



Dense Linear Algebra



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE

Mark Gates

mgates3@icl.utk.edu

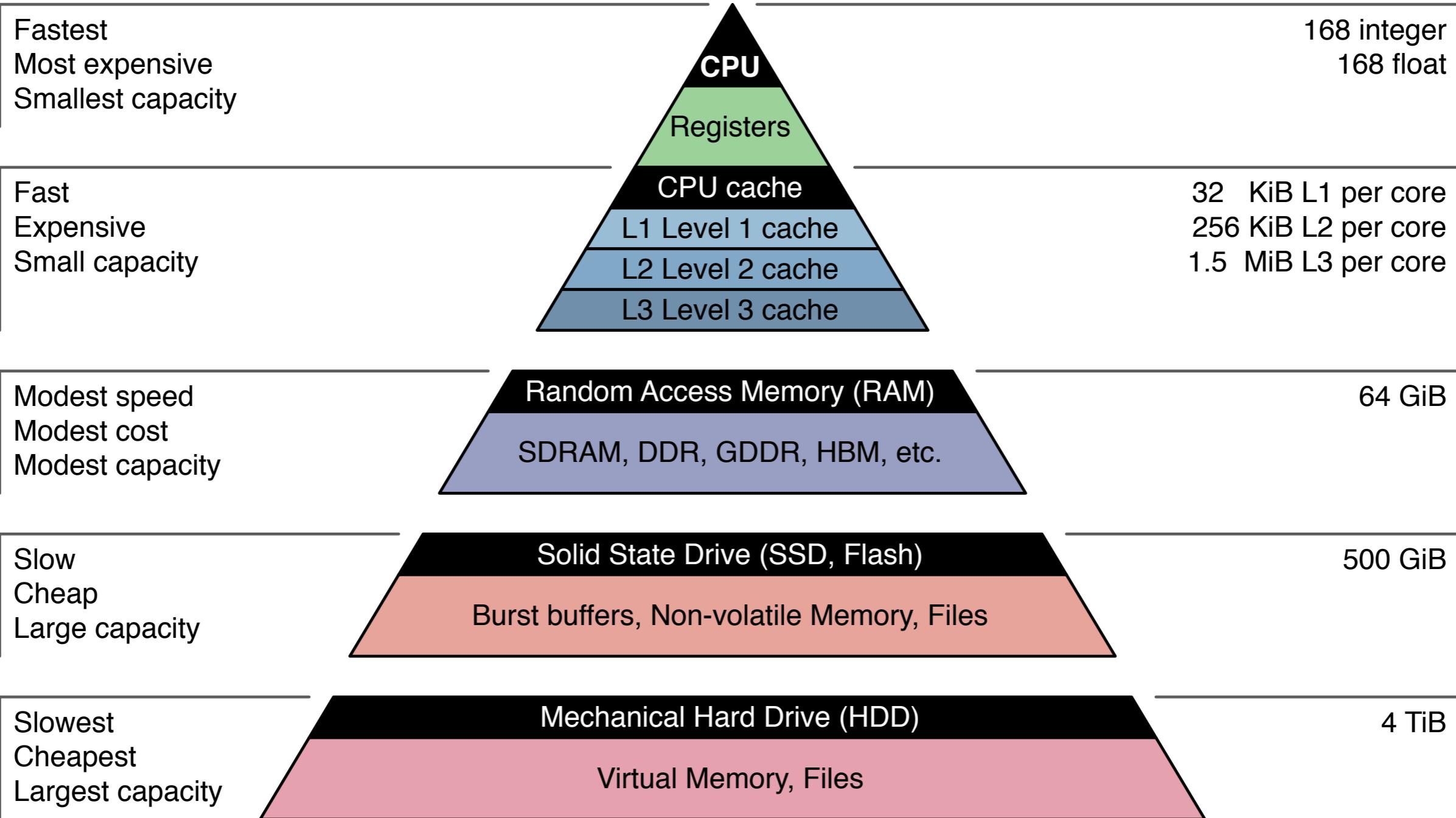
<http://www.icl.utk.edu/~mgates3/>

Outline

- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

Memory hierarchy

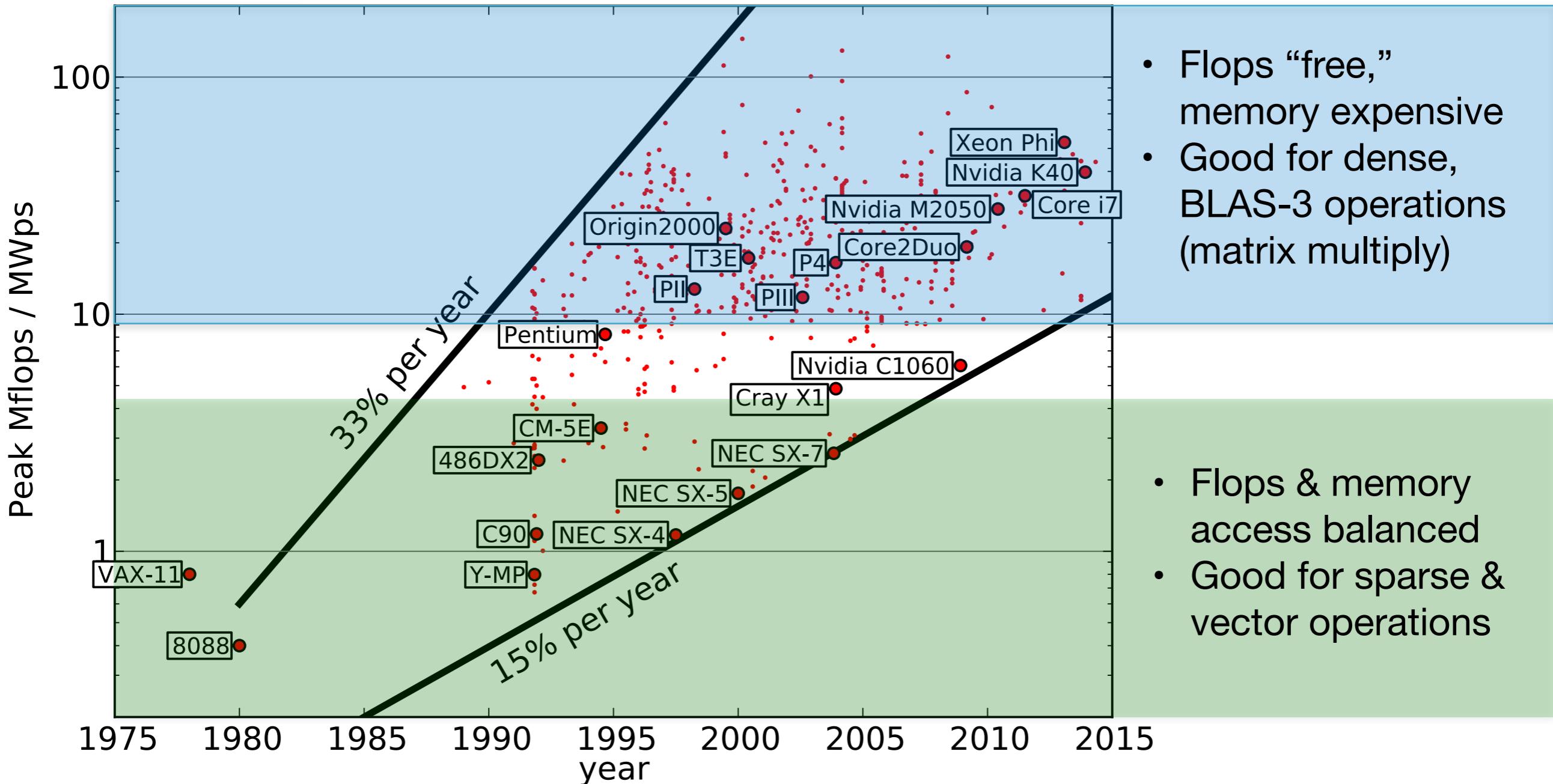
Examples (Haswell)



Adapted from illustration by Ryan Leng

Compute vs. memory speed

- Machine balance (# flops per RAM access)



Outline

- Legacy Software
 - **BLAS**
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

BLAS: Basic Linear Algebra Subroutines

- Level 1 BLAS – vector operations
 - $O(n)$ data and flops (floating point operations)
 - Memory bound:
 $O(1)$ flops per memory access

$$y = \alpha x + \beta y$$

- Level 2 BLAS – matrix-vector operations
 - $O(n^2)$ data and flops
 - Memory bound:
 $O(1)$ flops per memory access

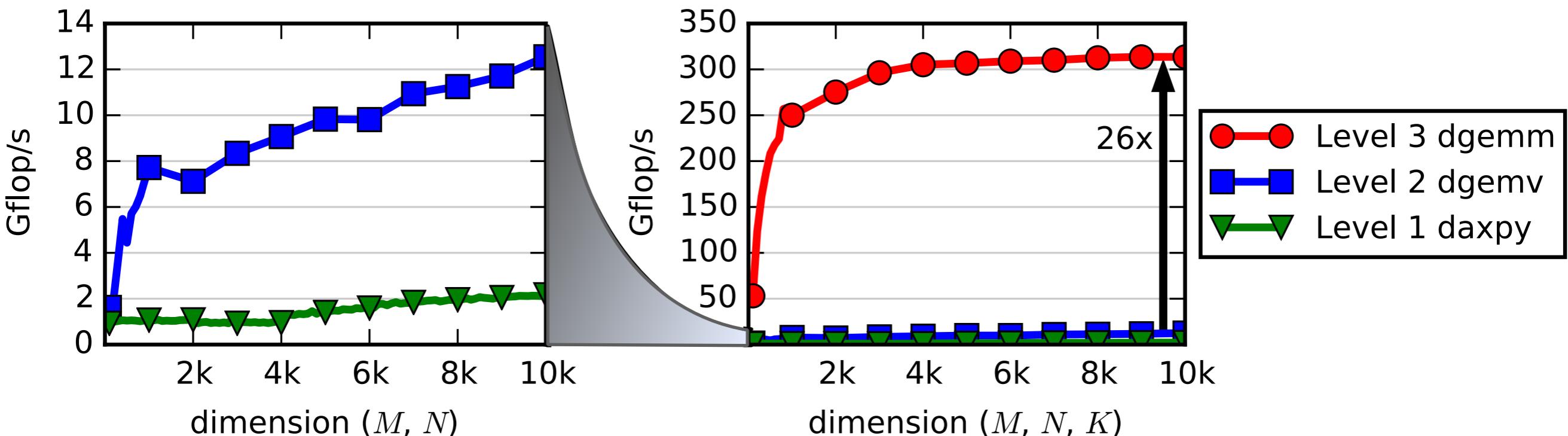
$$y = \alpha A x + \beta y$$

- Level 3 BLAS – matrix-matrix operations
 - $O(n^2)$ data, $O(n^3)$ flops
 - Surface-to-volume effect
 - Compute bound:
 $O(n)$ flops per memory access

$$C = \alpha A B + \beta C$$

BLAS: Basic Linear Algebra Subroutines

- Intel Sandy Bridge (2 socket E5-2670)
 - Peak 333 Gflop/s = 2.6 GHz * 16 cores * 8 double-precision flops/cycle †
 - Max memory bandwidth 51 GB/s ‡



† <http://stackoverflow.com/questions/15655835/flops-per-cycle-for-sandy-bridge-and-haswell-sse2-avx-avx2>

‡ http://ark.intel.com/products/64595/Intel-Xeon-Processor-E5-2670-20M-Cache-2_60-GHz-8_00-GTs-Intel-QPI

Memory bandwidth

- Intel Ark
 - <https://ark.intel.com/>
 - Search for model number (e.g., “E5-2697 v4” for Broadwell, see /proc/cpuinfo)

```
prompt> less /proc/cpuinfo
...
model name      : Intel(R) Xeon(R) CPU E5-2697 v4 @ 2.30GHz
...
cache size     : 46080 KB
siblings        : 18
physical id    : 0
core id         : 1
...
flags          : fpu vme de pse tsc msr pae mce cx8 apic sep mtrr pge mca cmov pat pse36 clflush dts
acpi mmx fxsr sse sse2 ss ht tm pbe syscall nx pdpe1gb rdtscp lm constant_tsc arch_perfmon pebs bts
rep_good nopl xtopology nonstop_tsc aperfmpf eagerfpu pni pclmulqdq dtes64 monitor ds_cpl vmx smx
est tm2 ssse3 fma cx16 xtpr pdcm pcid dca sse4_1 sse4_2 x2apic movbe popcnt tsc_deadline_timer aes
xsave avx f16c cdrand lahf_lm abm 3dnowprefetch ida arat epb pln pts dtherm intel_pt tpr_shadow vnmi
flexpriority ept vpid fsgsbase tsc_adjust bmi1 bmi2 avx2 smep bmi2 erms invpcid rtm cqm rdseed adx smap
xsaveopt cqmq_llc cqmq_occu_llc cqmq_mbm_total cqmq_mbm_local
...
```

siblings = # hyperthreads per socket (see ark)
each physical id / core id combination is unique core

supports SSE and AVX vector instructions

- Stream benchmark
 - <https://www.cs.virginia.edu/stream/>
 - Add -fopenmp or equivalent to Makefile CFLAGS and FFLAGS
 - Adjust STREAM_ARRAY_SIZE depending on cache size;
see instructions in code (I changed this to be dynamically allocated)

BLAS naming scheme

- One or two letter **data type** (precision)

- i = integer (e.g., index)
- s = single (float)
- d = double
- c = single-complex
- z = double-complex

- Two letter **matrix type** (BLAS 2, 3)

- ge = general nonsymmetric
- sy = symmetric ($A = A^T$)
- he = complex Hermitian ($A = A^H$)
- tr = triangular (L or U)
- Also banded and packed formats

- Two or more letter **function**, e.g.

- mv = matrix-vector product
- mm = matrix-matrix product
- etc.

Example

dgemm

$$C = AB + C$$



BLAS naming scheme

- One or two letter **data type**

- i = integer (e.g., index)
- s = single (float)
- d = double
- c = single-complex
- z = double-complex

- Two letter **matrix type** (BLAS 2, 3)

- ge = general nonsymmetric
- sy = symmetric ($A = A^T$)
- he = complex Hermitian ($A = A^H$)
- tr = triangular (L or U)
- Also banded and packed formats

- Two or more letter **function**, e.g.

- mv = matrix-vector product
- mm = matrix-matrix product
- etc.

BLAS 1 examples

sdot	$result = x^T y$ (single)
ddot	$result = x^T y$ (double)
cdotc	$result = x^H y$ (single-complex)
cdotu	$result = x^T y$ (single-complex)
zdotc	$result = x^H y$ (double-complex)
zdotu	$result = x^T y$ (double-complex)

_axpy $y = \alpha x + y$

_scal $y = \alpha y$

_copy $y = x$

_swap $x \leftrightarrow y$

_nrm2 $result = \|x\|_2$

_asum $result = \sum_i |x_i|$

i_amax $index = \text{argmax}_i |x_i|$

_rot Apply Given's rotation:

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} c & s \\ -s & c \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

BLAS naming scheme

- One or two letter **data type**

- i = integer (e.g., index)
- s = single (float)
- d = double
- c = single-complex
- z = double-complex

- Two letter **matrix type** (BLAS 2, 3)

- ge = general nonsymmetric
- sy = symmetric ($A = A^T$)
- he = complex Hermitian ($A = A^H$)
- tr = triangular (L or U)
- Also banded and packed formats

- Two or more letter **function**, e.g.

- mv = matrix-vector product
- mm = matrix-matrix product
- etc.

BLAS 2 examples

_gemv	$y = Ax + y,$	A general
_symv	$y = Ax + y,$	A symmetric
_hemv	$y = Ax + y,$	A Hermitian
_ger	$C = xy^T + C,$	C general
_syr	$C = xx^T + C,$	C symmetric
_her	$C = xx^H + C,$	C Hermitian
_trmv	$x = Ax,$	A triangular
_trsv	<i>solve</i> $Ax = b,$	A triangular

BLAS 3 examples

_gemm	$C = AB + C,$	all general
_symm	$C = AB + C,$	A symmetric
_hemm	$C = AB + C,$	A Hermitian
_syrk	$C = AA^T + C,$	C symmetric
_herk	$C = AA^H + C,$	C Hermitian
_trmm	$X = AX \text{ or } XA,$	A triangular
_trsm	<i>solve</i> $AX = B,$	A triangular

Why is BLAS so important?

- BLAS are efficient, portable, parallel, and widely available
 - Vendors (Intel, AMD, IBM, Cray, NVIDIA, ...) provide highly optimized versions
 - Open source libraries available (ATLAS, OpenBLAS)
- Commonly used for high quality linear algebra and HPC software
- Performance of many applications depends on underlying BLAS

Outline

- Legacy Software
 - BLAS
 - **LINPACK**
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

LINPACK

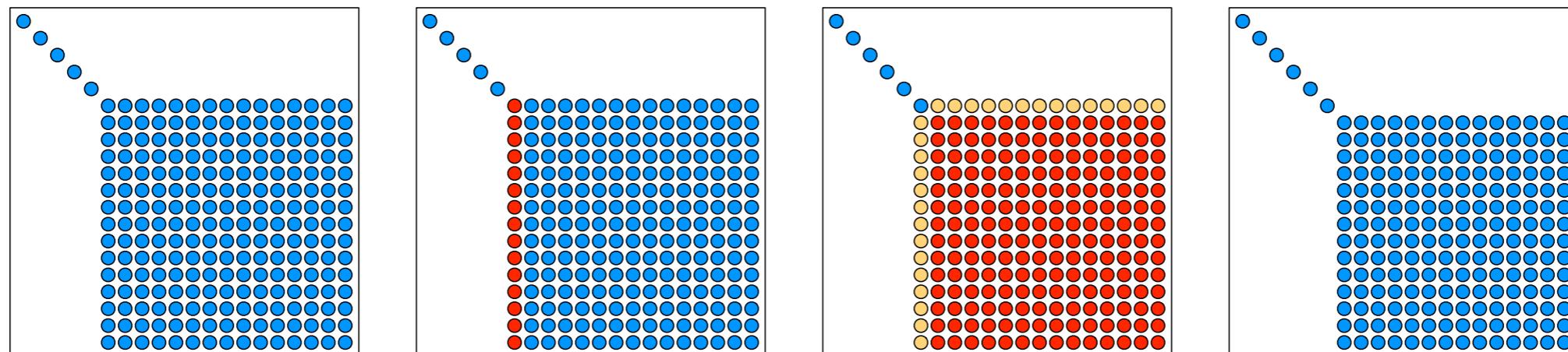
LINear algebra PACKage

- Math software library
 - solving dense & banded linear systems
 - singular value decomposition (SVD)
- Written in 1970s using Fortran 66
- Aiming for software portability and efficiency
- Level 1 BLAS
- Four primary contributors:
 - **Jack Dongarra**, Argonne
 - **Jim Bunch**, U. California, San Diego
 - **Cleve Moler**, New Mexico
 - **Pete Stewart**, U. of Maryland
- Superseded by LAPACK

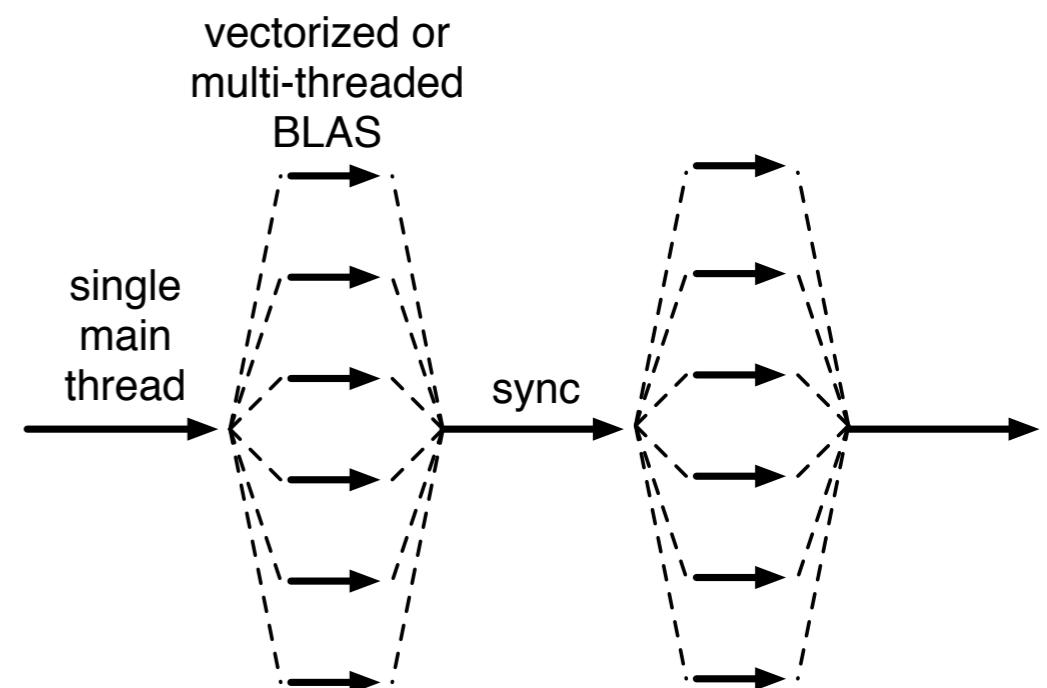
Computing in 1974

- High performance computers
 - IBM 370/195
 - CDC 7600
 - Univac 1110
 - DEC PDP-10
 - Honeywell 6030
- Fortran 66
- EISPACK released in 1974
 - Eigenvalue and singular value problems
 - Translation of Algol into Fortran
 - Did not use BLAS
- Level 1 BLAS — vector operations (1979)
 - Aimed at vector supercomputer architectures
- LINPACK released in 1979
 - About time of Cray 1

LU Factorization in LINPACK



- Factor one column at a time
 - i_amax and _scal
- Update each column of trailing matrix, one column at a time
 - _axpy
- Level 1 BLAS
- Bulk synchronous
 - Single main thread
 - Parallel work in BLAS
 - “Fork-and-join” model



Outline

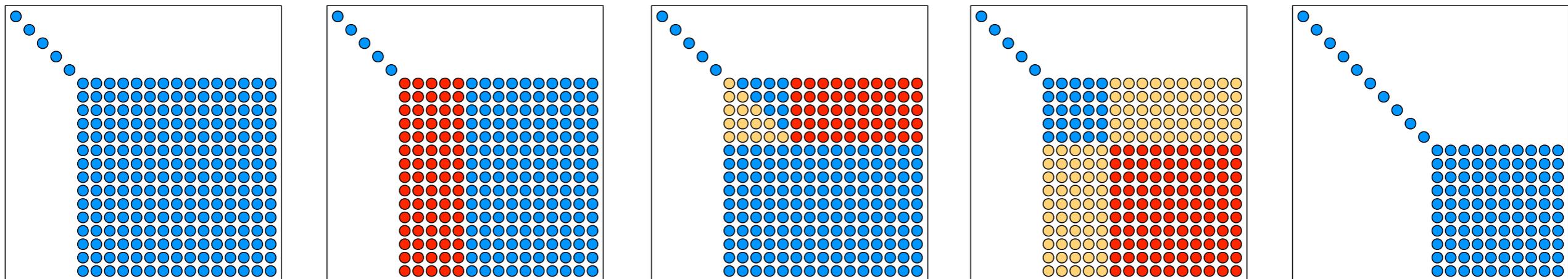
- Legacy Software
 - BLAS
 - LINPACK
 - **LAPACK**
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

LAPACK

Linear Algebra PACKage

- Reformulate linear algebra algorithms as **blocked algorithms** using BLAS-3
- Linear systems
 - Nearly 100% BLAS-3
- Singular value and symmetric eigenvalue problems
 - About 50% BLAS-3
- Nonsymmetric eigenvalue problems
 - About 80% BLAS-3
- 4 data types: single, double, single-complex, double-complex

