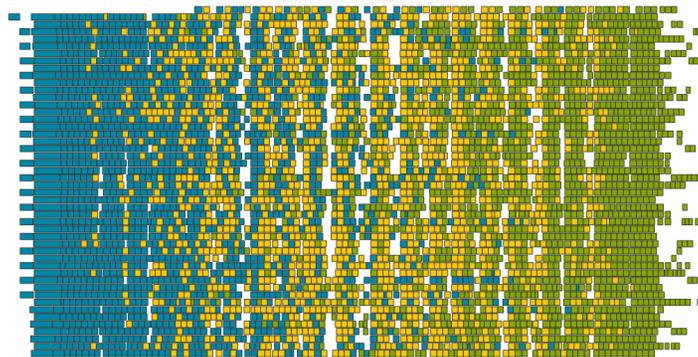
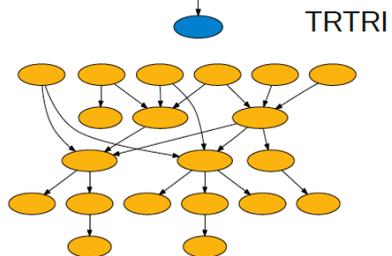
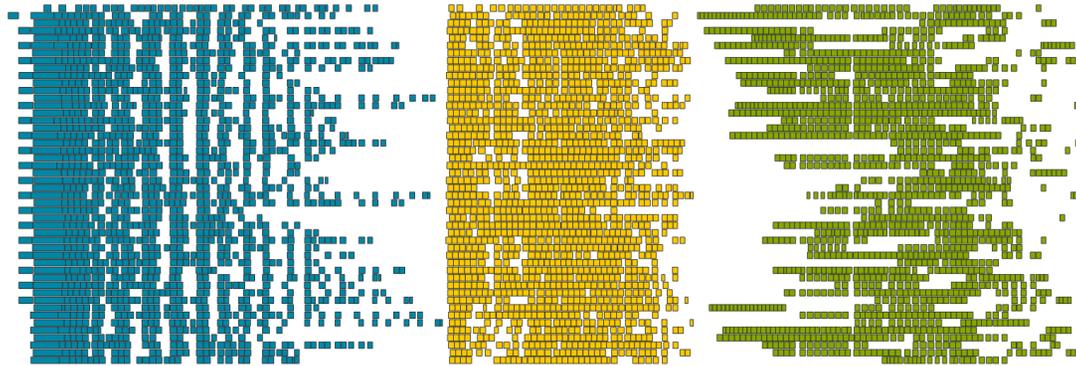
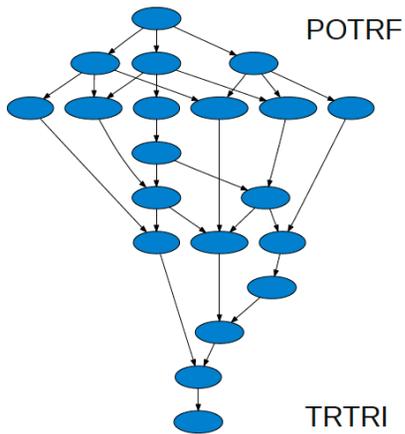


DGESDD on 16 Sandy Bridge cores



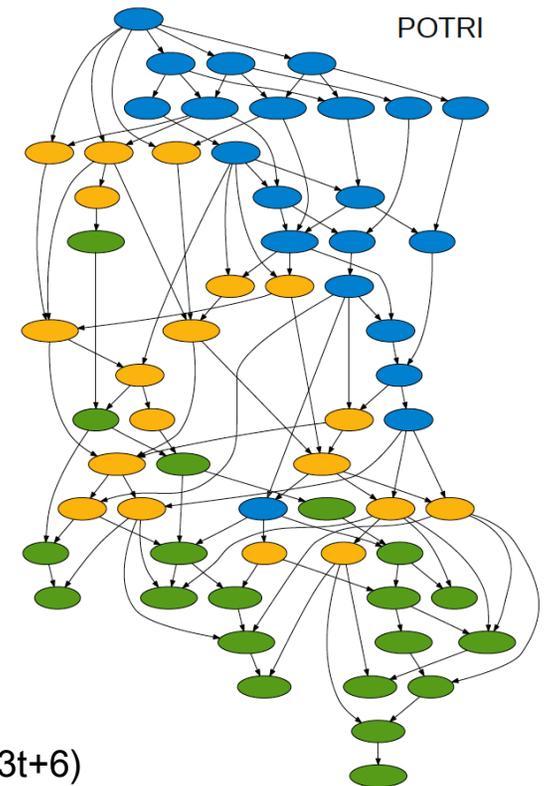
Pipelining: Cholesky Inversion

3 Steps: Factor, Invert L, Multiply L's



48 cores
POTRF, TRTRI and LAUUM.
The matrix is 4000 x 4000, tile size is 200 x 200,

POTRF+TRTRI+LAUUM: $25(7t-3)$
Cholesky Factorization alone: $3t-2$



Pipelined: $18(3t+6)$

Random Butterfly Pivoting (RBP)

- **To solve $Ax = b$:**
 - **Compute $A_r = U^T A V$, with U and V random matrices**
 - **Factorize A_r without pivoting (GENP)**
 - **Solve $A_r y = U^T b$ and then Solve $x = Vy$**
- **U and V are Recursive Butterfly Matrices**
 - **Randomization is cheap ($O(n^2)$ operations)**
 - **GENP is fast (“Cholesky” speed, take advantage of the GPU)**
 - **Accuracy is in practice similar to GEPP (with iterative refinement), but...**

A **butterfly matrix** is defined as any n -by- n matrix of the form:

$$B = \frac{1}{\sqrt{2}} \begin{pmatrix} R & S \\ R & -S \end{pmatrix}$$

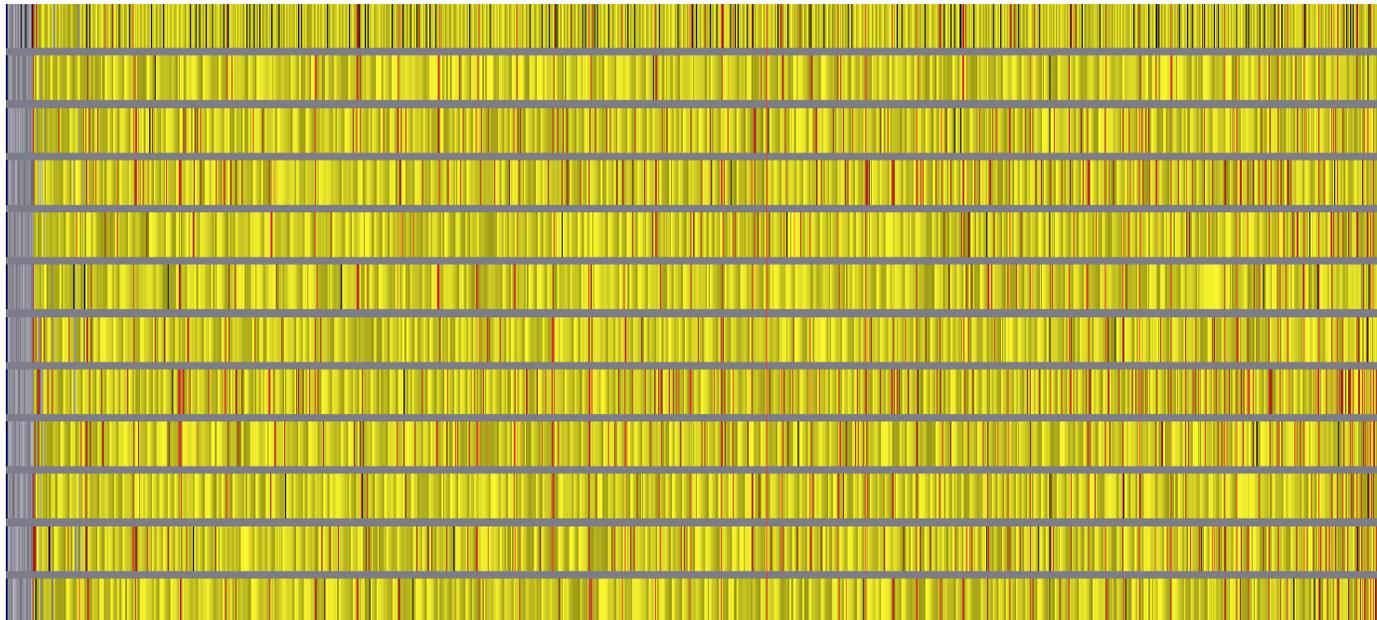
where R and S are random diagonal matrices.

$$B = \begin{pmatrix} \text{red} & \text{green} \\ \text{red} & \text{green} \end{pmatrix}$$

Think of this as a preconditioner step.