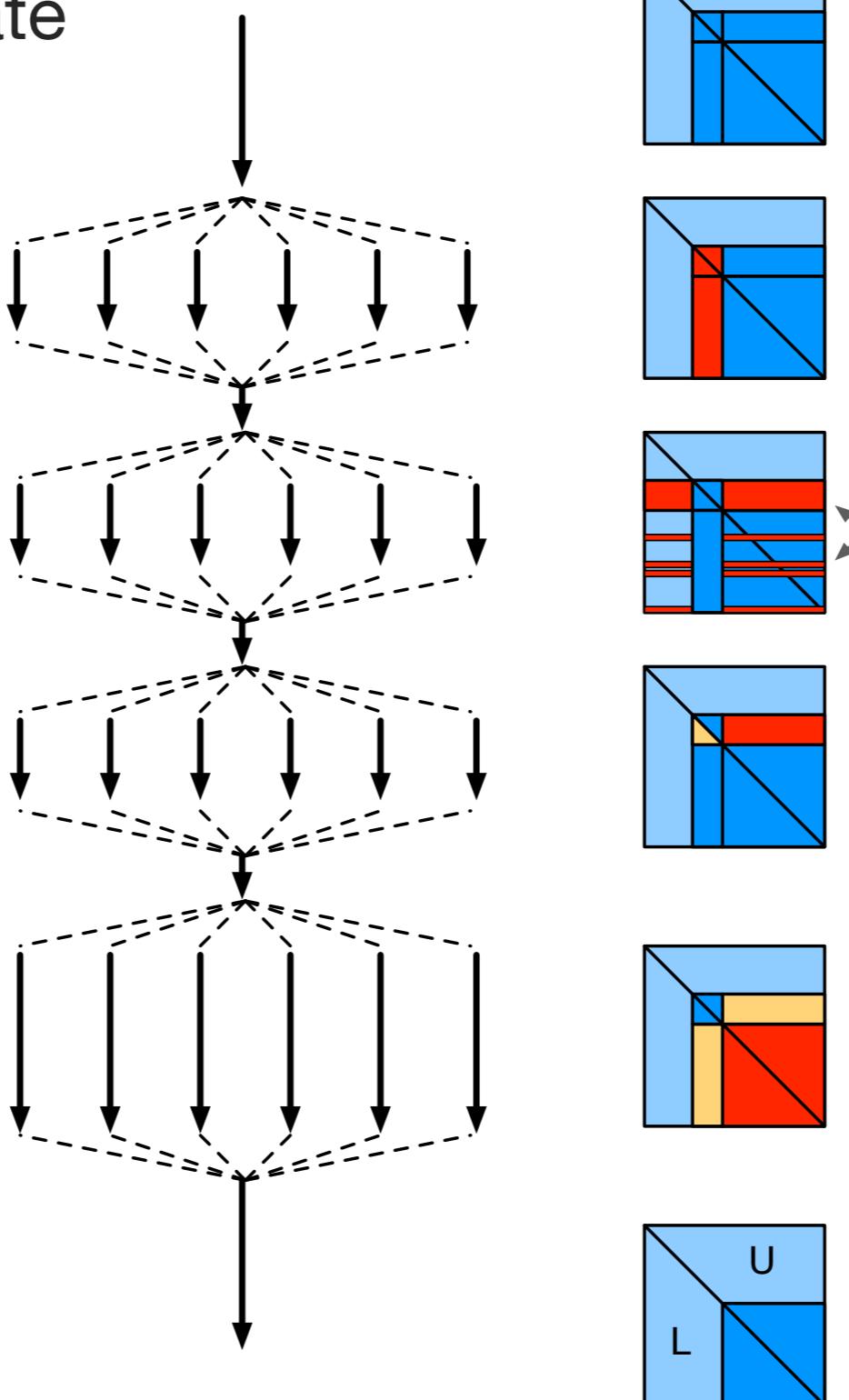
LU factorization in LAPACK



- Factor panel of n_b columns
 - getf2, unblocked BLAS-2 code
- Level 3 BLAS update block-row of U
 - trsm
- Level 3 BLAS update trailing matrix
 - gemm
 - Aimed at machines with cache hierarchy
- Bulk synchronous

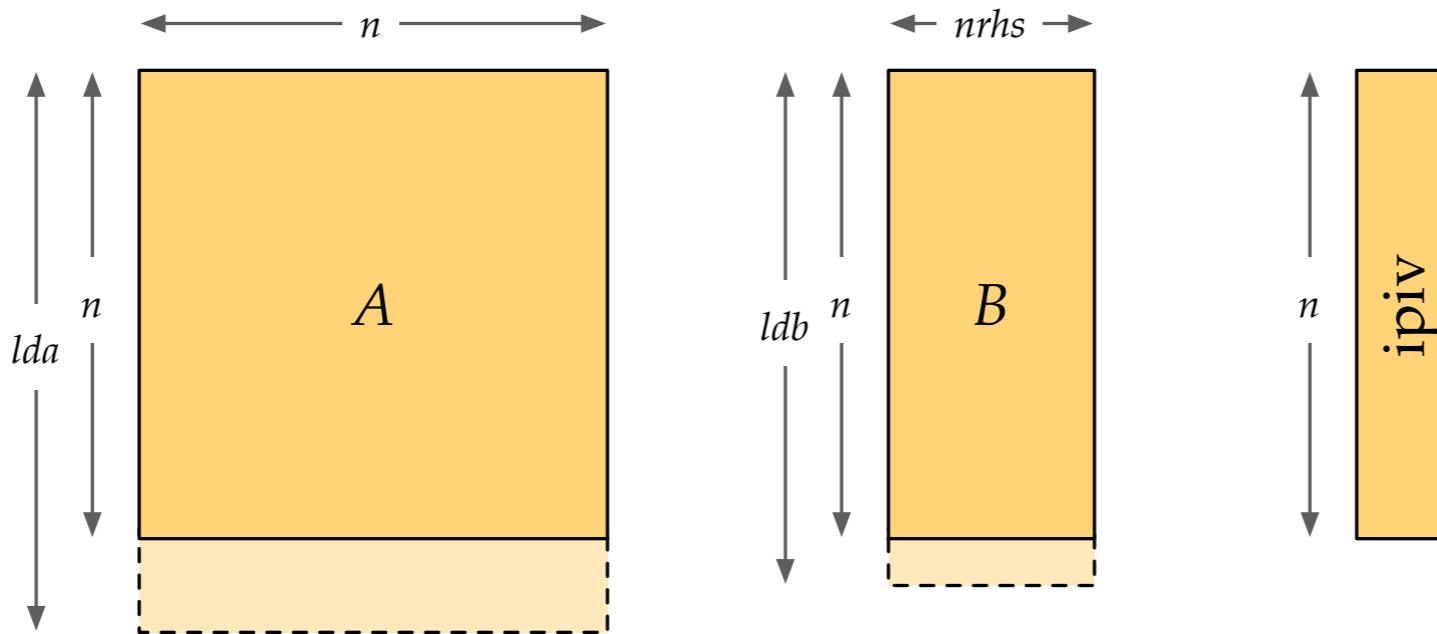
Parallelism in LAPACK

- Most flops in gemm update
 - $2/3 n^3$ term
 - Easily parallelized using multi-threaded BLAS
 - Done in any reasonable software
- Other operations lower order
 - Potentially expensive if not parallelized

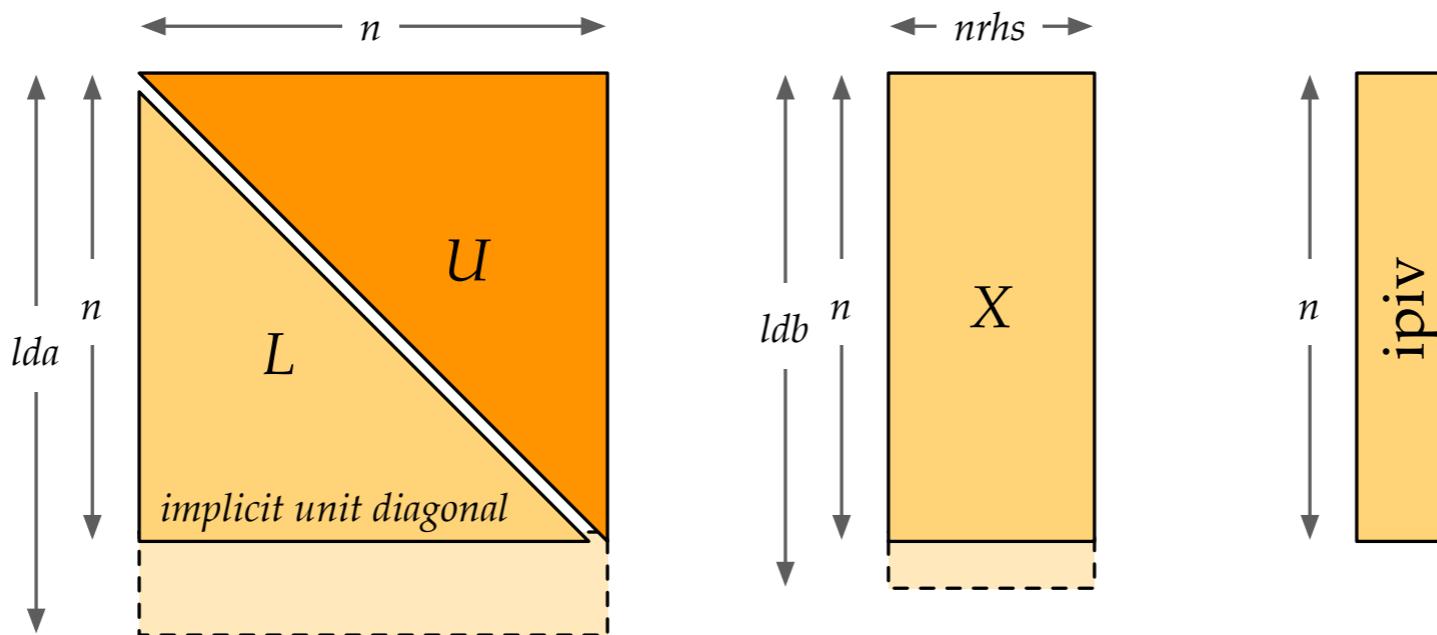


LAPACK routine, solve $AX = B$

- dgesv(n, nrhs, A, lda, ipiv, B, ldb, info)
- input:



- output:



info (error code)
= 0: no error
< 0: invalid argument
> 0: numerical error
(e.g., singular)

L, U overwrite A
 X overwrites B
Matrices stored column-wise

LAPACK naming

- Similar to BLAS naming
 - Data type + matrix type + function
- Additional matrix types:
 - po = real symmetric / complex Hermitian **positive definite** (SPD / HPD, all $\lambda > 0$)
 - or = **orthogonal** ($A^T A = I$)
 - un = **unitary** ($A^H A = I$)
 - Also banded, packed, etc. formats

LAPACK routines

- ~480 routines x 4 data types (s, d, c, z)
- **Driver routines:** solve an entire problem
 - sv Solve linear system: $Ax = b$
 - gesv general non-symmetric (LU)
 - posv symmetric positive definite (Cholesky)
 - sysv symmetric indefinite (LDL^T)
 - Also packed, banded, tridiagonal storage
 - ls Linear least squares: $Ax \approx b$, A is rectangular
 - gglse linear equality-constrained least squares
 - ggg1m general Gauss-Markov linear model
 - ev Eigenvalue decomposition: $Ax = \lambda x$ and $Ax = \lambda Mx$
 - syevd symmetric and symmetric generalized
 - geev non-symmetric and non-symmetric generalized
 - svd Singular value decomposition: $A = U\Sigma V^H$
 - gesvd standard and generalized
 - gesdd D&C (faster)

<http://www.icl.utk.edu/~mgates3/docs/>
<http://www.icl.utk.edu/~mgates3/docs/lapack.html>

LAPACK routines

- **Computation routines:** one step of problem

- trf triangular factorization (LU, Cholesky, LDL^T)
- trs triangular solve
- qrf orthogonal QR factorization
- mqr multiply by Q
- gqr generate Q
- etc.

- **Auxiliary routines,** begin with “la”

- lan_ matrix norm (one, inf, Frobenius, max); see SVD for 2-norm
- lascl scale matrix
- lacpy copy matrix
- laset set matrix & diagonal to constants
- etc.

<http://www.icl.utk.edu/~mgates3/docs/>
<http://www.icl.utk.edu/~mgates3/docs/lapack.html>

Outline

- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - **ScaLAPACK**
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

ScaLAPACK

Scalable Linear Algebra PACKage

- Distributed memory
- Message Passing
 - Clusters of SMPs
 - Supercomputers
- Dense linear algebra
- Modules
 - PBLAS: Parallel BLAS
 - BLACS: Basic Linear Algebra Communication Subprograms



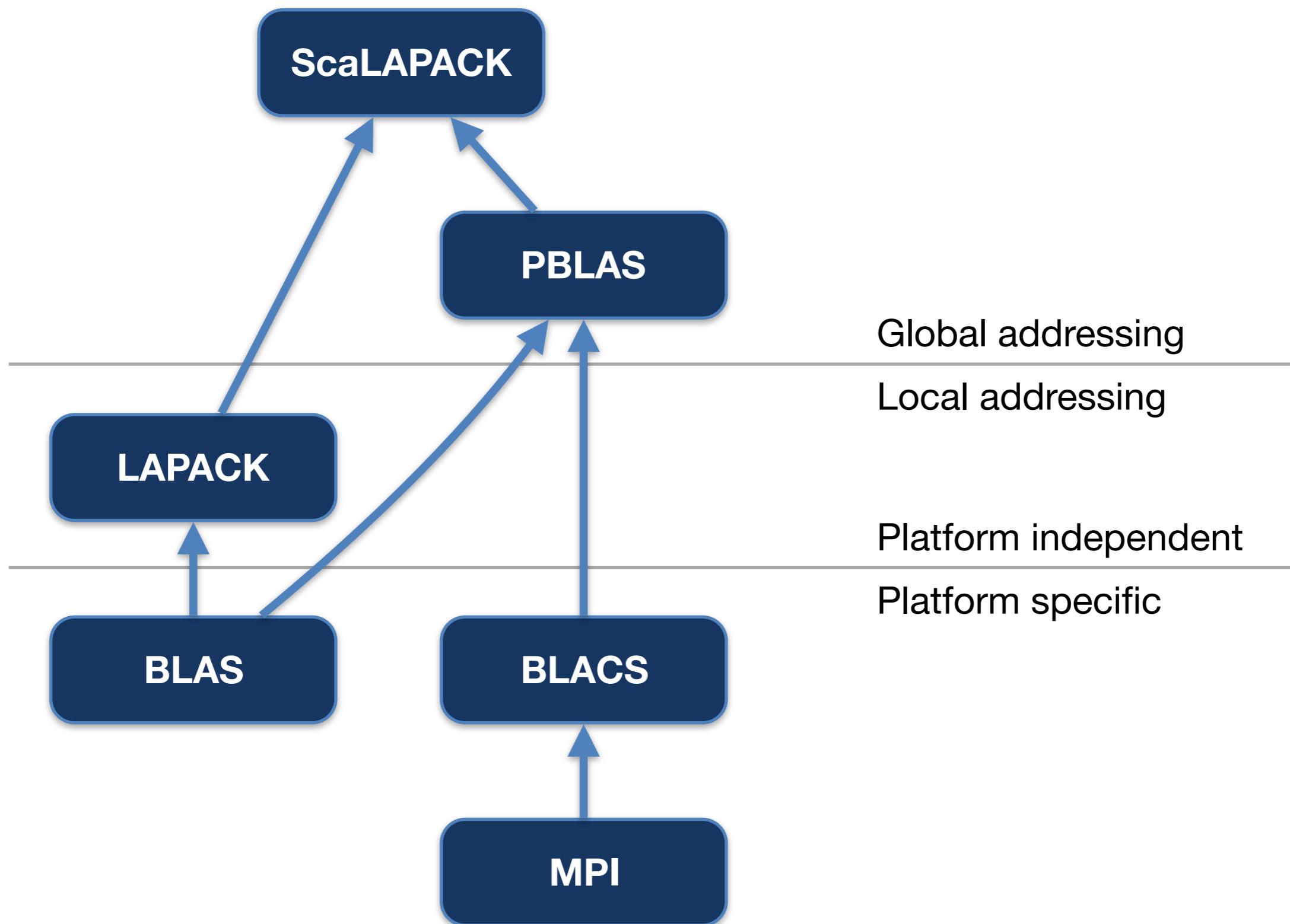
Advancing the Era of Accelerated Computing

PBLAS

- Similar to BLAS in functionality and naming
- Built on BLAS and BLACS
- Provide global view of matrix
- LAPACK: `dge____(m, n, A(ia, ja), lda, ...)`
 - Submatrix offsets implicit in pointer
- ScaLAPACK: `pdge____(m, n, A, ia, ja, descA, ...)`
 - Pass submatrix offsets and matrix descriptor

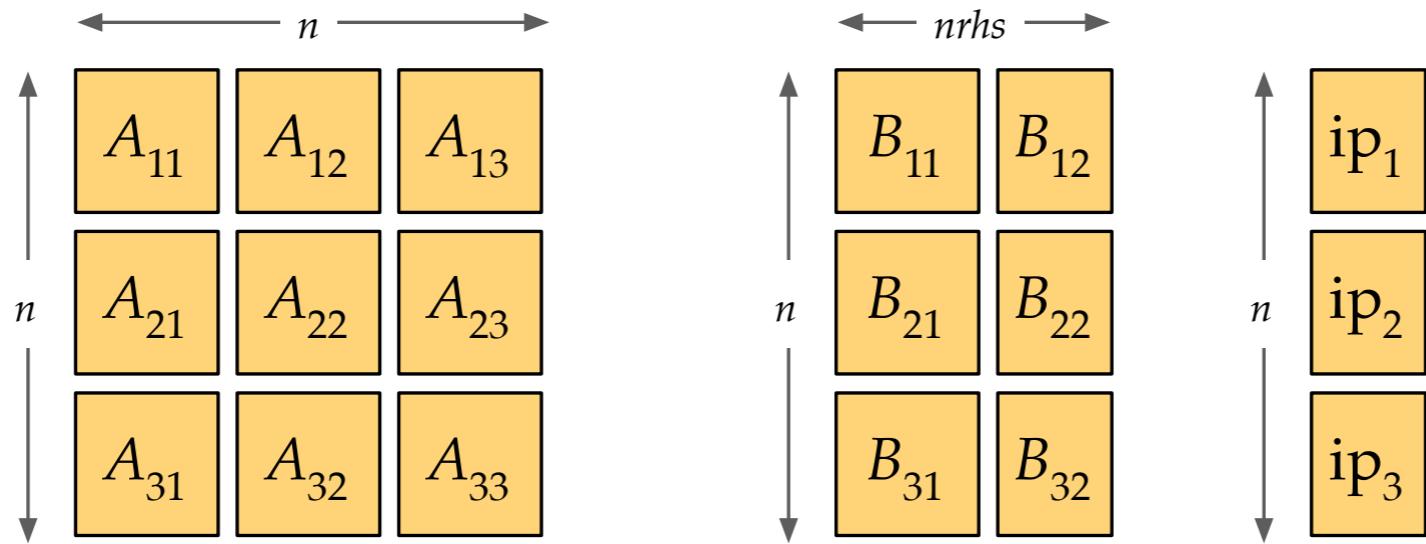


ScaLAPACK structure



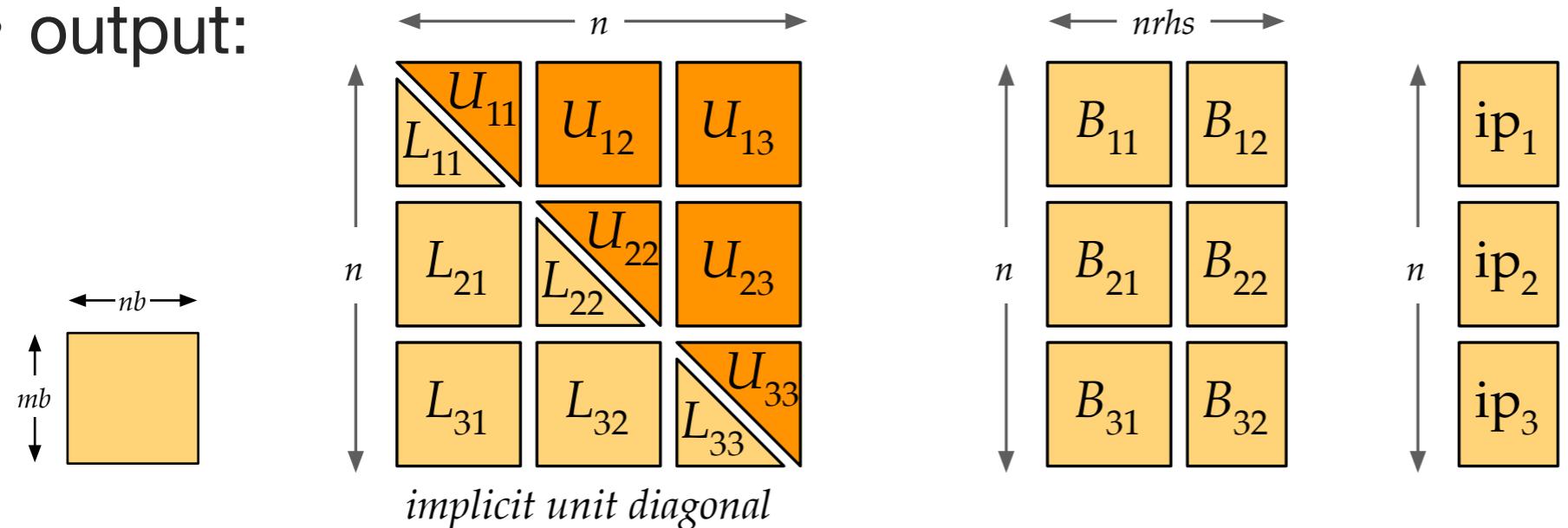
ScaLAPACK routine, solve $AX = B$

- LAPACK: `dgesv(n, nrhs, A, lda, ipiv, B, ldb, info)`
- ScaLAPACK: `pdgesv(n, nrhs, A, ia, ja, descA, ipiv, B, ib, jb, descB, info)`
- input:



**Global matrix
point of view**

- output:



info (error code)
= 0: no error
< 0: invalid argument
> 0: numerical error
(e.g., singular)

L, U overwrite A
 X overwrites B

2D block-cyclic layout

$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view

			n					
			14	15	16	17	18	
			24	25	26	27	28	
			34	35	36	37	38	
			44	45	46	47	48	
			54	55	56	57	58	
			64	65	66	67	68	
			74	75	76	77	78	
			84	85	86	87	88	
			94	95	96	97	98	
			51	52	53	61	62	63
			71	72	73	81	82	83
			91	92	93	94	95	96

Local process point of view

			q = 3 processes			
			Process 1, 1	Process 1, 2	Process 1, 3	
	11	14	17	12	15	18
	31	34	37	32	35	38
	51	54	57	52	55	58
	71	74	77	72	75	78
	91	94	97	92	95	98
			Process 2, 1	Process 2, 2	Process 2, 3	
	21	24	27	22	25	28
	41	44	47	42	45	48
	61	64	67	62	65	68
	81	84	87	82	85	88

2D block-cyclic layout

$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view

			n		
			11	12	13
			21	22	23
			31	32	33
			44	42	43
			51	52	53
			61	62	63
			71	72	73
			81	82	83
			91	92	93
			14	15	16
			24	25	26
			34	35	36
			44	45	46
			54	55	56
			64	65	66
			74	75	76
			84	85	86
			94	95	96
			17	18	
			27	28	
			37	38	
			47	48	
			57	58	
			67	68	
			77	78	
			87	88	
			97	98	

Local process point of view

			$q = 3$ processes		
			Process 1, 1	Process 1, 2	Process 1, 3
			11	14	17
			31	34	37
			51	54	57
			71	74	77
			91	94	97
			Process 2, 1	Process 2, 2	Process 2, 3
			21	24	27
			41	44	47
			61	64	67
			81	84	87
			Process 2, 1	Process 2, 2	Process 2, 3
			21	24	27
			41	44	47
			61	64	67
			81	84	87
			21	24	27
			41	44	47
			61	64	67
			81	84	87

2D block-cyclic layout

$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view

			n		
			11	12	13
			21	22	23
			31	32	33
			44	42	43
			51	52	53
			61	62	63
			71	72	73
			81	82	83
			91	92	93
			14	15	16
			24	25	26
			34	35	36
			44	45	46
			54	55	56
			64	65	66
			74	75	76
			84	85	86
			94	95	96
			17	18	
			27	28	
			37	38	
			47	48	
			57	58	
			67	68	
			77	78	
			87	88	
			97	98	

Local process point of view

			$q = 3$ processes		
			Process 1, 1	Process 1, 2	Process 1, 3
			11	14	17
			31	34	37
			51	54	57
			71	74	77
			91	94	97
			Process 2, 1	Process 2, 2	Process 2, 3
			21	24	27
			41	44	47
			61	64	67
			81	84	87
			22	25	28
			42	45	48
			62	65	68
			82	85	88
			23	26	
			43	46	
			63	66	
			83	86	