Goal: Transform A into a matrix that would be sufficiently “random” so that, with a probability close to 1, pivoting is not needed.

PLASMA RBT execution trace



- with $n=2000$, $nb=250$ on 12-core AMD Opteron -

Partial randomization (i.e. gray) is inexpensive.

Factorization without pivoting is scalable without synchronizations.

Mixed Precision Methods

- **Mixed precision, use the lowest precision required to achieve a given accuracy outcome**
 - **Improves runtime, reduce power consumption, lower data movement**
 - **Reformulate to find correction to solution, rather than solution; Δx rather than x .**

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

$$\boxed{x_{i+1} - x_i} = -\frac{f(x_i)}{f'(x_i)}$$

Idea Goes Something Like This...

- **Exploit 32 bit floating point as much as possible.**
 - **Especially for the bulk of the computation**
- **Correct or update the solution with selective use of 64 bit floating point to provide a refined results**
- **Intuitively:**
 - **Compute a 32 bit result,**
 - **Calculate a correction to 32 bit result using selected higher precision and,**
 - **Perform the update of the 32 bit results with the correction using high precision.**



Mixed-Precision Iterative Refinement

- Iterative refinement for dense systems, $Ax = b$, can work this way.

```
L U = lu(A) O(n3)
x = L\ (U\b) O(n2)
r = b - Ax O(n2)
WHILE || r || not small enough
    z = L\ (U\r) O(n2)
    x = x + z O(n1)
    r = b - Ax O(n2)
END
```

- Wilkinson, Moler, Stewart, & Higham provide error bound for SP fl pt results when using DP fl pt.



Mixed-Precision Iterative Refinement

- Iterative refinement for dense systems, $Ax = b$, can work this way.

$L U = \text{lu}(A)$	SINGLE	$O(n^3)$
$x = L \backslash (U \backslash b)$	SINGLE	$O(n^2)$
$r = b - Ax$	DOUBLE	$O(n^2)$
WHILE $\ r \ $ not small enough		
$z = L \backslash (U \backslash r)$	SINGLE	$O(n^2)$
$x = x + z$	DOUBLE	$O(n^1)$
$r = b - Ax$	DOUBLE	$O(n^2)$
END		

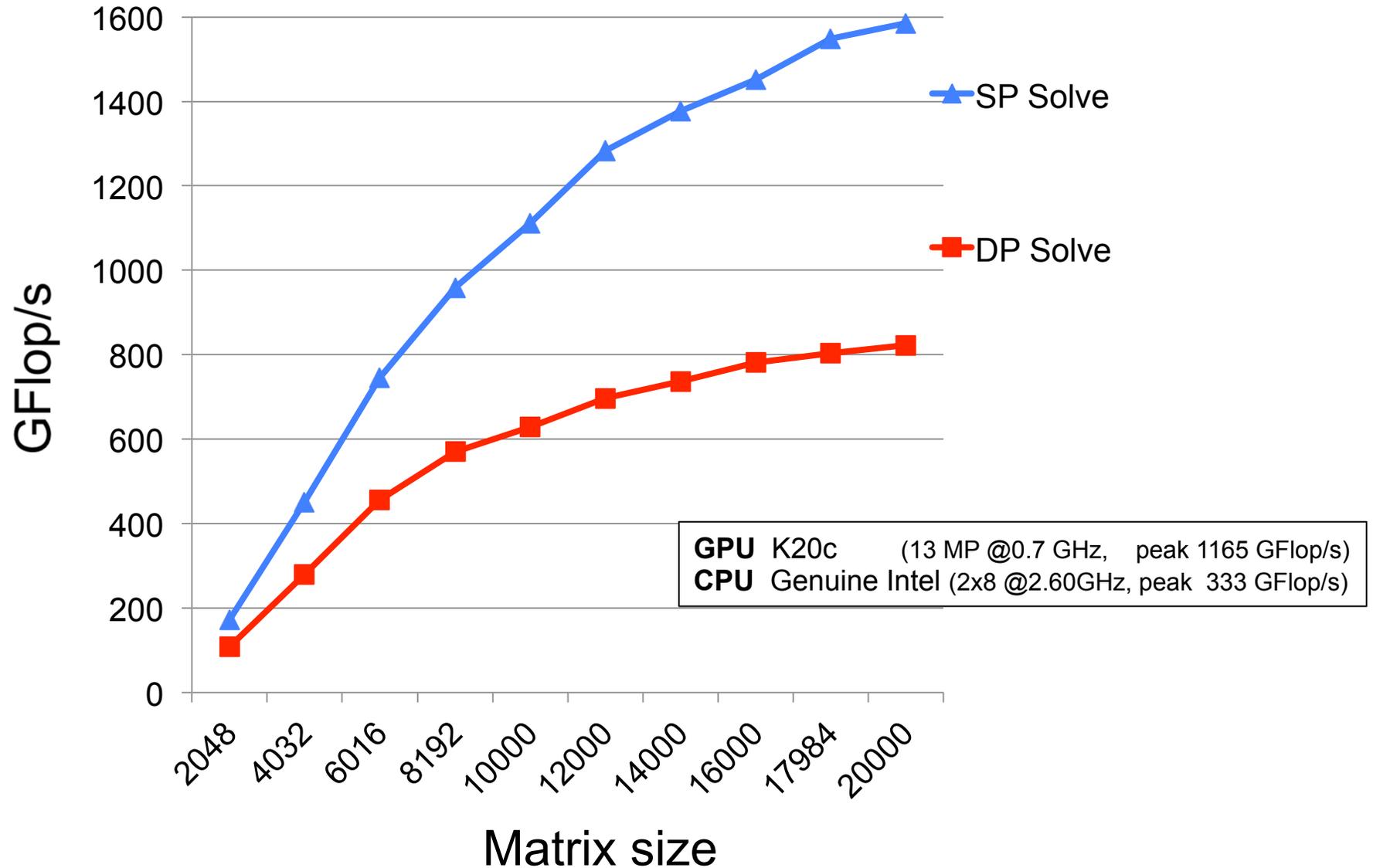
- Wilkinson, Moler, Stewart, & Higham provide error bound for SP fl pt results when using DP fl pt.
- It can be shown that using this approach we can compute the solution to 64-bit floating point precision.

- Requires extra storage, total is 1.5 times normal;
- $O(n^3)$ work is done in lower precision
- $O(n^2)$ work is done in high precision
- Problems if the matrix is ill-conditioned in sp; $O(10^8)$