2D block-cyclic layout

$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view

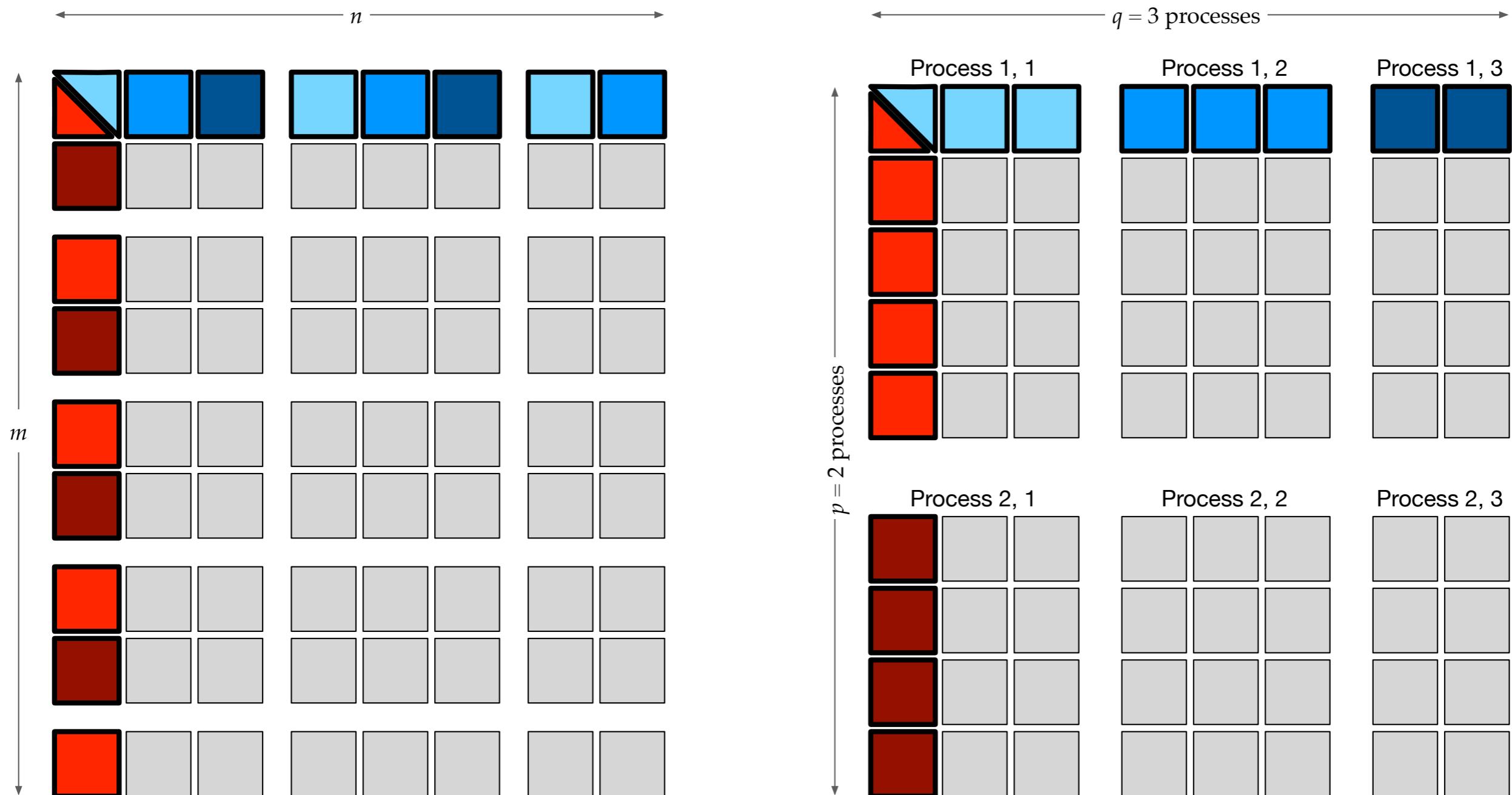
			n		
			11	12	13
m	14	15	16	17	18
	21	22	23	27	28
	31	32	33	37	38
	44	42	43	47	48
	51	52	53	57	58
	61	62	63	67	68
	71	72	73	77	78
	81	82	83	87	88
	91	92	93	97	98
	94	95	96		

Local process point of view

$q = 3$ processes						
Process 1, 1			Process 1, 2			
Process 1, 1	11	14	17	12	15	18
	31	34	37	32	35	38
	51	54	57	52	55	58
	71	74	77	72	75	78
	91	94	97	92	95	98
Process 2, 1			Process 2, 2			
Process 2, 1	21	24	27	22	25	28
	41	44	47	42	45	48
	61	64	67	62	65	68
	81	84	87	82	85	88
Process 2, 3						

LU factorization

- Bold is LU panel (reds) and triangular solve (blues)
- Gray is trailing matrix update

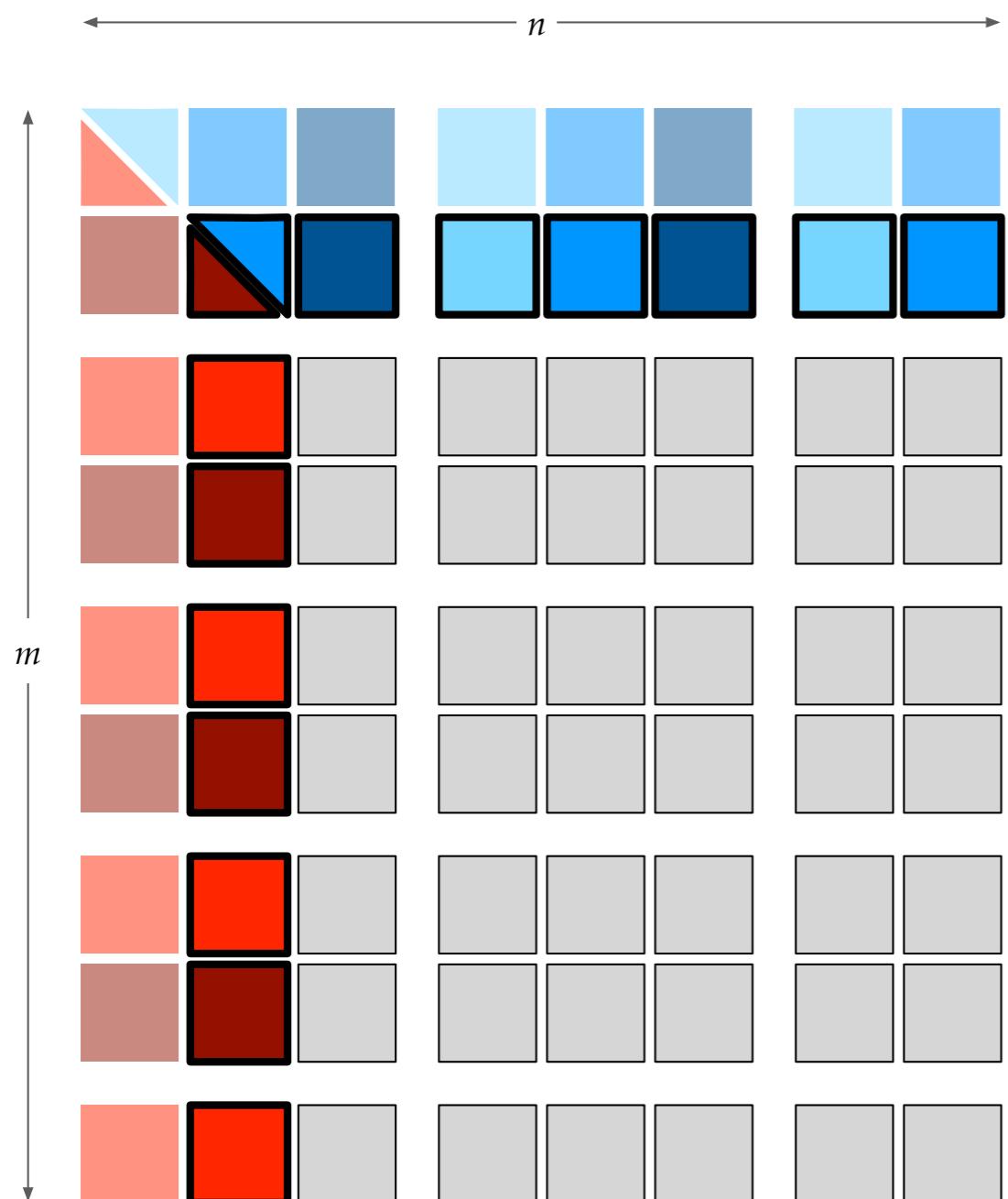


LU factorization

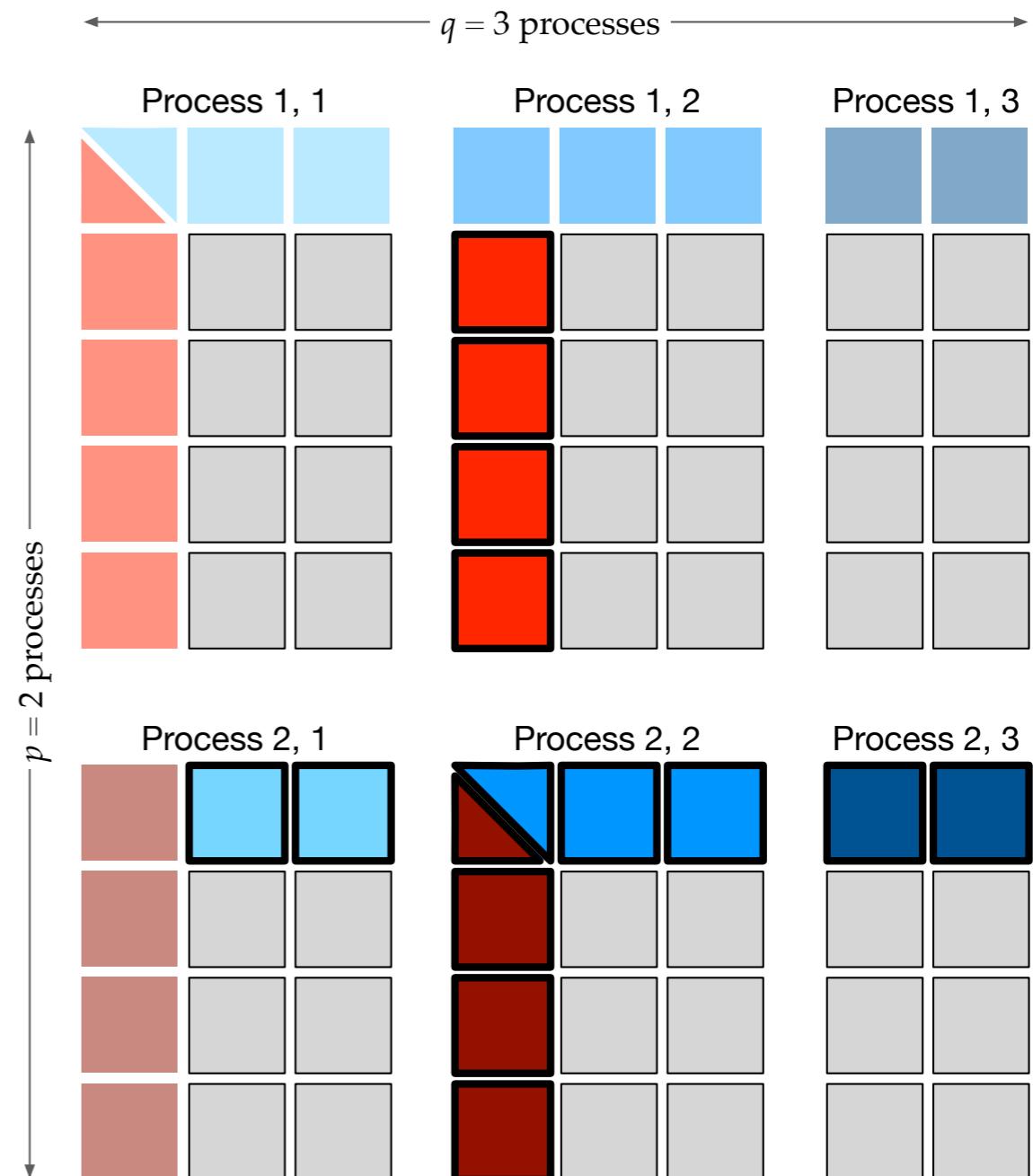
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

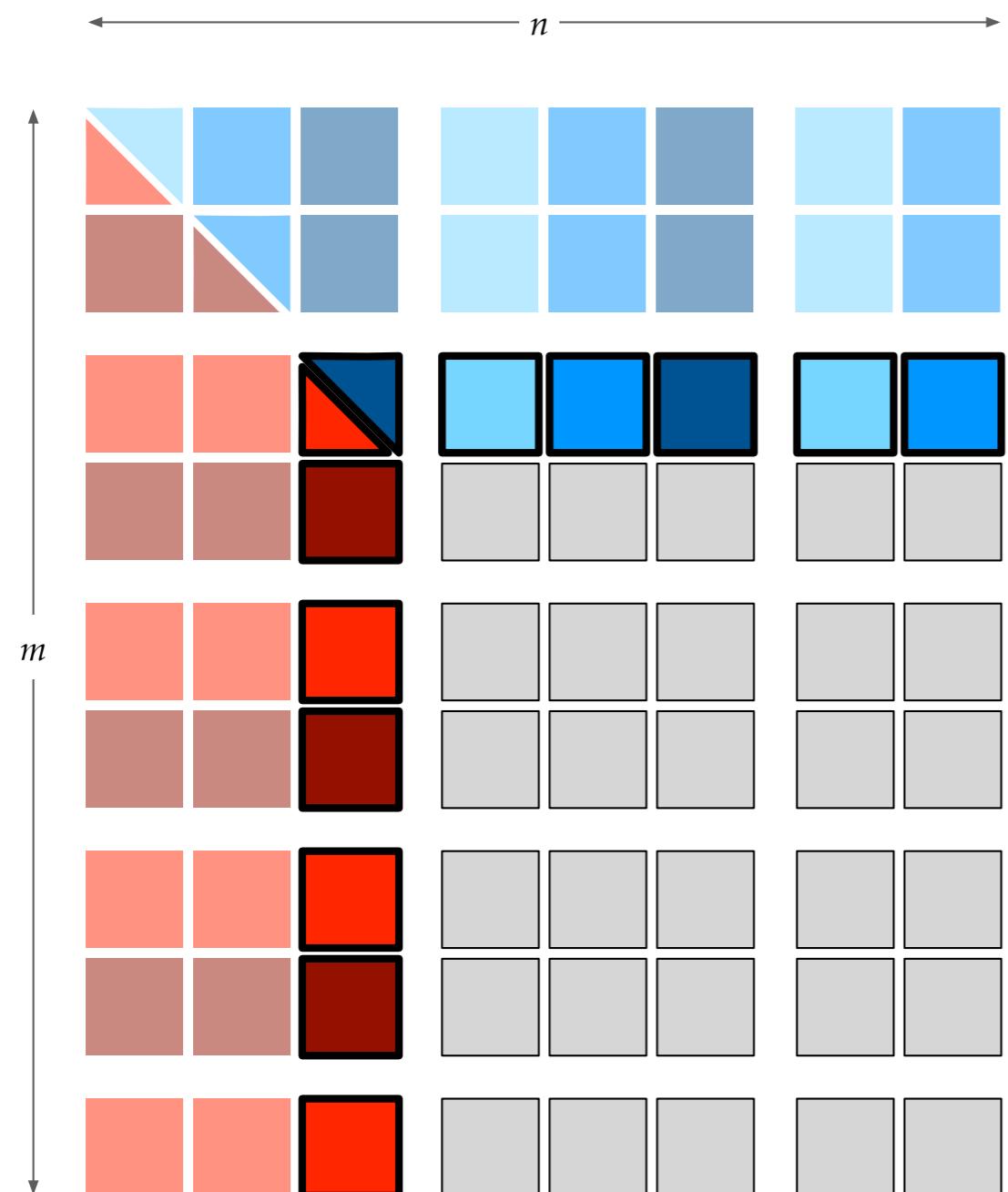


LU factorization

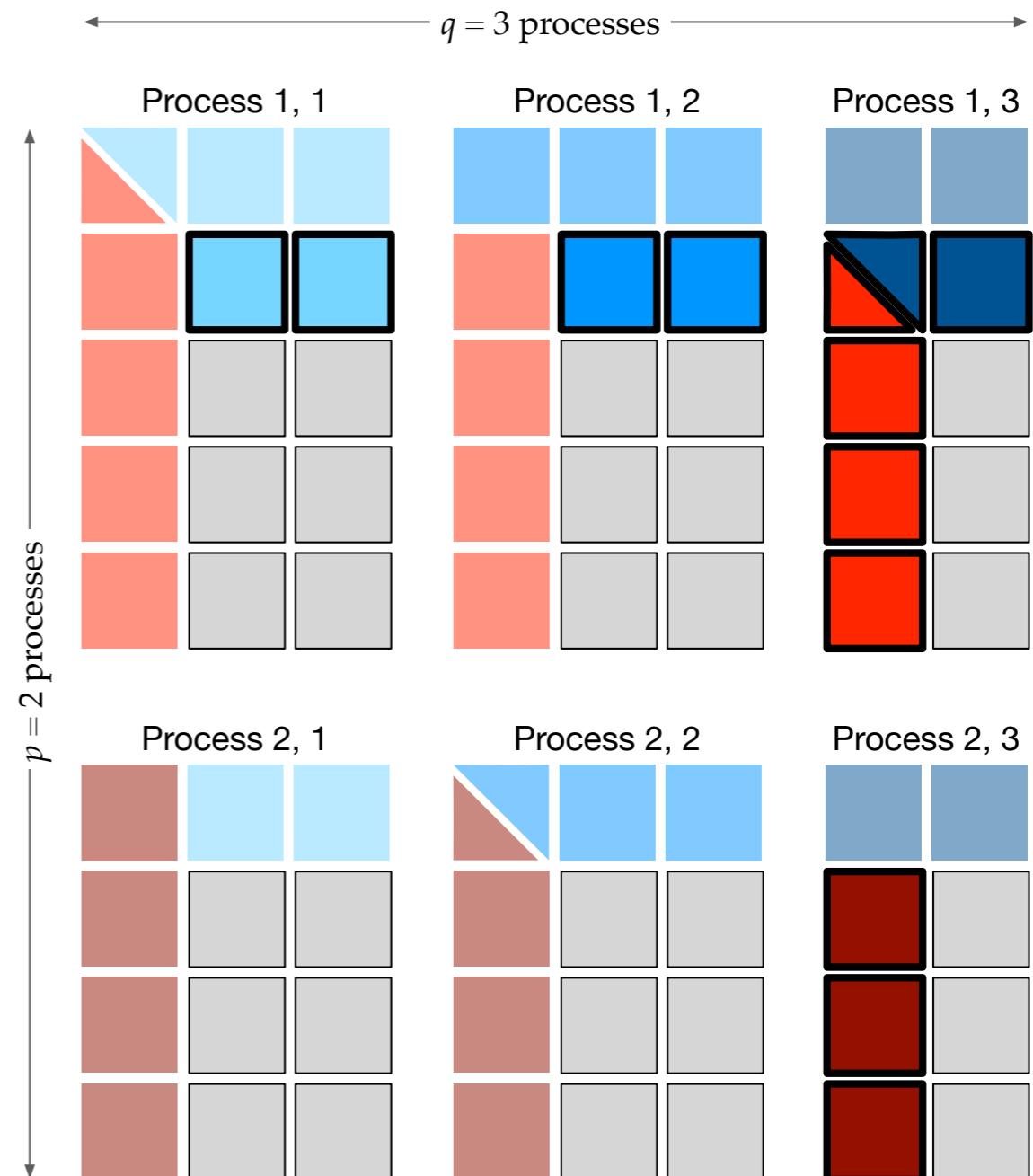
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

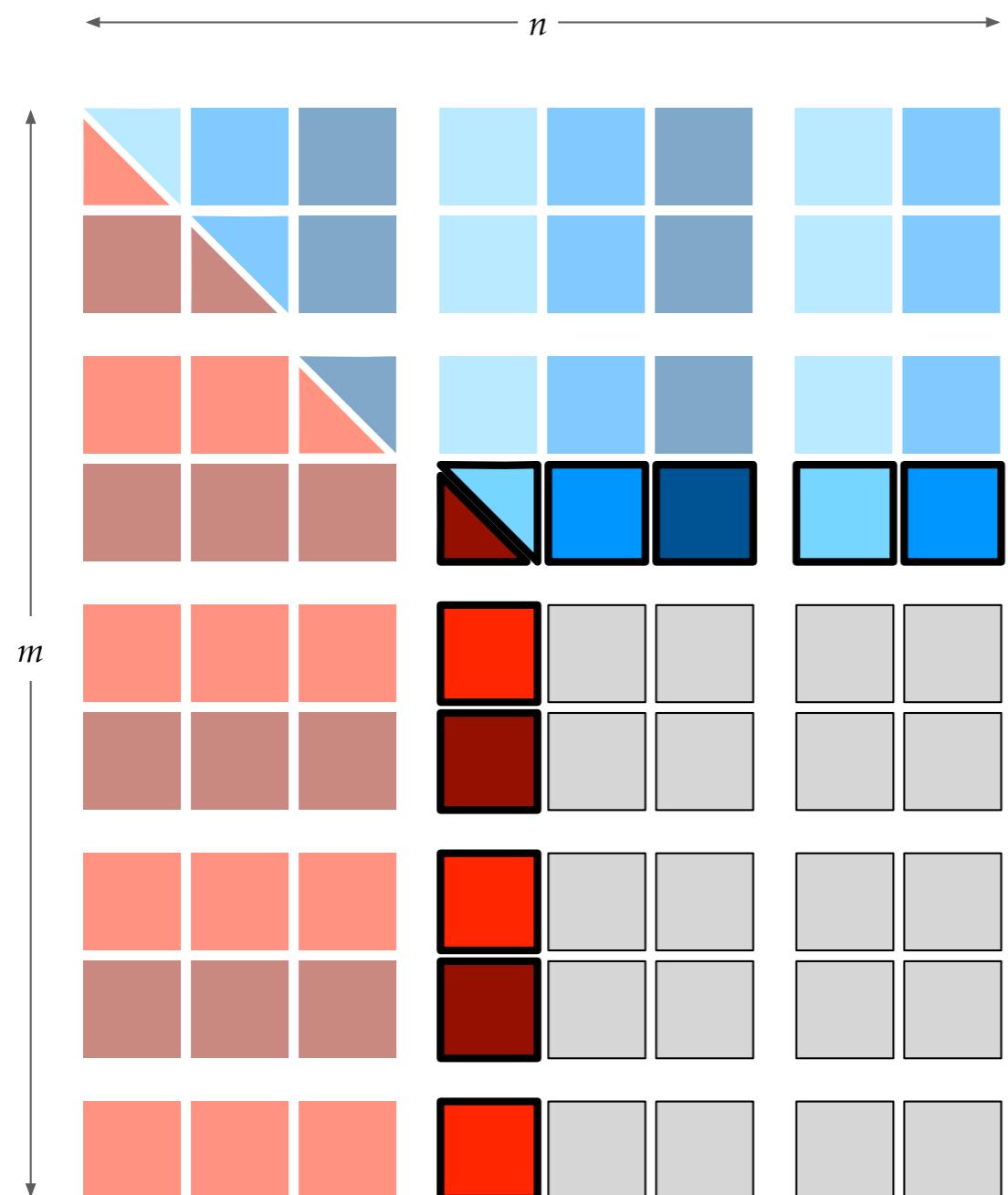


LU factorization

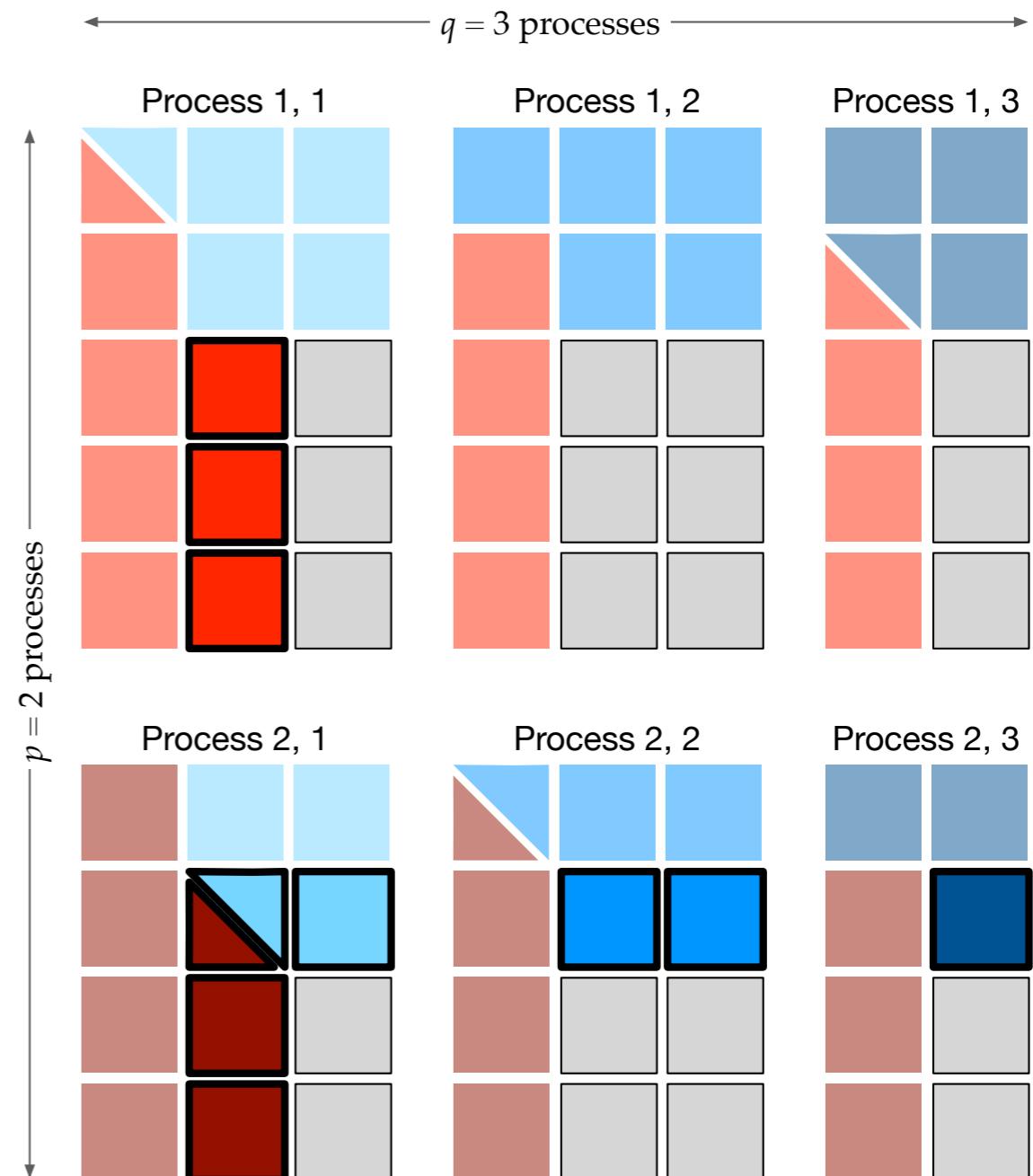
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