Mixed precision iterative refinement

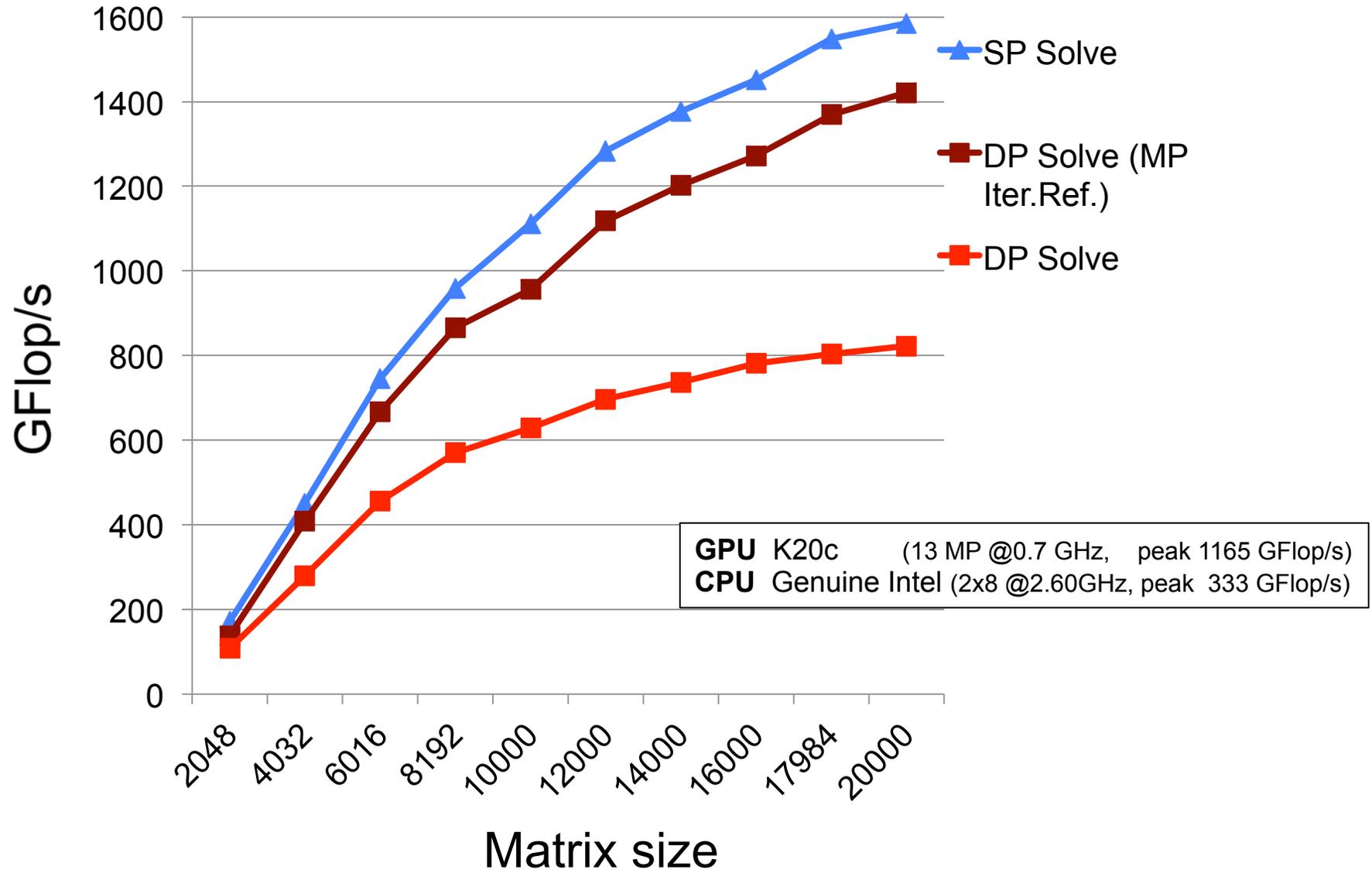
Solving general dense linear systems using mixed precision iterative refinement





Mixed precision iterative refinement

Solving general dense linear systems using mixed precision iterative refinement





Critical Issues at Peta & Exascale for Algorithm and Software Design

- **Synchronization-reducing algorithms**
 - Break Fork-Join model
- **Communication-reducing algorithms**
 - Use methods which have lower bound on communication
- **Mixed precision methods**
 - 2x speed of ops and 2x speed for data movement
- **Autotuning**
 - Today's machines are too complicated, build "smarts" into software to adapt to the hardware
- **Fault resilient algorithms**
 - Implement algorithms that can recover from failures/bit flips
- **Reproducibility of results**
 - Today we can't guarantee this. We understand the issues, but some of our "colleagues" have a hard time with this.



Collaborators / Software / Support

- ◆ **PLASMA**
<http://icl.cs.utk.edu/plasma/>
- ◆ **MAGMA**
<http://icl.cs.utk.edu/magma/>
- ◆ **Quark (RT for Shared Memory)**
<http://icl.cs.utk.edu/quark/>
- ◆ **PaRSEC (Parallel Runtime Scheduling and Execution Control)**
<http://icl.cs.utk.edu/parsec/>



- ◆ Collaborating partners
University of Tennessee, Knoxville
University of California, Berkeley
University of Colorado, Denver

MAGMA



PLASMA