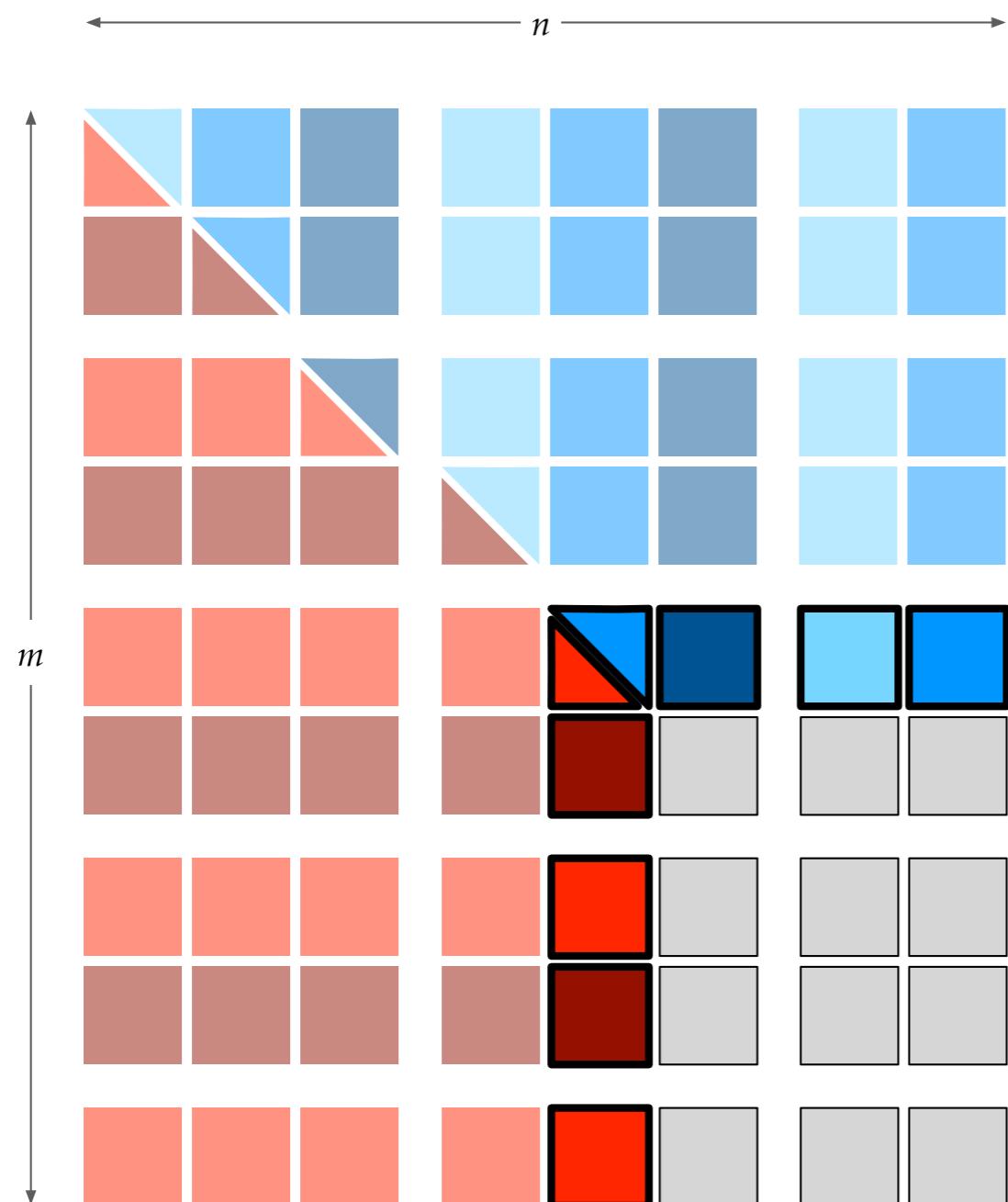


LU factorization

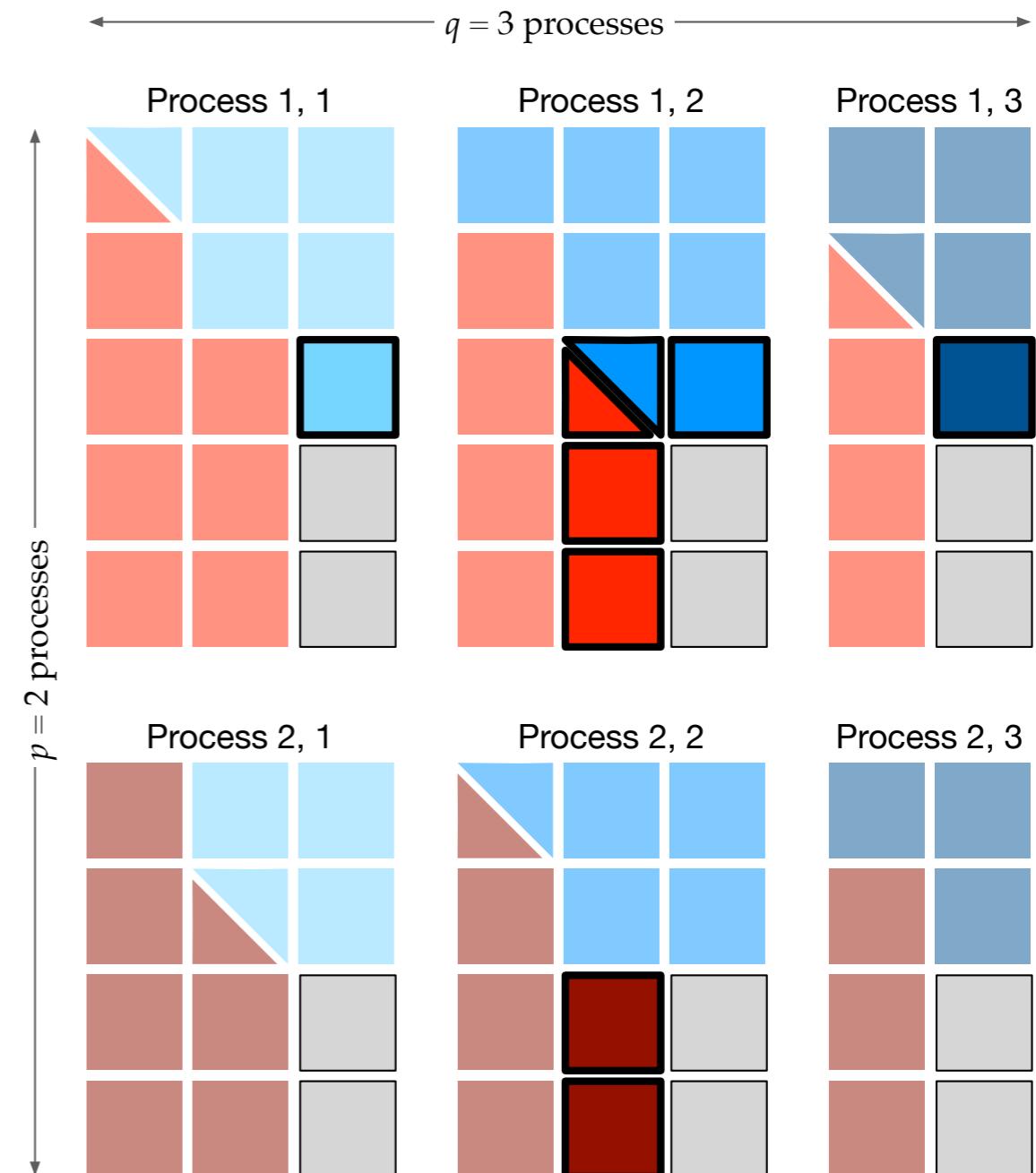
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

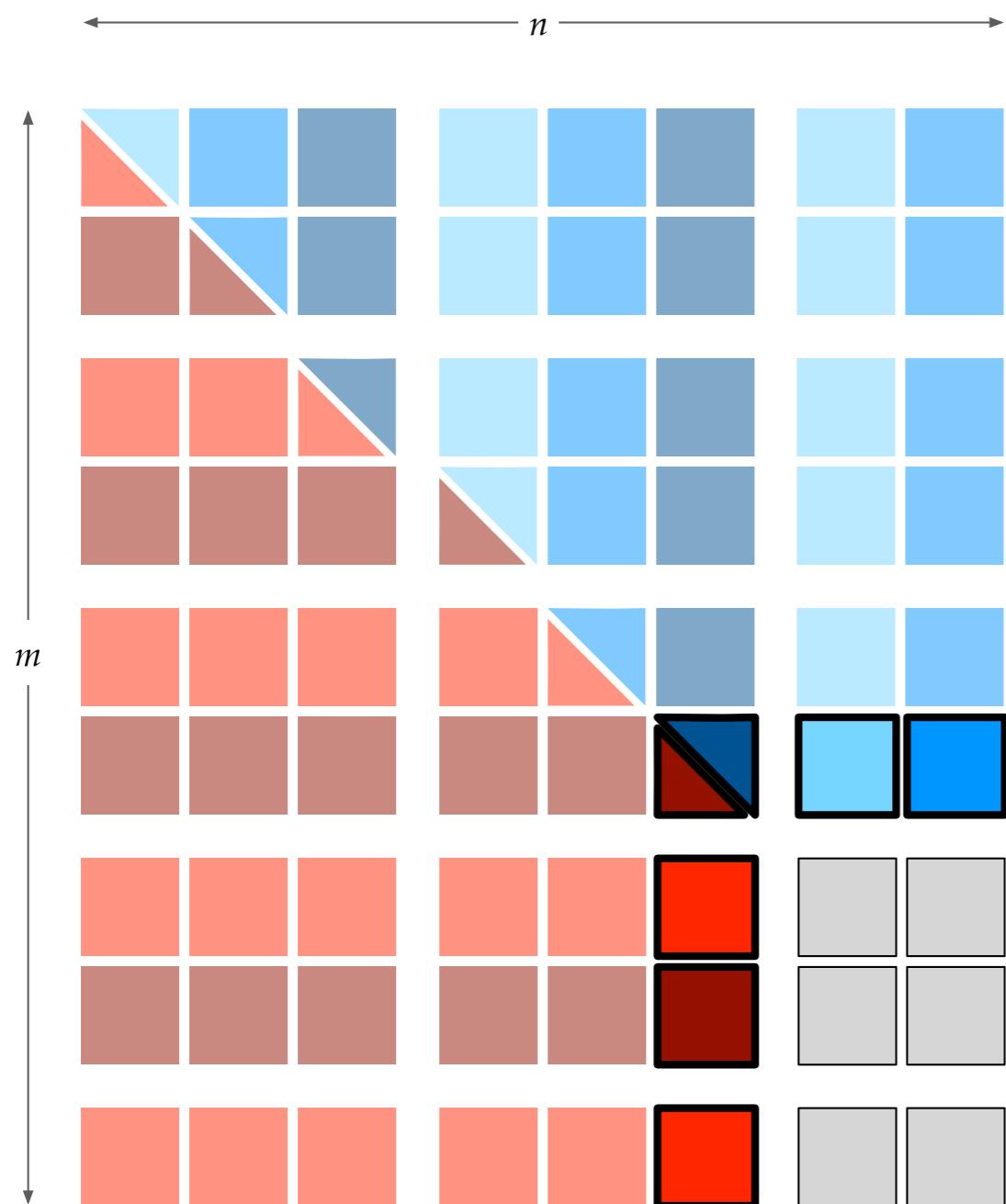


LU factorization

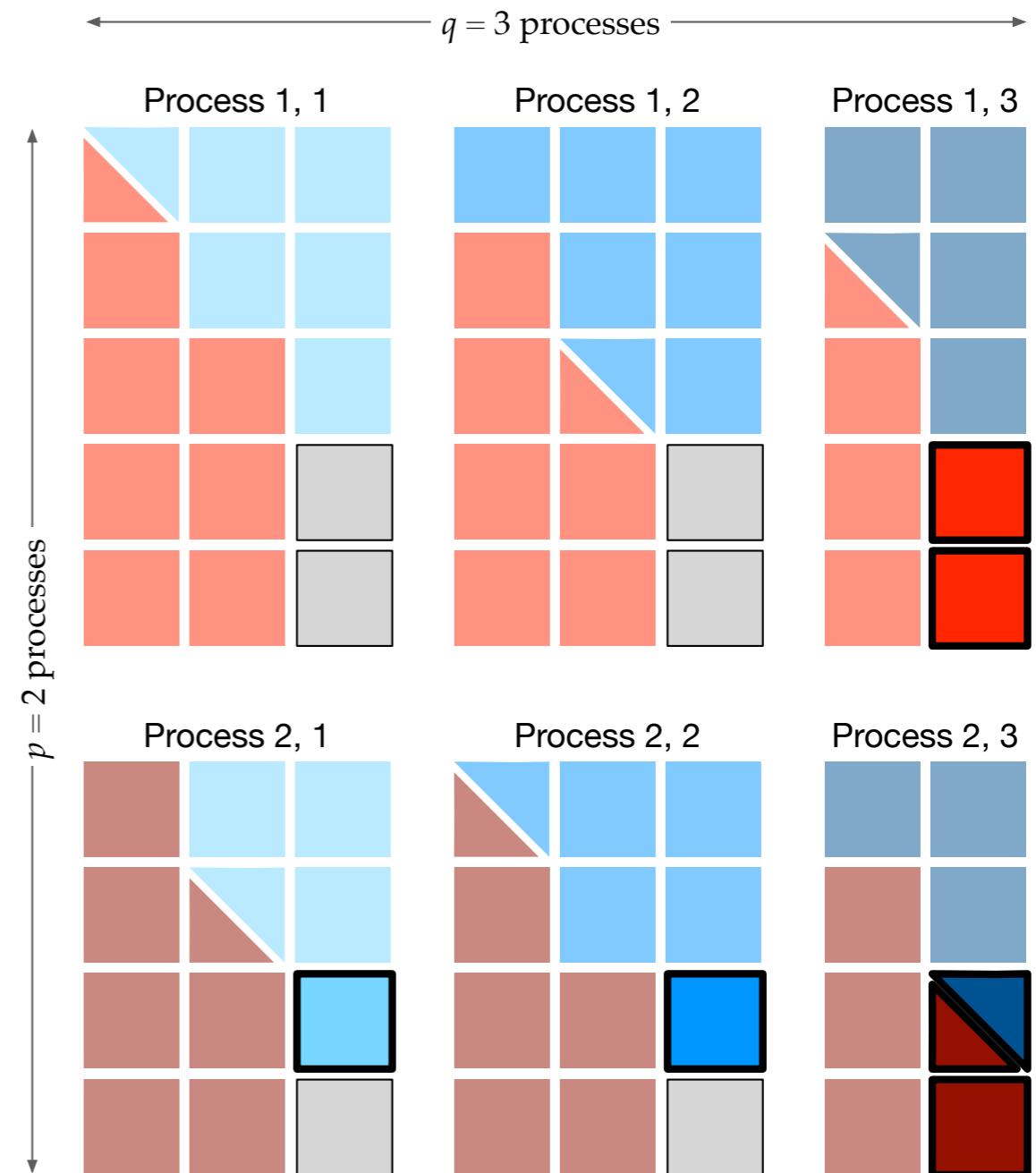
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

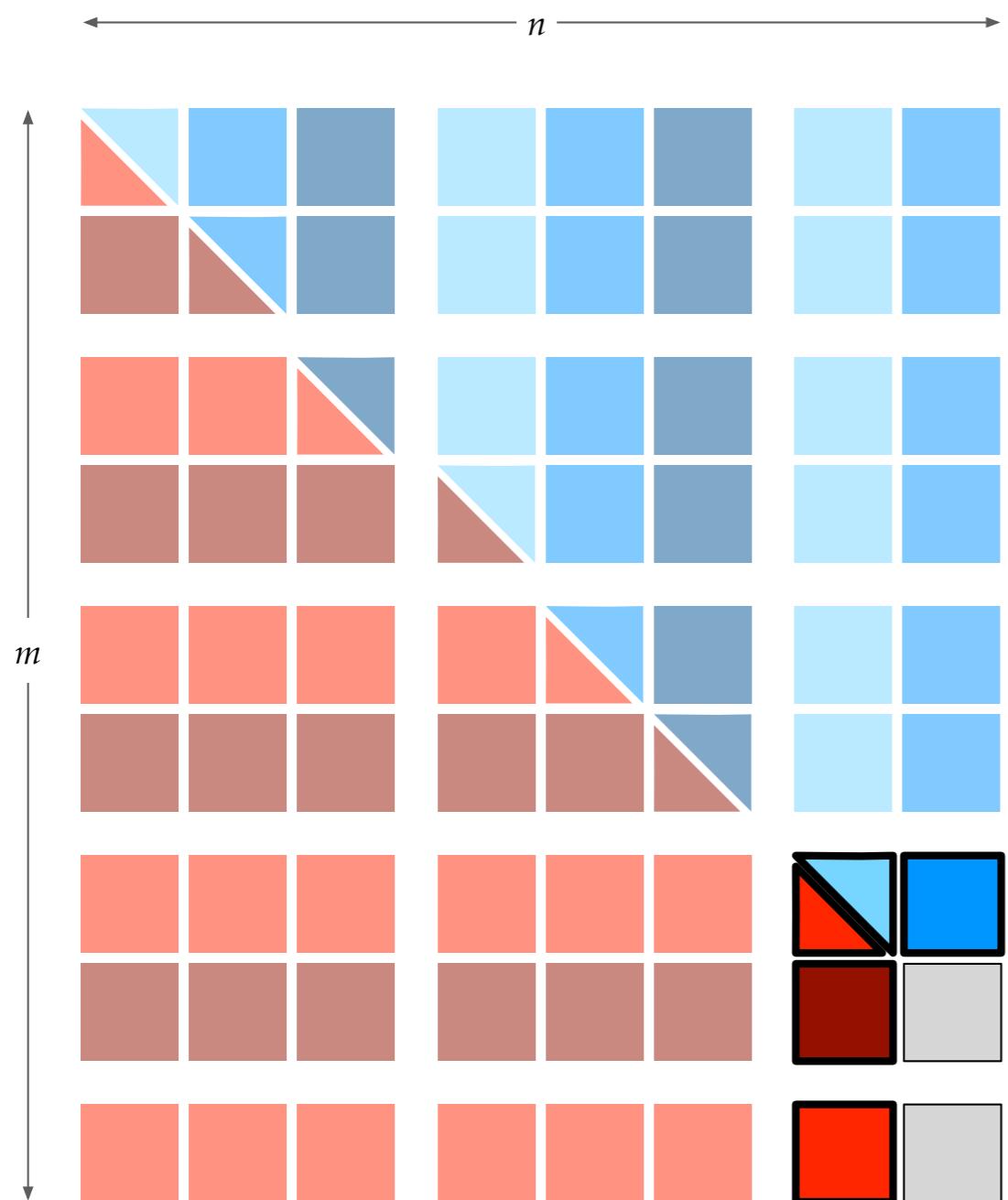


LU factorization

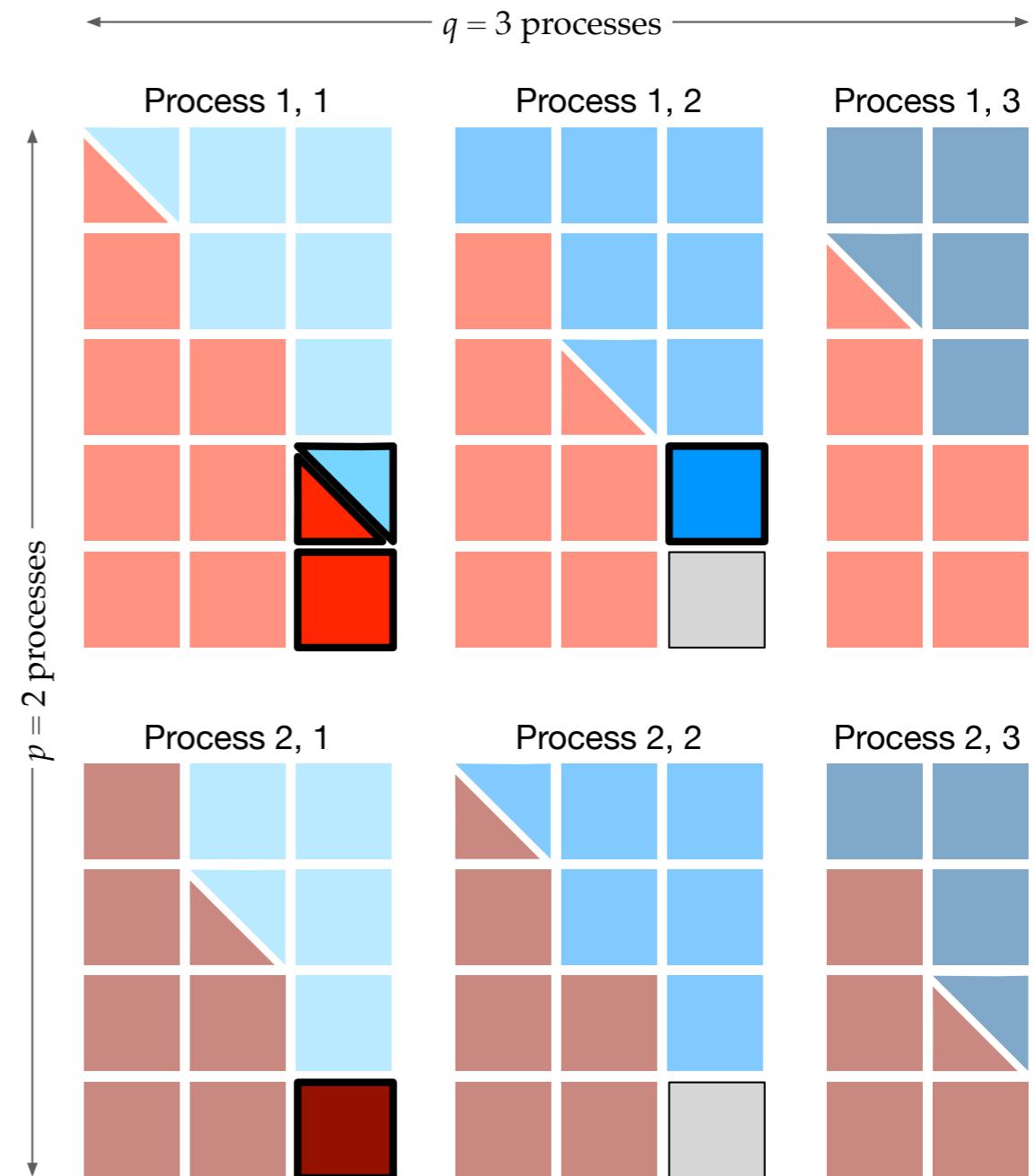
m × *n* matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

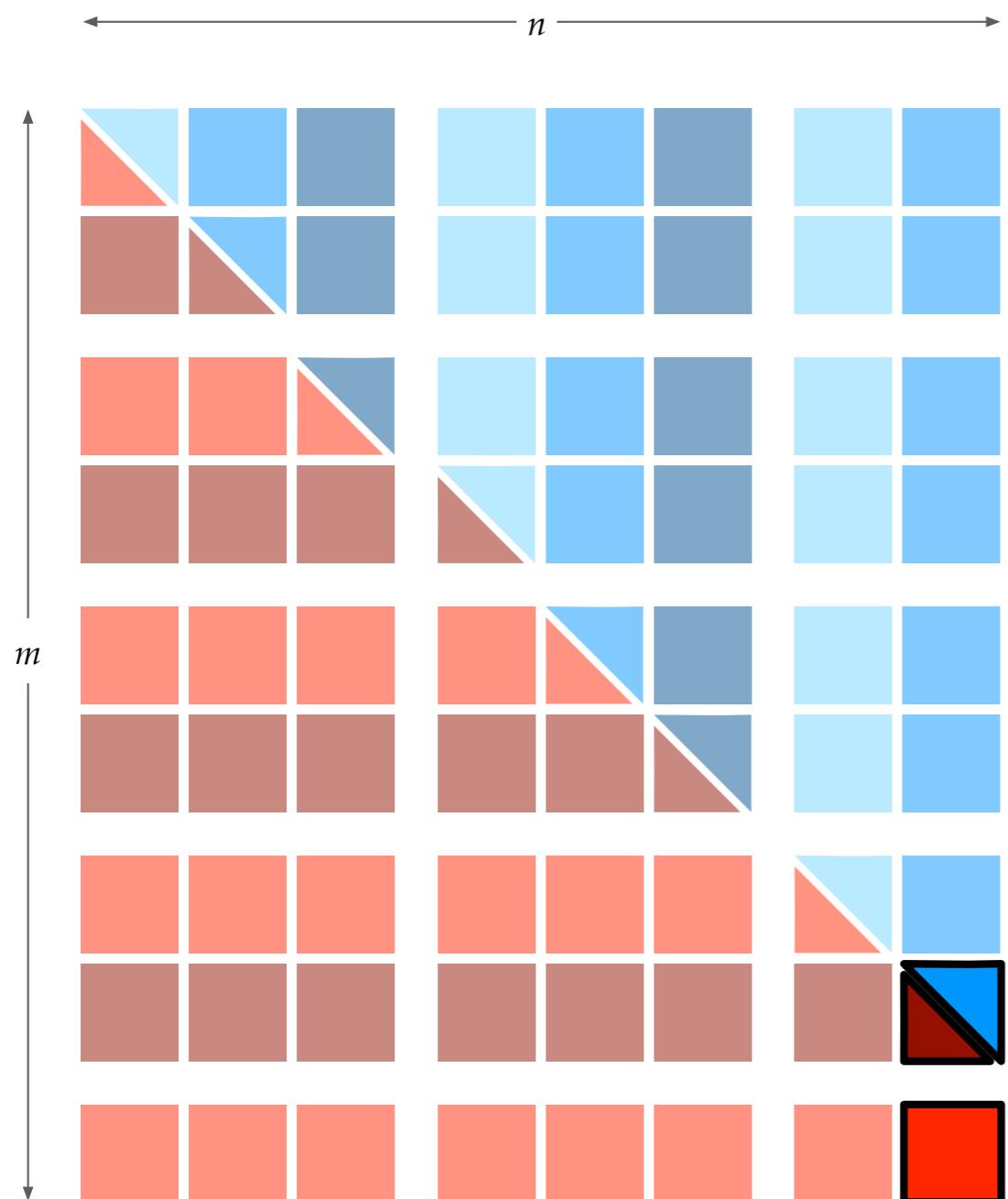


LU factorization

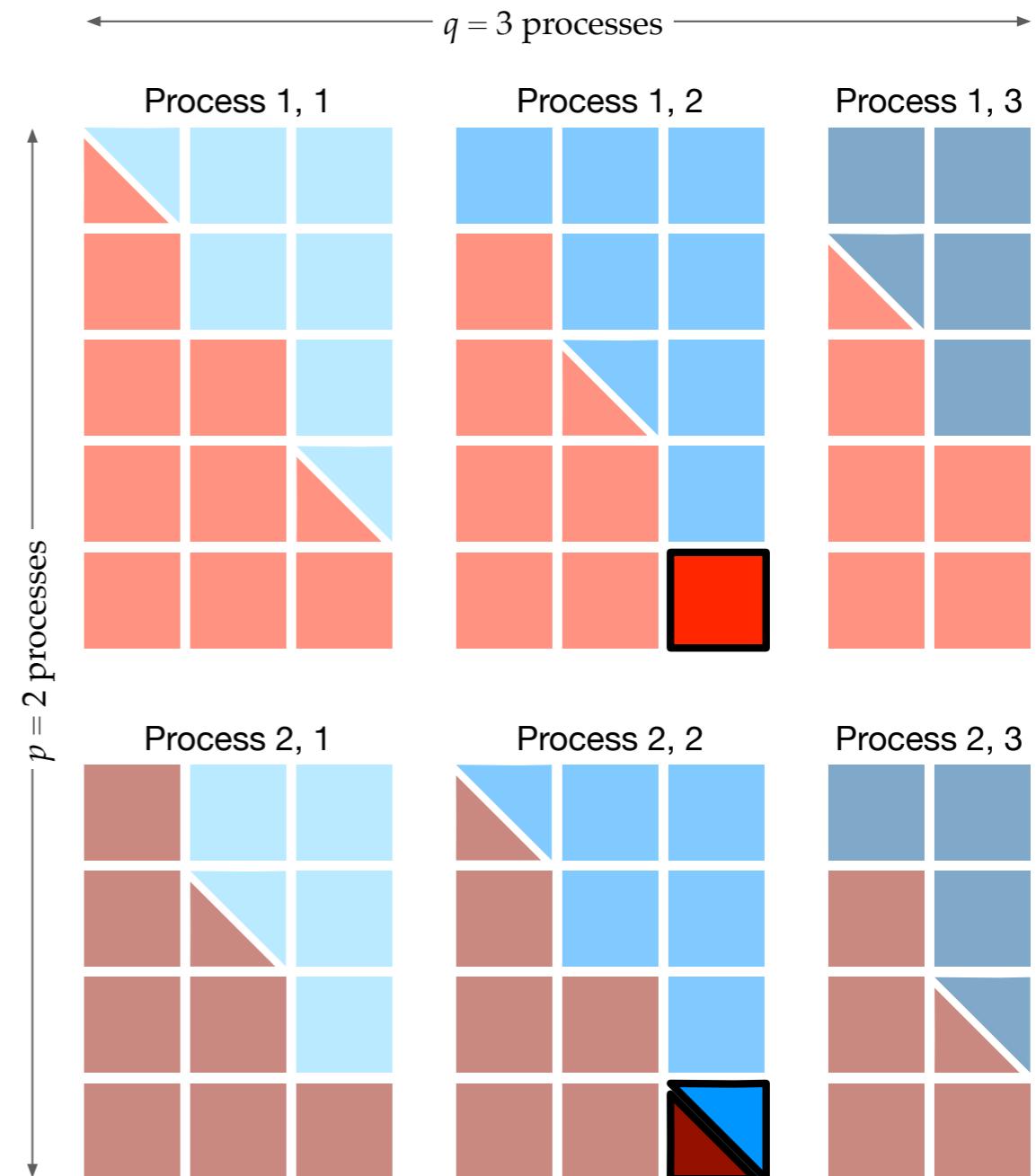
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view

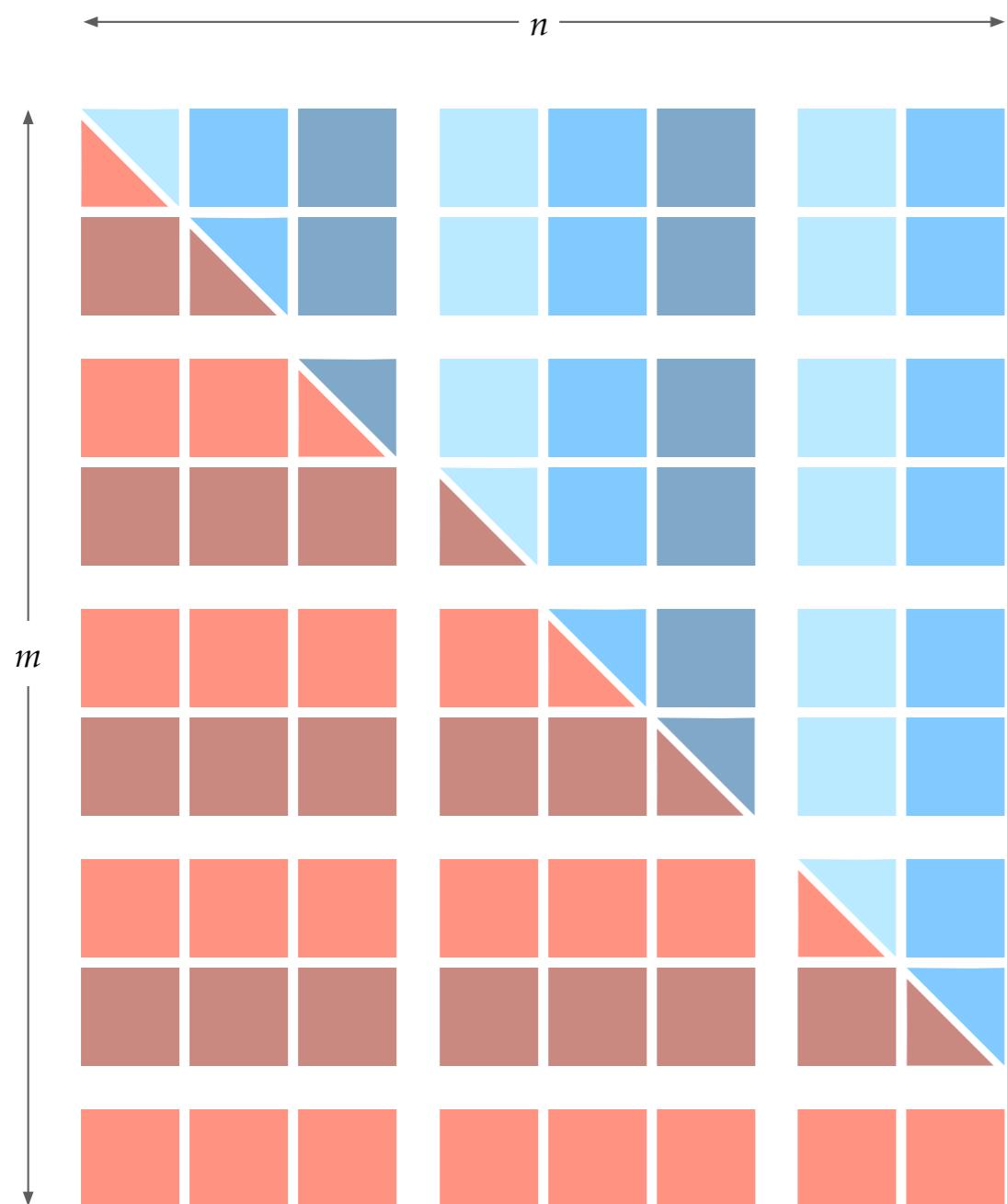


LU factorization

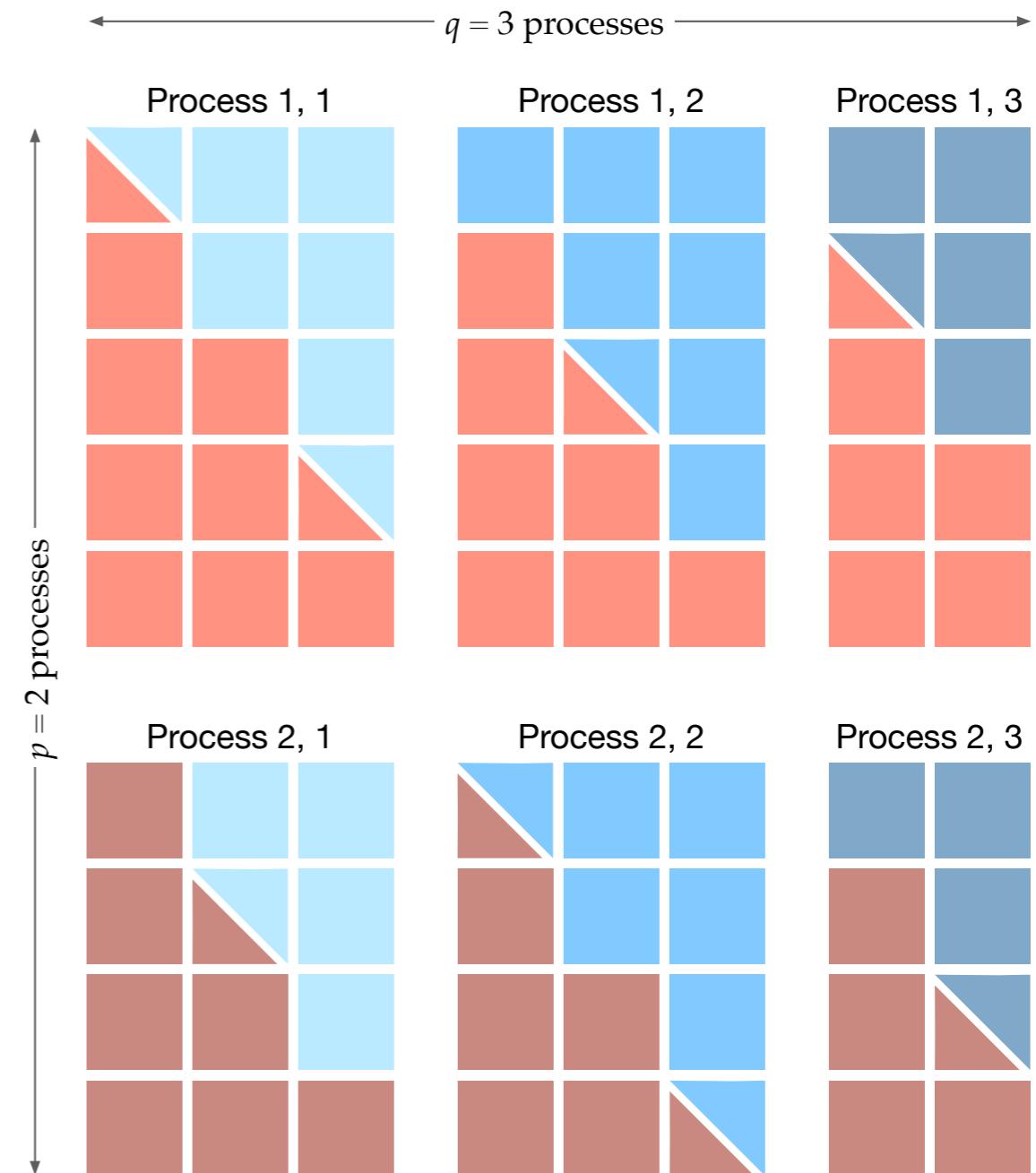
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view



Local process point of view



Why 2D block cyclic?

- Why not simple 2D distribution?

$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view

n			
m	11	12	13
	21	22	23
	31	32	33
	44	42	43
	51	52	53
	61	62	63
	64	65	66
	74	75	76
	84	85	86
	94	95	96

Local process point of view

q processes				
Process 1, 1	Process 1, 2	Process 1, 3		
11	12	13		
21	22	23		
31	32	33		
44	42	43		
51	52	53		
61	62	63		
71	72	73		
81	82	83		
91	92	93		
processes	Process 2, 1	Process 2, 2		
	61	62	63	
	71	72	73	
	81	82	83	
	91	92	93	
	Process 2, 3	67	68	
	64	65	66	
	74	75	76	
	84	85	86	
	94	95	96	
	77	78	87	88
	87	88	97	98

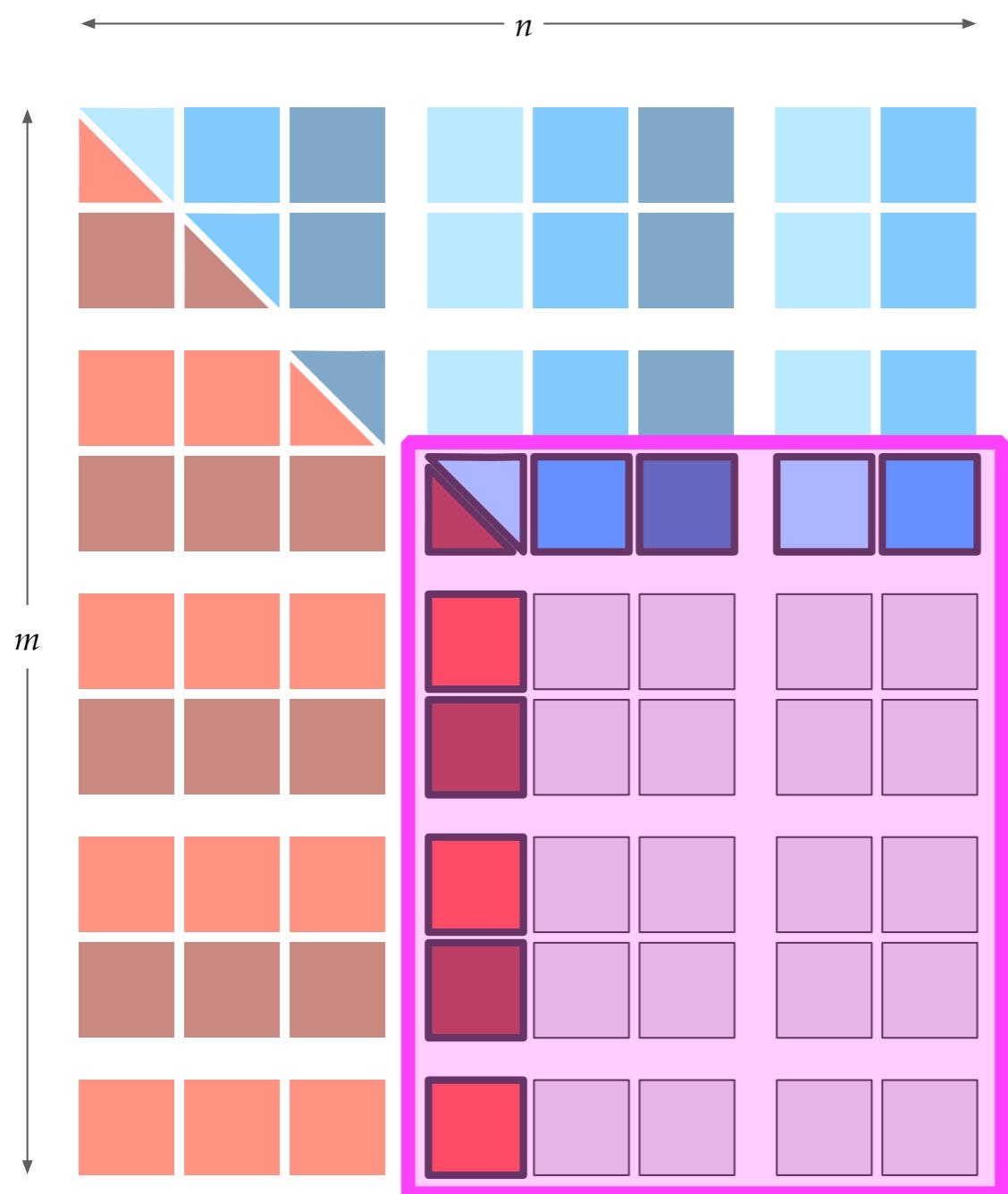
Why 2D block cyclic?

- Better load balancing!

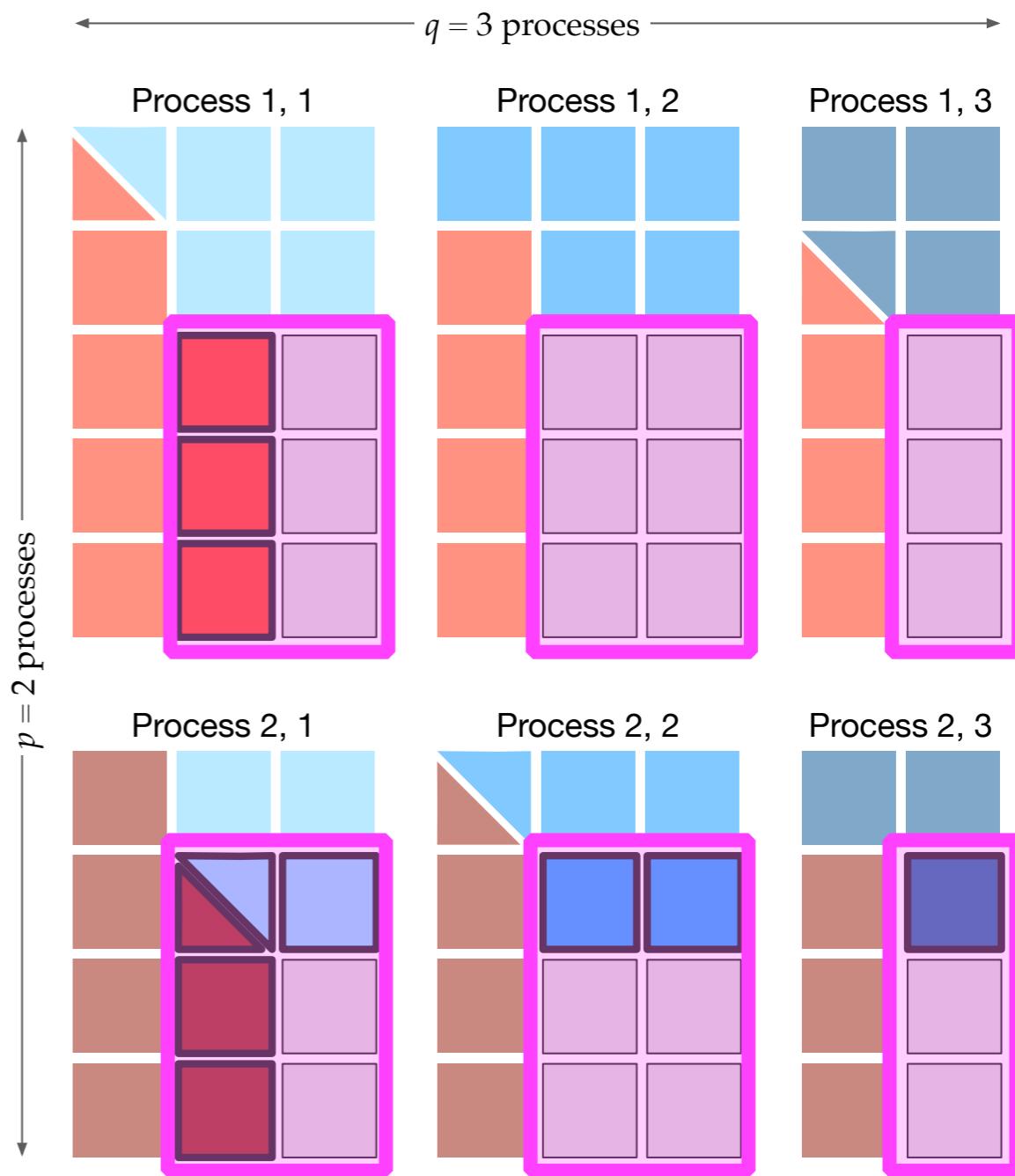
$m \times n$ matrix

$p \times q = 2 \times 3$ process grid

Global matrix view

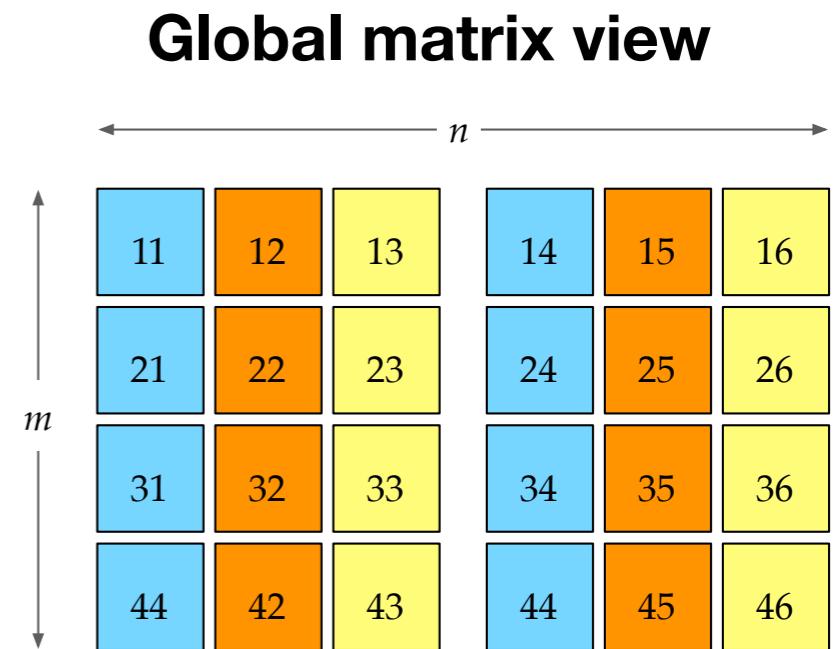


Local process point of view

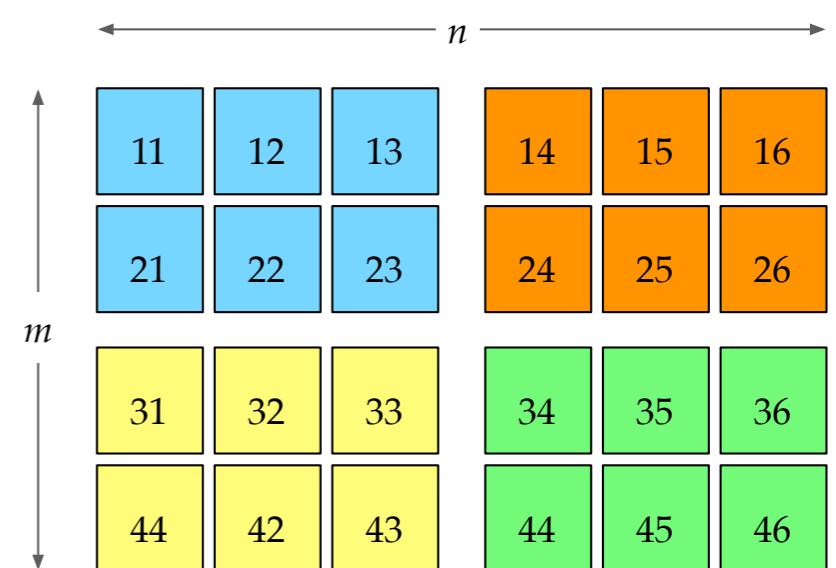


Other distributions

- 2D block cyclic not always required or best
- 1D block column cyclic (let $p = 1$) or
1D block row cyclic (let $q = 1$)
 - MAGMA and SLATE use for distributing across multiple GPUs within node



- 1D and 2D block (not cyclic)
 - Fine if matrix doesn't shrink during computation
 - gemm
 - Jacobi SVD

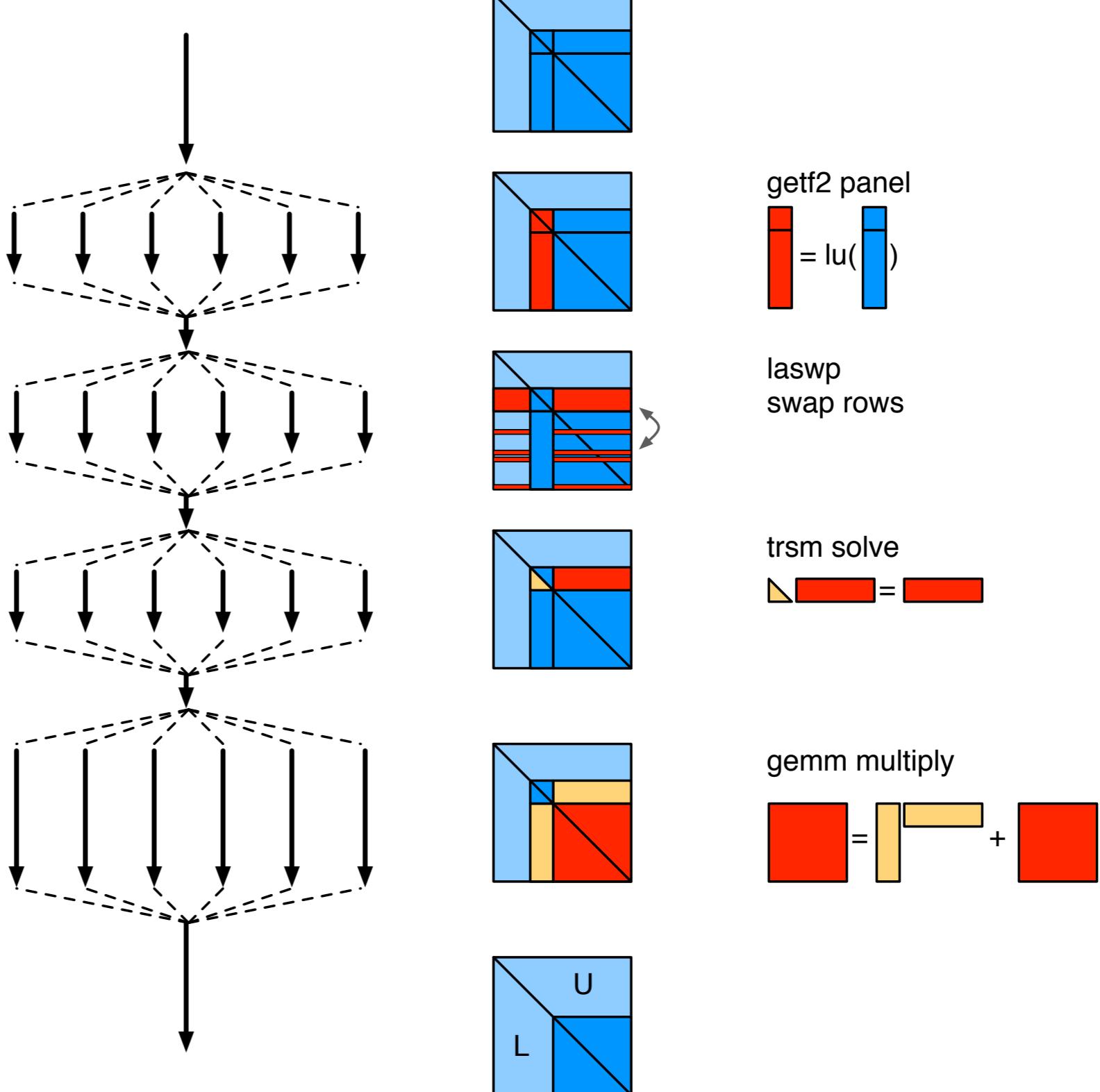


ScaLAPACK routines

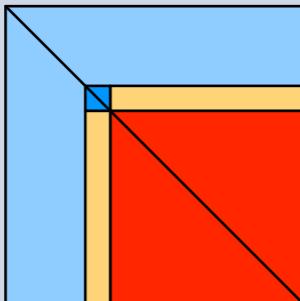
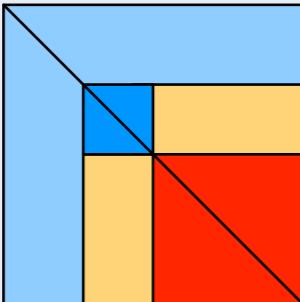
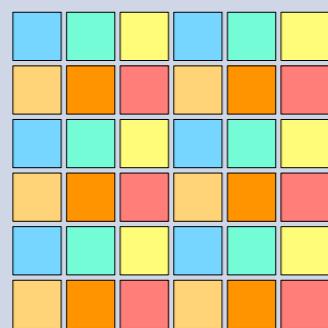
- ~140–160 routines x 4 data types (s, d, c, z)
- **Driver routines:** solve an entire problem
 - sv Solve linear system: $Ax = b$
 - gesv general non-symmetric (LU)
 - posv symmetric positive definite (Cholesky)
 - ~~sysv symmetric indefinite (LDL^T)~~
 - Also ~~packed~~, banded, tridiagonal storage
 - ls Linear least squares: $Ax \approx b$
 - ~~gglse linear equality-constrained least squares~~
 - ~~ggglm general Gauss-Markov linear model~~
 - ev Eigenvalue decomposition: $Ax = \lambda x$ and $Ax = \lambda Mx$
 - syevd symmetric and symmetric generalized
 - ~~geev~~ non-symmetric (only gehrd Hessenberg reduction, no geev)
 - svd Singular value decomposition: $A = U\Sigma V^H$
 - gesvd standard ~~and generalized~~ (also not optimized for tall matrices)
 - ~~gesdd D&C (faster)~~

Parallelism in ScaLAPACK

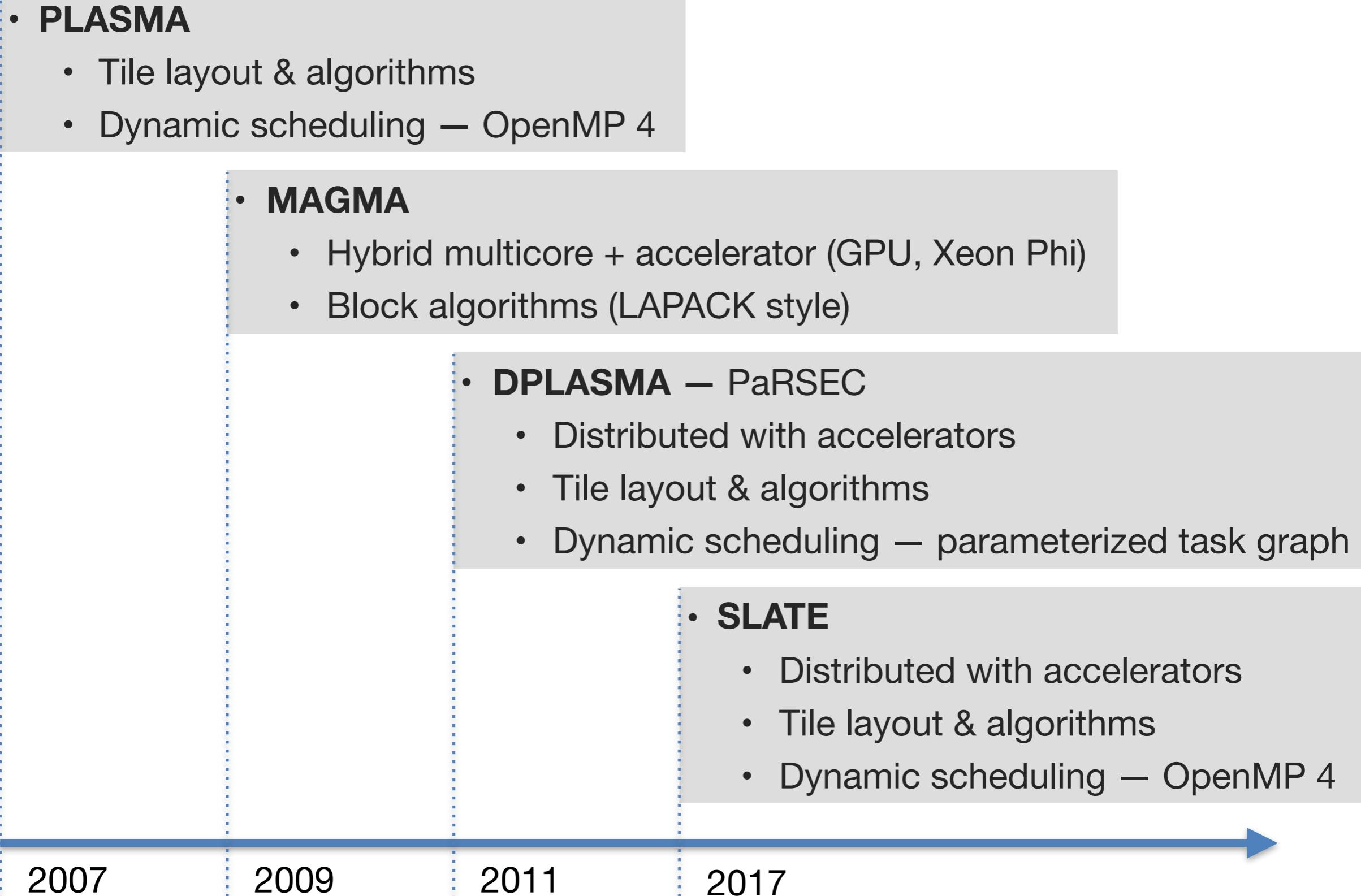
- Similar to LAPACK
- Bulk-synchronous
- Most flops in gemm update
 - $2/3 n^3$ term
 - Can use **sequential BLAS**,
 $p \times q = \# \text{ cores}$
= # MPI processes,
 $\text{num_threads} = 1$
 - Or **multi-threaded BLAS**,
 $p \times q = \# \text{ nodes}$
= # MPI processes,
 $\text{num_threads} = \# \text{ cores/node}$



Legacy software libraries

	<p>LINPACK (70s) vector operations</p>	Level 1 BLAS
	<p>LAPACK (80s) block operations</p>	Level 3 BLAS
	<p>ScaLAPACK (90s) 2D block cyclic distribution</p>	PBLAS BLACS MPI

Emerging software solutions

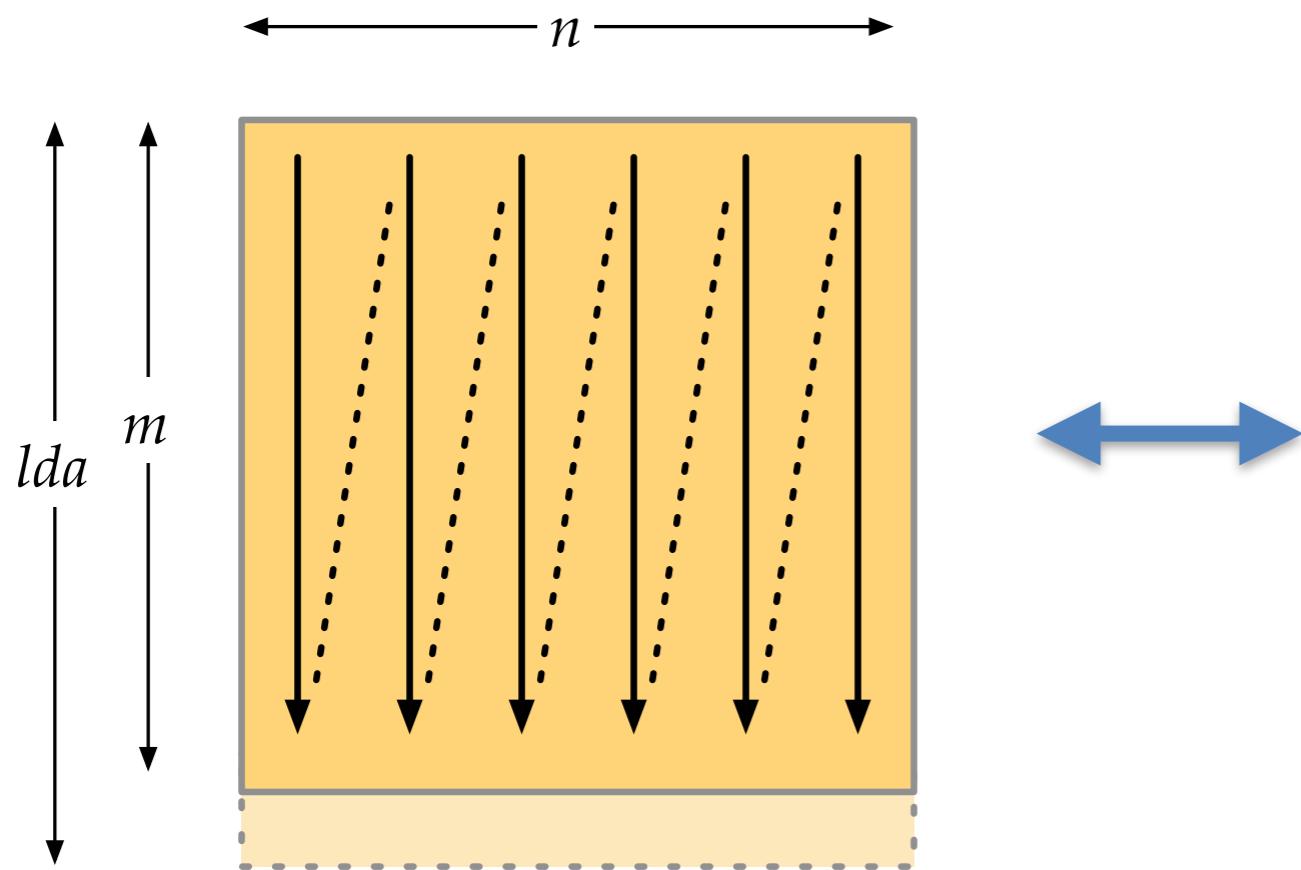


Outline

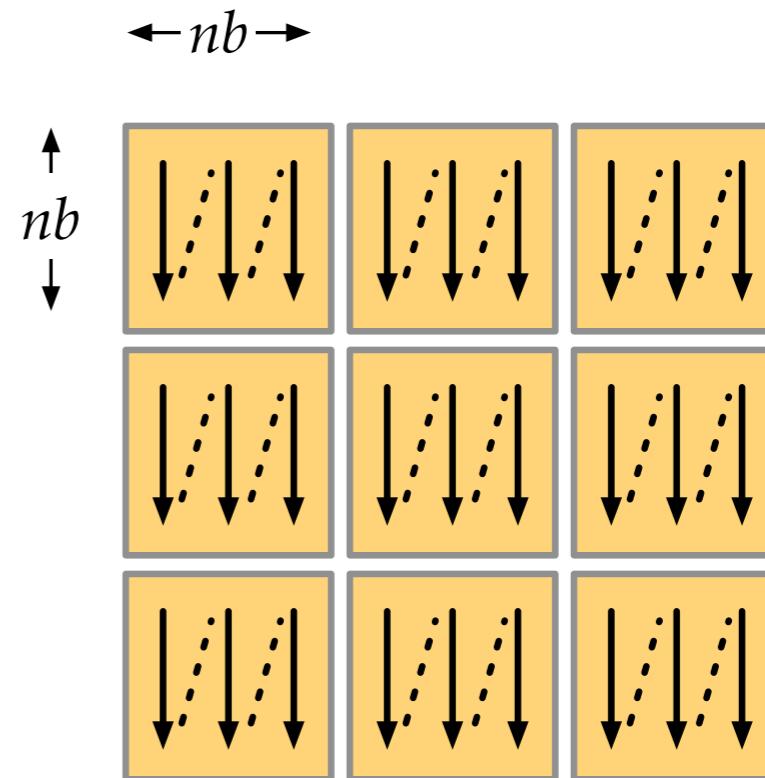
- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - **PLASMA using OpenMP**
 - DPLASMA and PaRSEC
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

Tile matrix layout

LAPACK column major



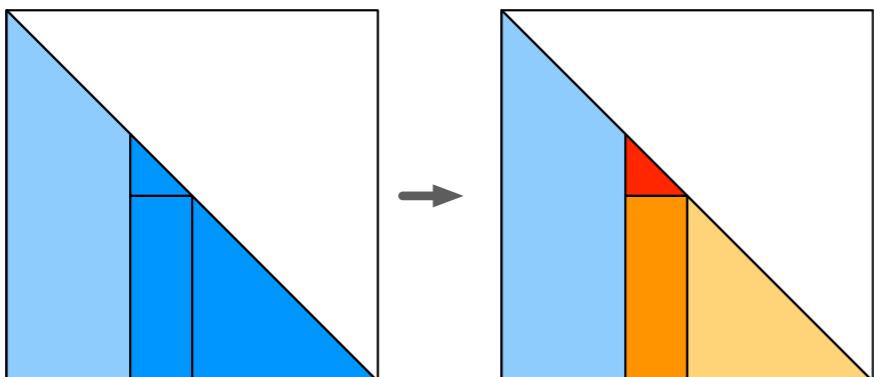
(D)PLASMA tile layout



- Tiled layout
 - Each tile is contiguous (column major)
 - Enables dataflow scheduling
 - Cache and TLB efficient (reduces conflict misses and false sharing)
 - MPI messaging efficiency (zero-copy communication)
 - In-place, parallel layout translation

Tile algorithms: Cholesky

LAPACK Algorithm (right looking)

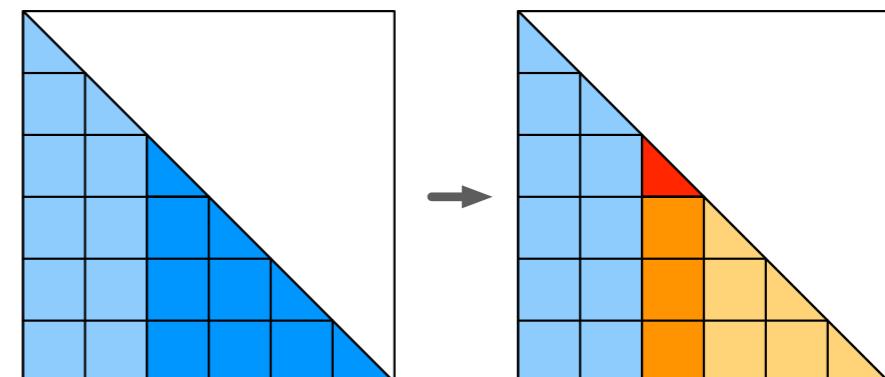


= $\text{chol}(\text{Blue Triangle})$

A = / trsm
B =
C =

= - $A^T \quad B^T \quad C^T$ herk

Tile Algorithm



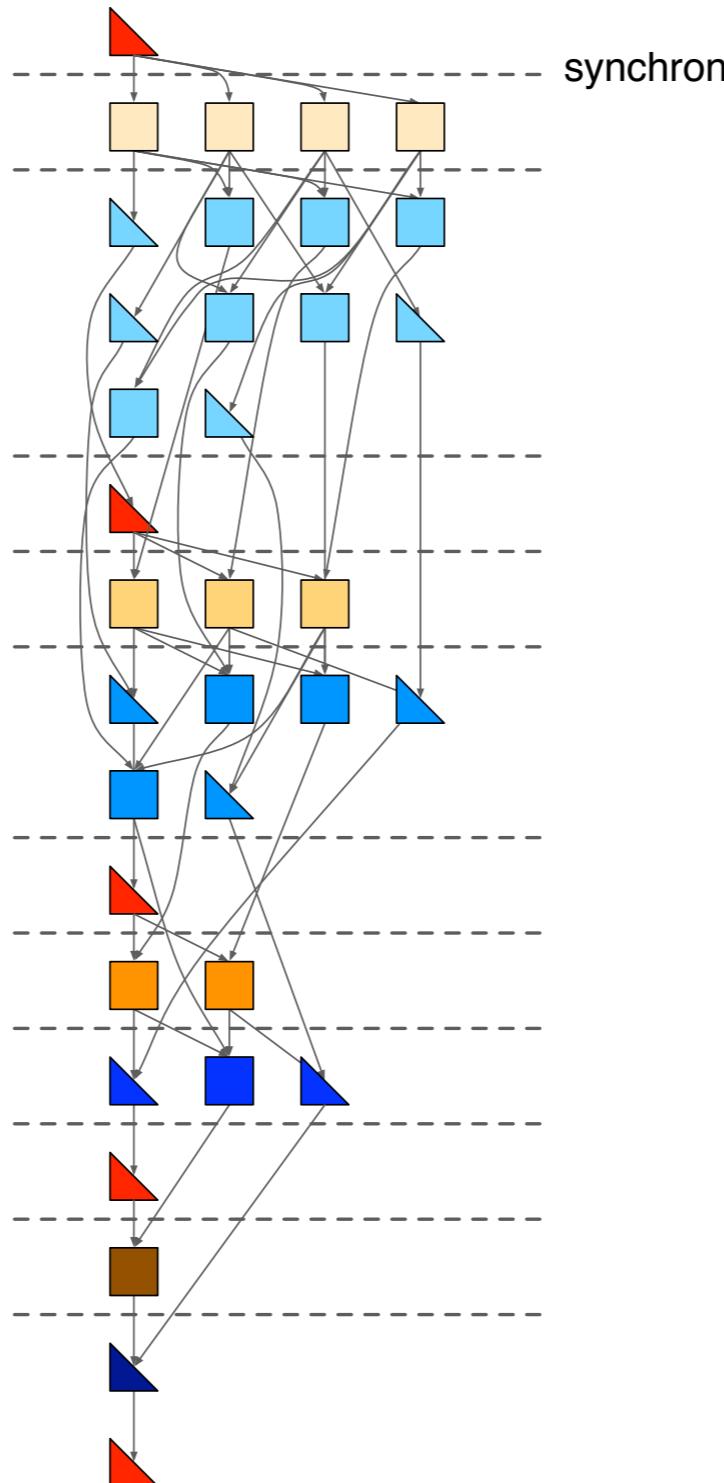
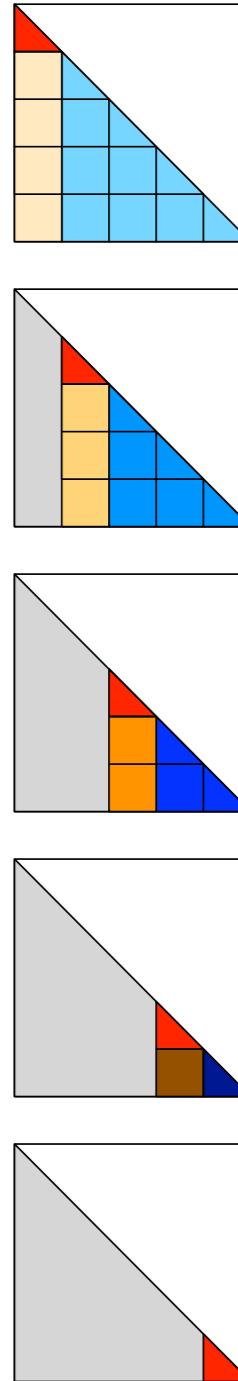
= $\text{chol}(\text{Blue Triangle})$

A = / trsm
 B = / trsm
 C = / trsm

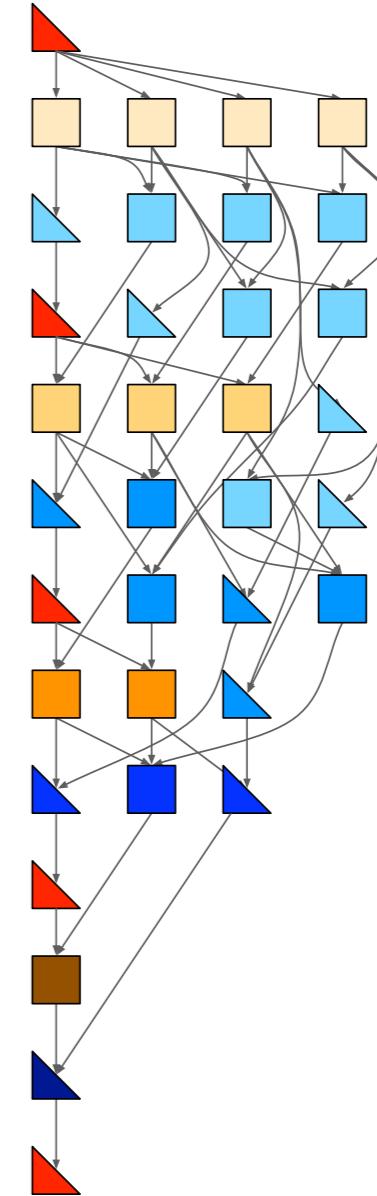
= - $A \quad A^T$ herk
 = - $B \quad A^T$ gemm
 = - $C \quad A^T$ gemm
 = - $B \quad B^T$ herk
 = - $B \quad C^T$ gemm
 = - $C \quad C^T$ herk

Track dependencies – Directed acyclic graph (DAG)

Fork-join schedule on 4 cores
with artificial synchronizations

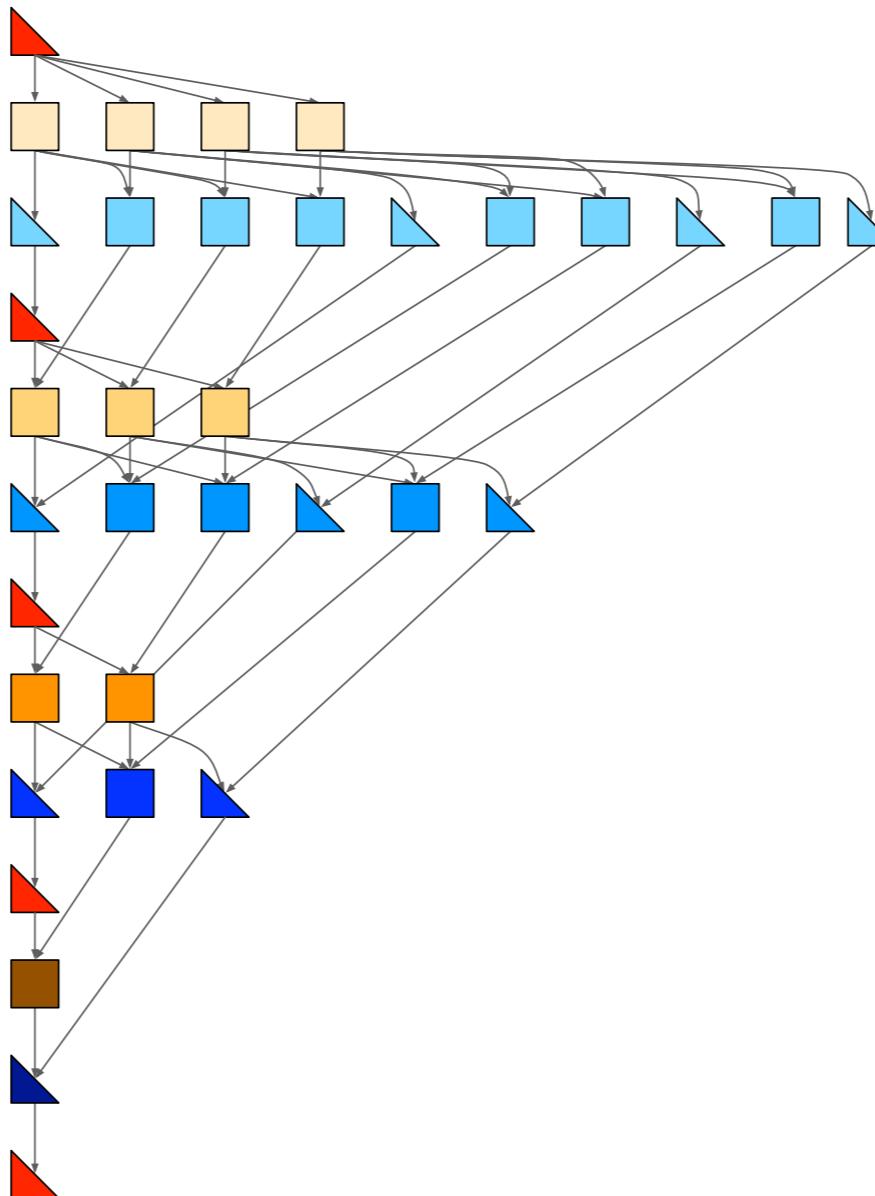
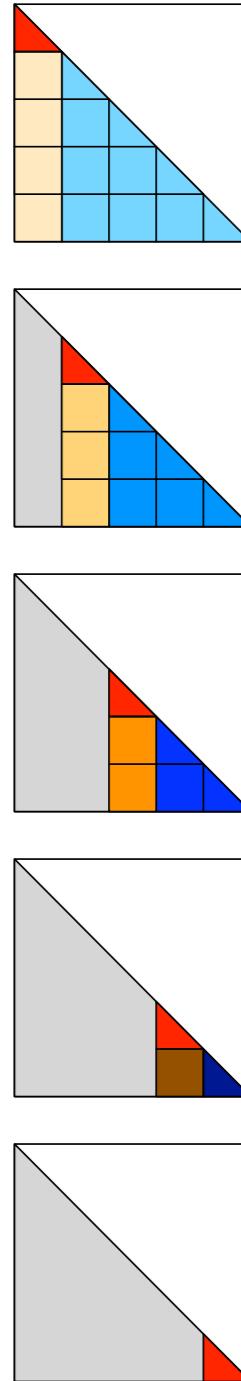


Reordered without
synchronizations



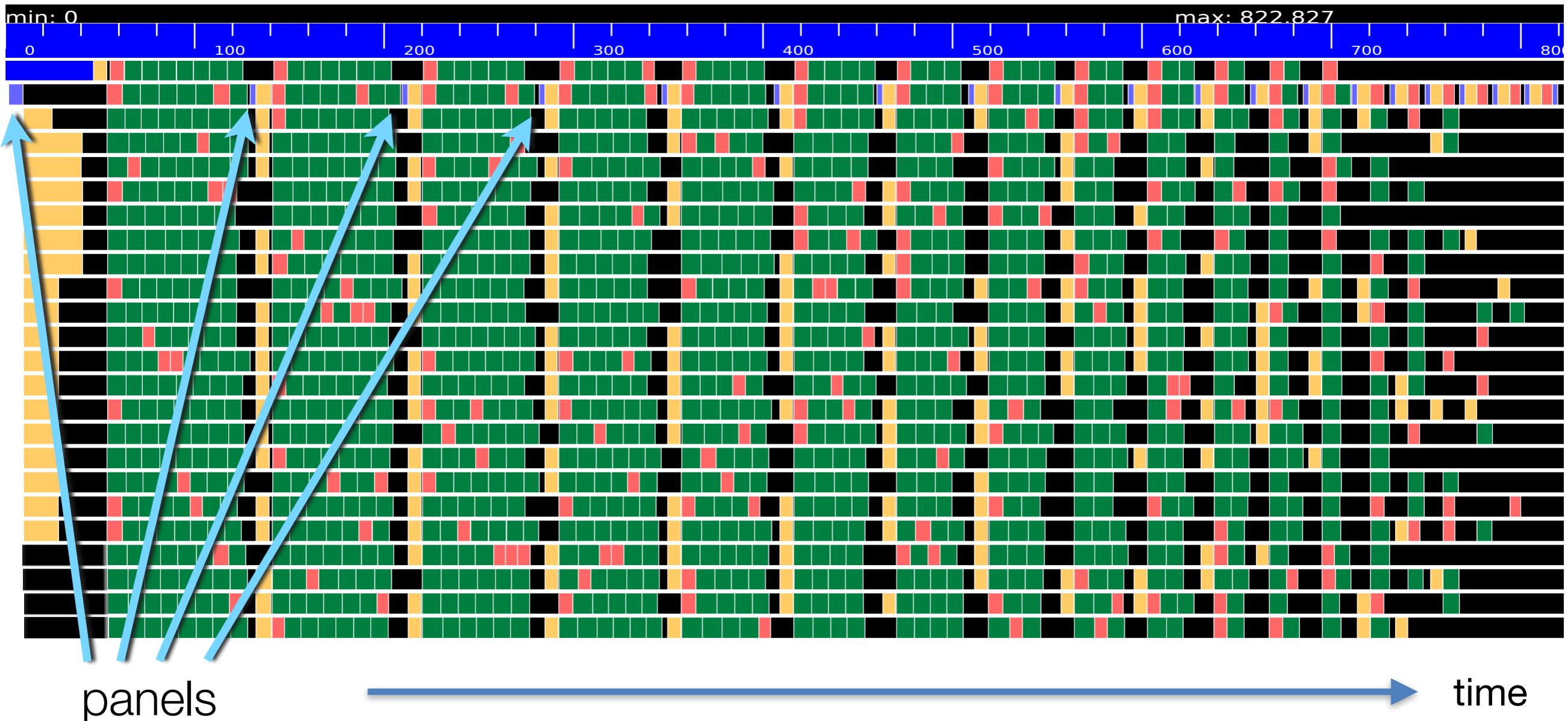
Track dependencies – Directed acyclic graph (DAG)

Critical path with infinite number of processors (need 10 here)
yields minimum possible time



Execution trace

- LAPACK-style fork-join leave cores idle



potrf

trsm

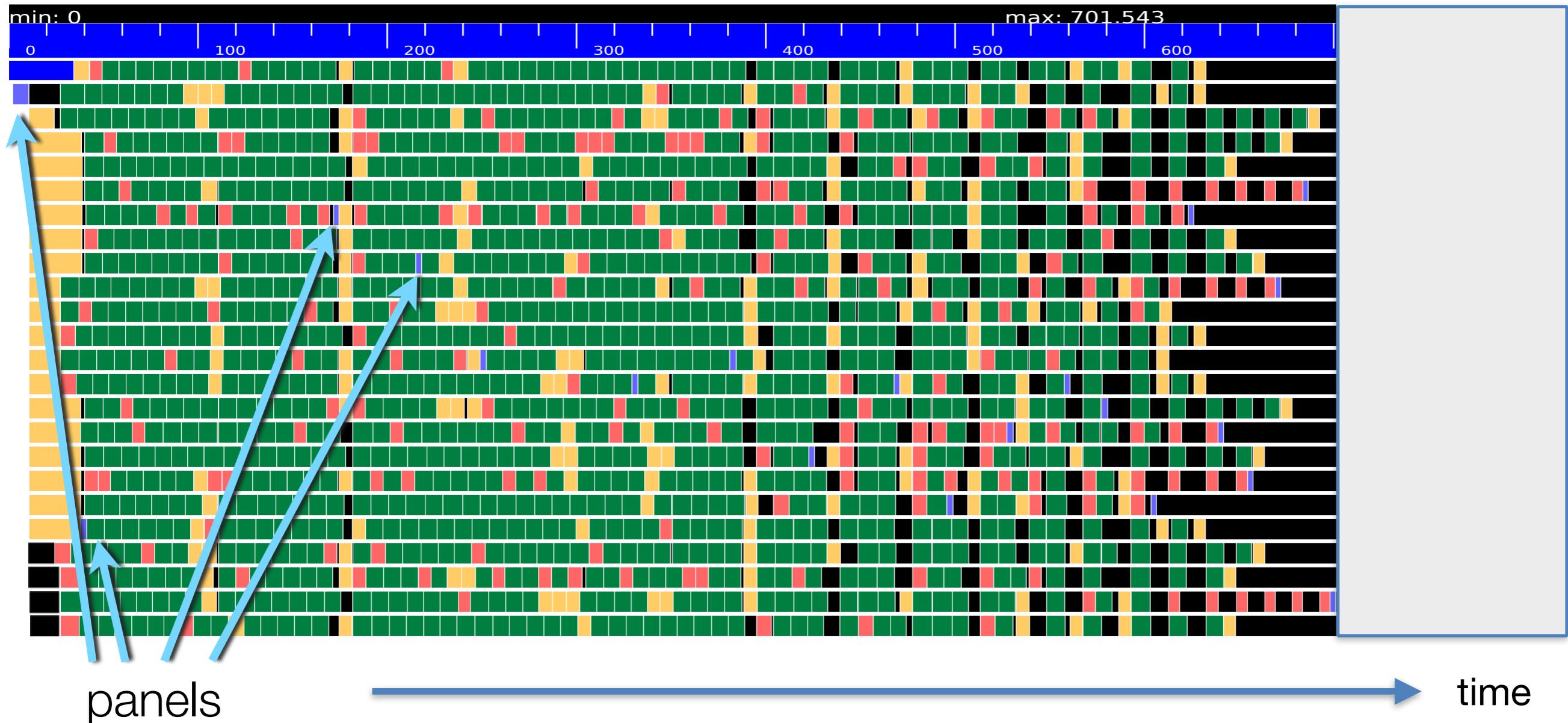
syrk

gemm

idle

Execution trace

- PLASMA squeezes out idle time



24 cores
Matrix is 8000×8000 , tile size is 400×400 .

potrf

trsm

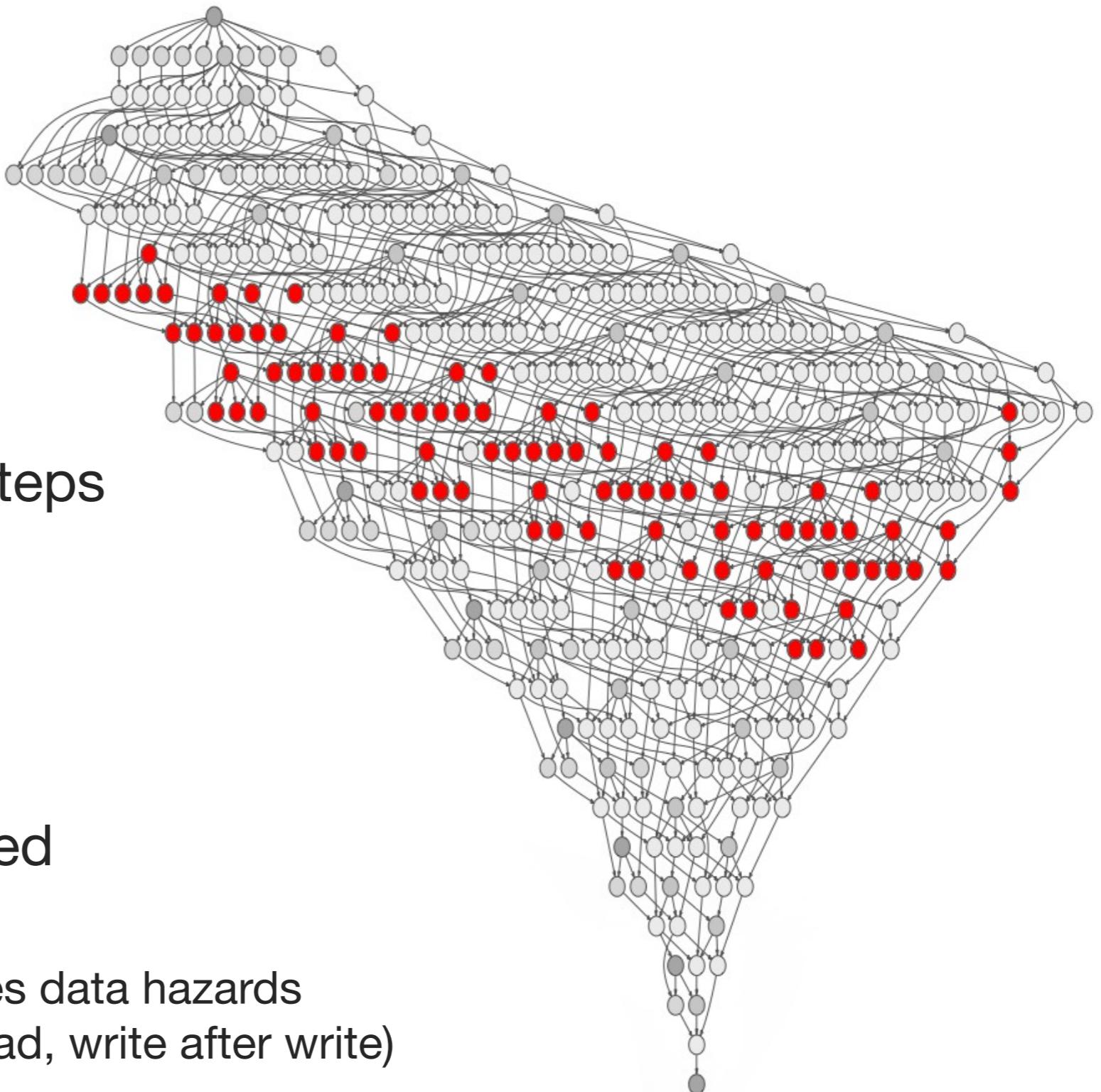
syrk

gemm

idle

Dataflow scheduling

- Exploit parallelism
- Load balance
- Remove artificial synchronization between steps
- Maximize data locality
- Parallel correctness inherited from serial tile algorithm
 - Runtime automatically resolves data hazards (read after write, write after read, write after write)



Dynamic scheduling runtimes

- Jade
- SMPSS / OMPSs
- StarPU
- QUARK
- SuperGlue & DuctTEiP
- OpenMP 4
- Stanford University
- Barcelona Supercomputer Center
- INRIA Bordeaux
- University of Tennessee
- Uppsala University



May 2008

OpenMP 3.0

#pragma omp task

April 2009



GCC 4.4

July 2013

OpenMP 4.0

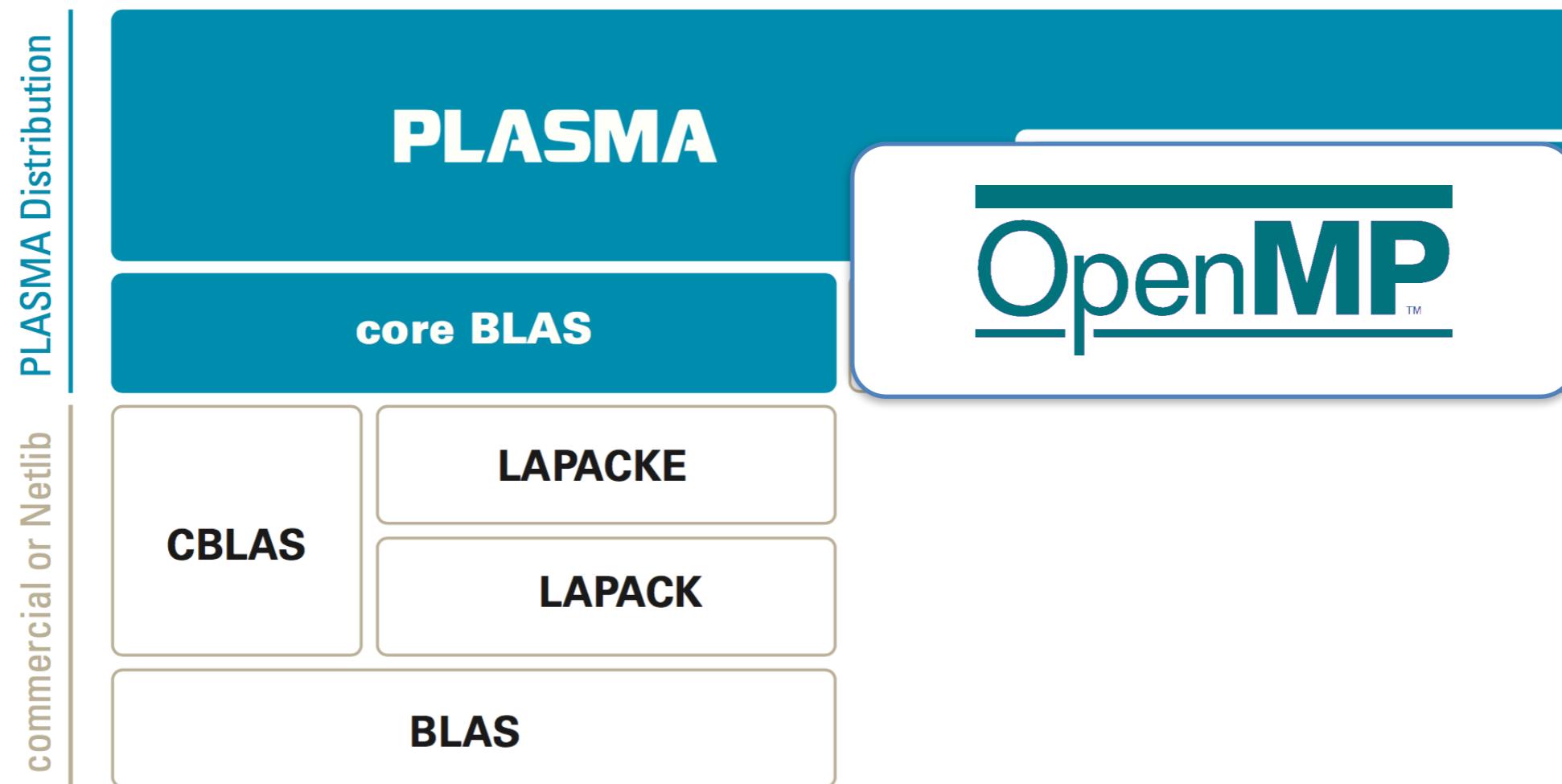
#pragma omp task depend

April 2014



GCC 4.9

PLASMA architecture



Cholesky: pseudo-code

```
// Sequential Tile Cholesky
for k = 1 ... ntiles {
    potrf( Akk )
    for i = k+1 ... ntiles {
        trsm( Akk, Aik )
    }
    for i = k+1 ... ntiles {
        syrk( Aik, Aii )
        for j = i+1 ... ntiles {
            gemm( Ajk, Aik, Aij )
        }
    }
}
```

```
// PLASMA OpenMP Tile Cholesky
for k = 1 .. ntiles {
    omp task potrf( ... )
    for i = k+1 .. ntiles {
        omp task trsm( ... )
    }
    for i = k+1 .. ntiles {
        omp task syrk( ... ) {
            for j = i+1 .. ntiles
                omp task gemm( ... )
        }
    }
}
```

```

#pragma omp parallel
#pragma omp master
{
    for (k = 0; k < nt; k++) {
        #pragma omp task depend(inout:A(k,k)[0:nb*nb])
        info = LAPACKE_dpotrf_work(
            LAPACK_COL_MAJOR,
            lapack_const(PlasmaLower),
            nb, A(k,k), nb);

        for (m = k+1; m < nt; m++) {
            #pragma omp task depend(in:A(k,k)[0:nb*nb]) \
                depend(inout:A(m,k)[0:nb*nb])
            cblas_dtrsm(
                CblasColMajor,
                CblasRight, CblasLower,
                CblasTrans, CblasNonUnit,
                nb, nb,
                1.0, A(k,k), nb,
                A(m,k), nb);
        }
        for (m = k+1; m < nt; m++) {
            #pragma omp task depend(in:A(m,k)[0:nb*nb]) \
                depend(inout:A(m,m)[0:nb*nb])
            cblas_dsyrk(
                CblasColMajor,
                CblasLower, CblasNoTrans,
                nb, nb,
                -1.0, A(m,k), nb,
                1.0, A(m,m), nb);

            for (n = k+1; n < m; n++) {
                #pragma omp task depend(in:A(m,k)[0:nb*nb]) \
                    depend(in:A(n,k)[0:nb*nb]) \
                    depend(inout:A(m,n)[0:nb*nb])
                cblas_dgemm(
                    CblasColMajor,
                    CblasNoTrans, CblasTrans,
                    nb, nb, nb,
                    -1.0, A(m,k), nb,
                    A(n,k), nb,
                    1.0, A(m,n), nb);
            }
        }
    }
}

```

Cholesky: OpenMP

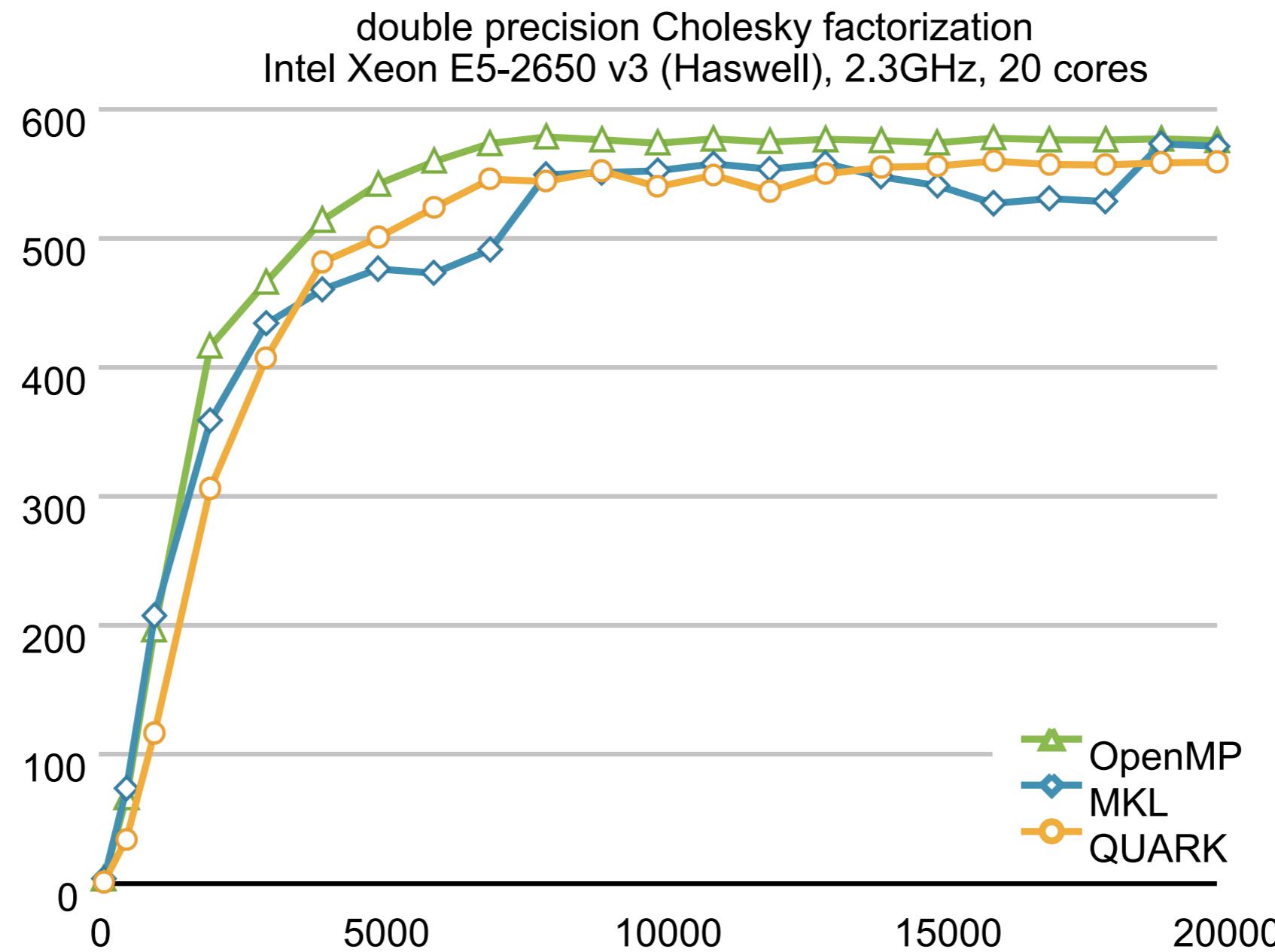
```

#pragma omp task depend(in:A(m,k)[0:nb*nb]) \
    depend(in:A(n,k)[0:nb*nb]) \
    depend(inout:A(m,n)[0:nb*nb])

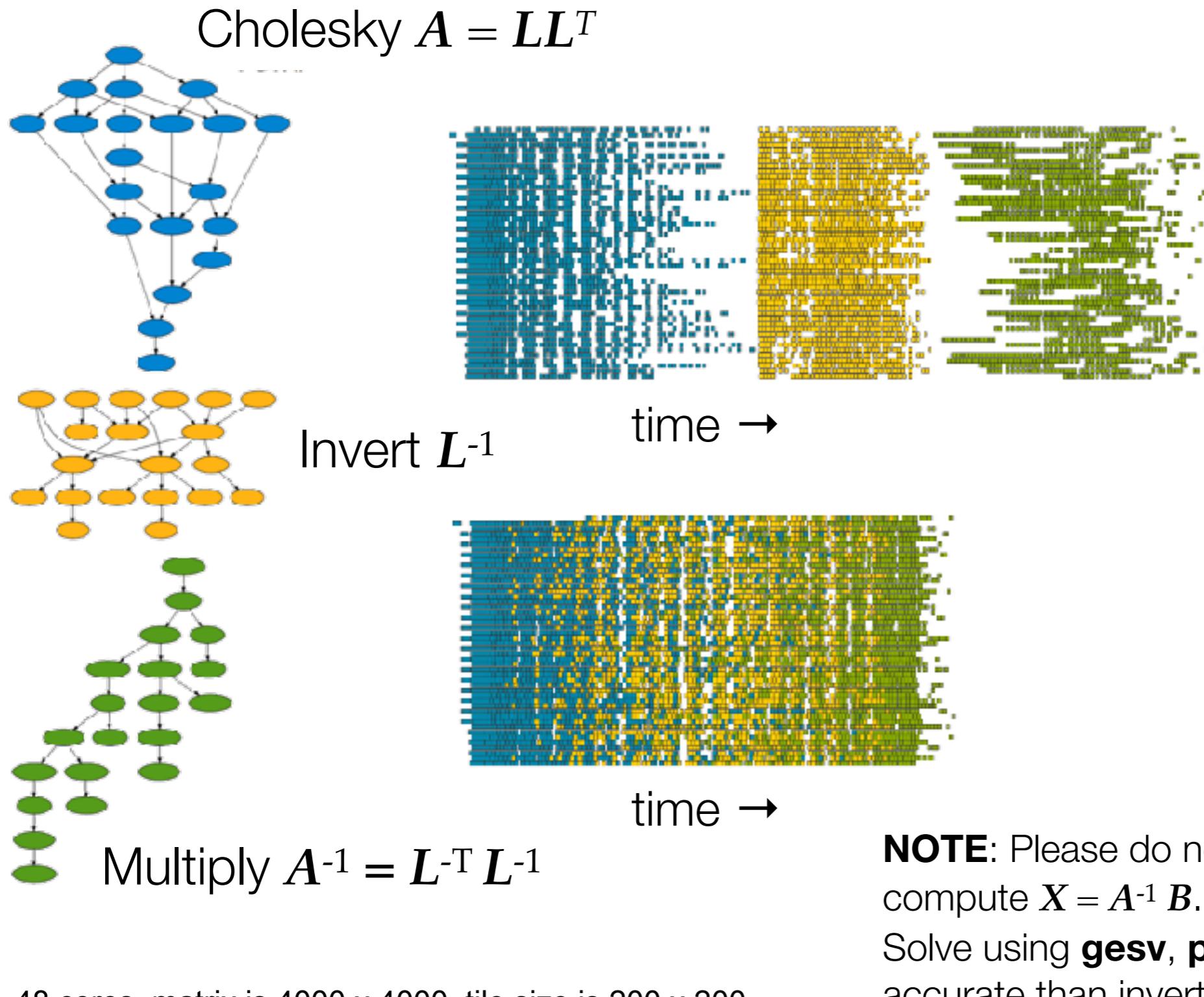
cblas_dgemm(
    CblasColMajor,
    CblasNoTrans, CblasTrans,
    nb, nb, nb,
    -1.0, A(m,k), nb,
    A(n,k), nb,
    1.0, A(m,n), nb);

```

Cholesky performance



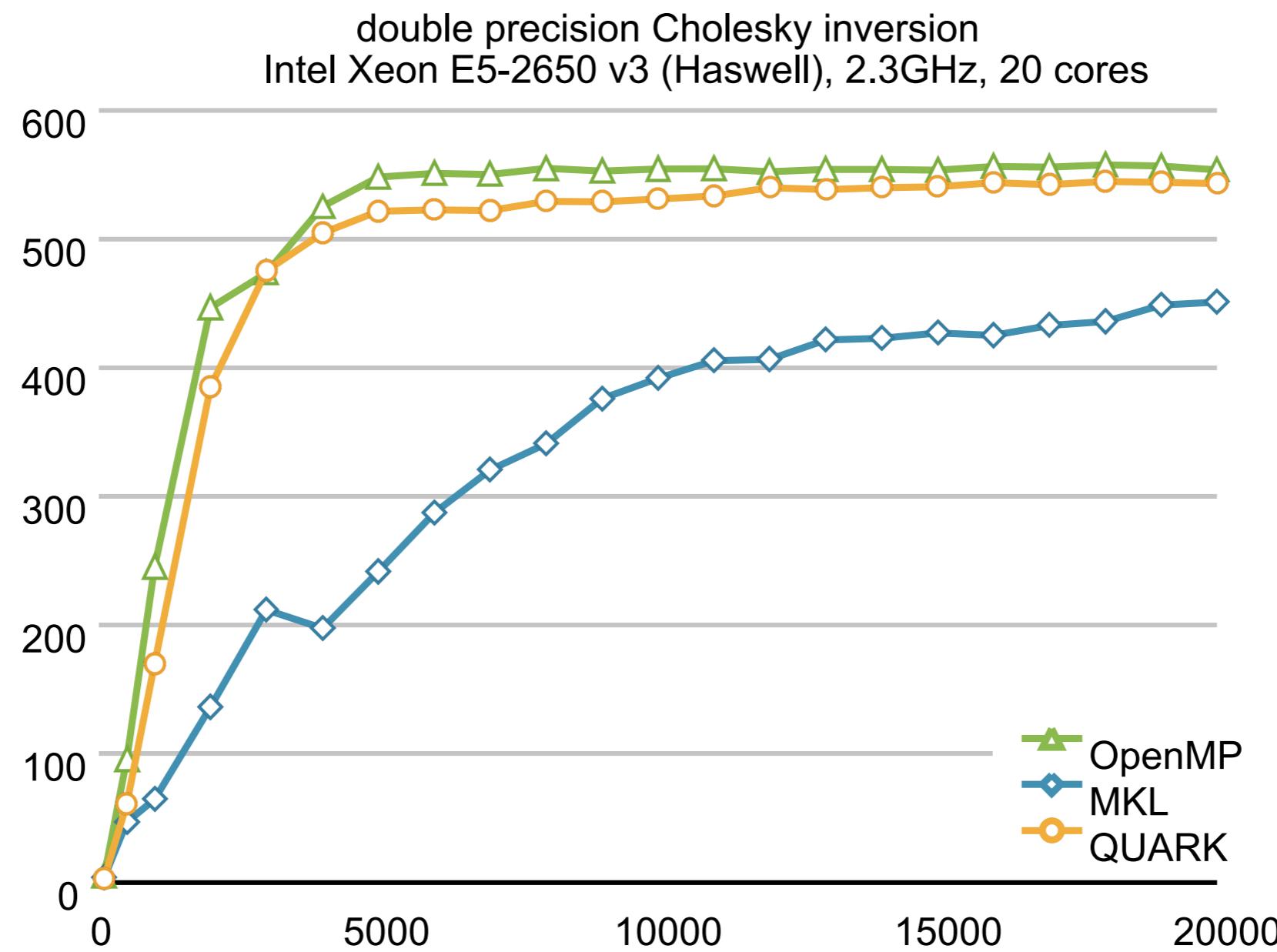
Merging DAGs



48 cores, matrix is 4000 x 4000, tile size is 200 x 200.

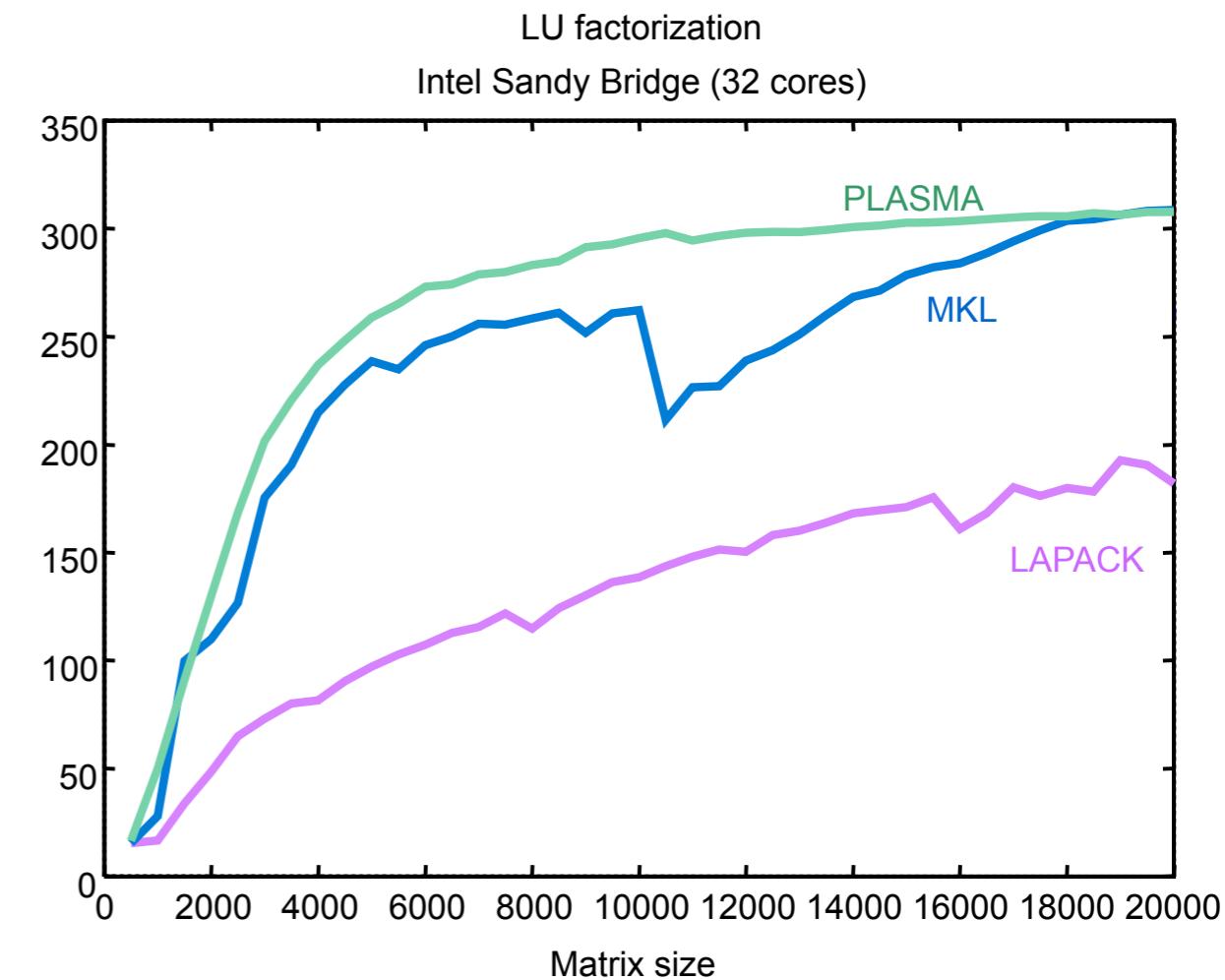
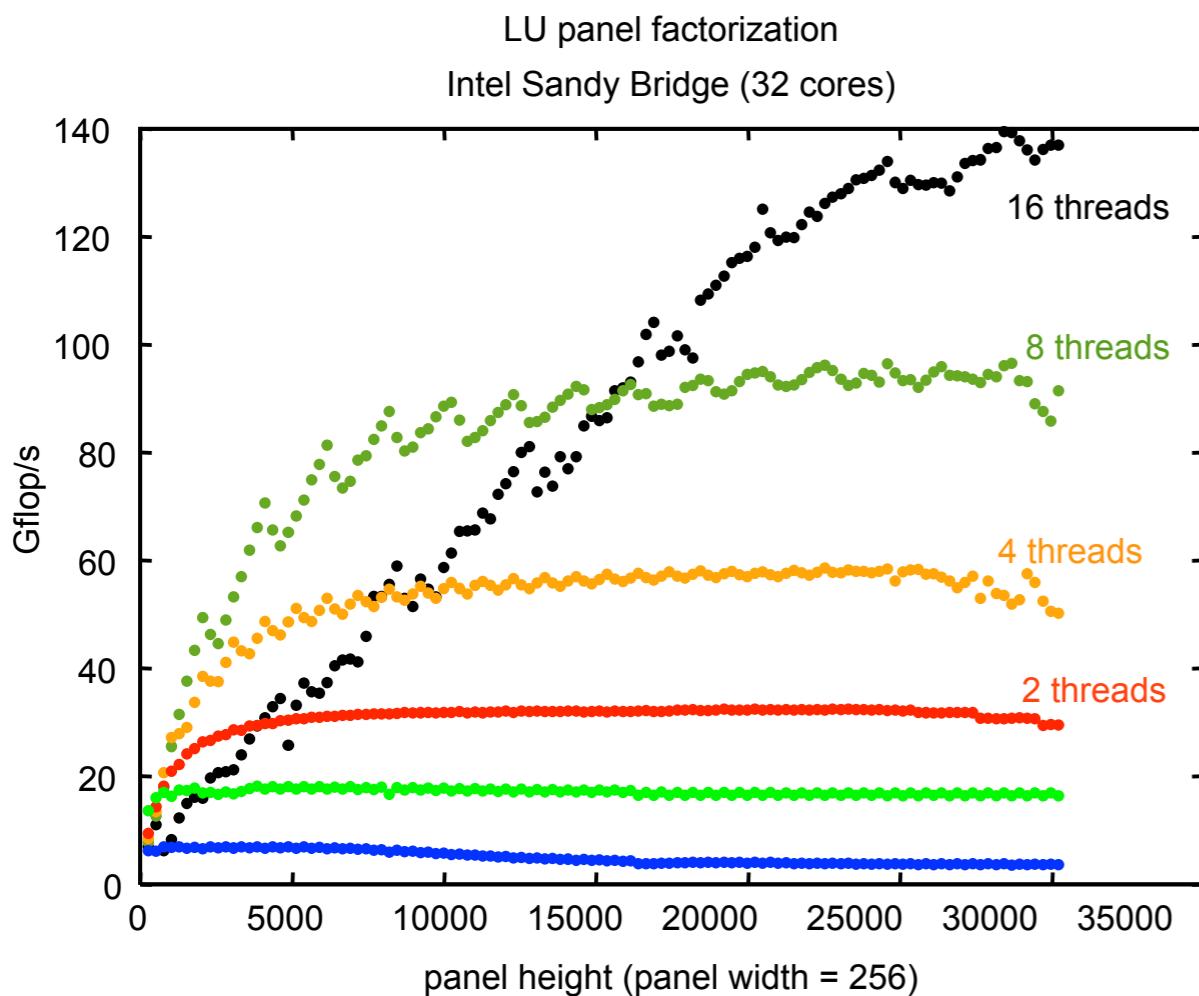
NOTE: Please do not explicitly invert matrices to compute $\mathbf{X} = \mathbf{A}^{-1} \mathbf{B}$.
Solve using **gesv**, **posv**, or **sysv**; faster and more accurate than inverting and multiplying

Cholesky inversion performance



LU performance

- Multi-threaded panel
- Factorization reaches peak performance quickly

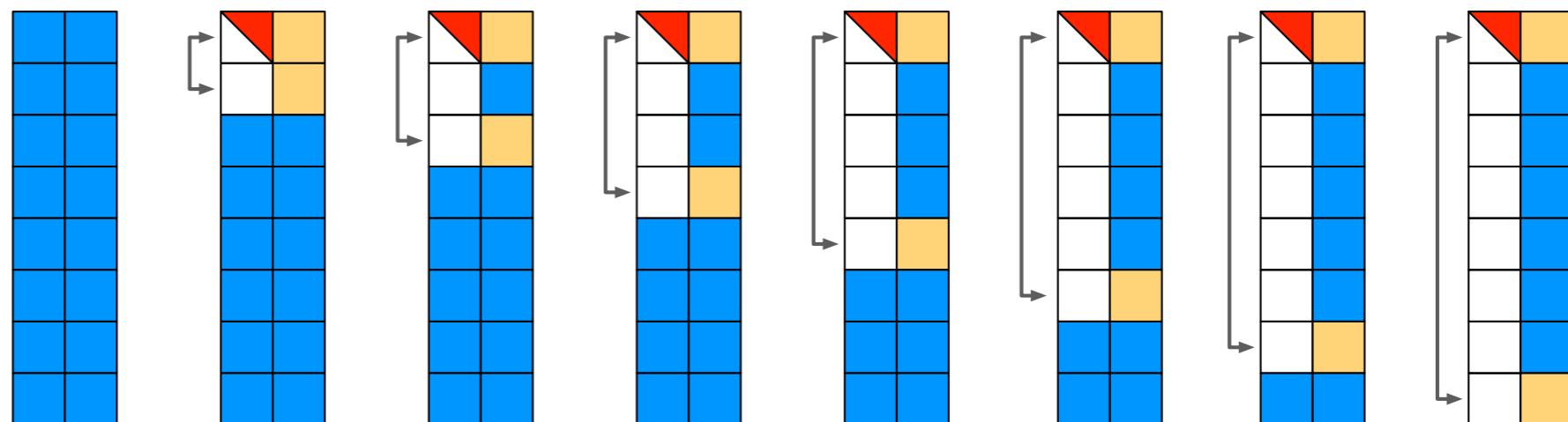


S. Donfack, J. Dongarra, M. Faverge, M. Gates, J. Kurzak, P. Luszczek, I. Yamazaki, *A survey of recent developments in parallel implementations of Gaussian elimination*, Concurrency and Computation: Practice and Experience, 27(5):1292–1309, 2015. [DOI: 10.1002/cpe.3110](https://doi.org/10.1002/cpe.3110)

PLASMA QR

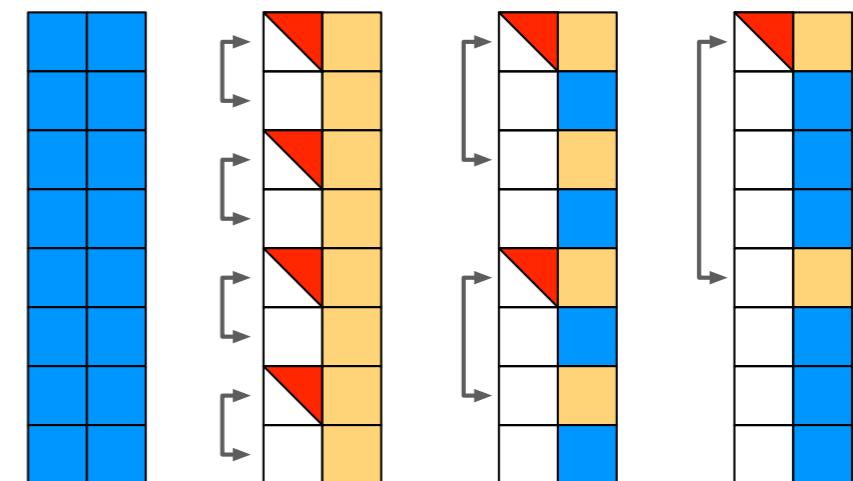
- Tile QR

- Great for square matrices
- Great for multicore
- Pairwise reductions
“domino”



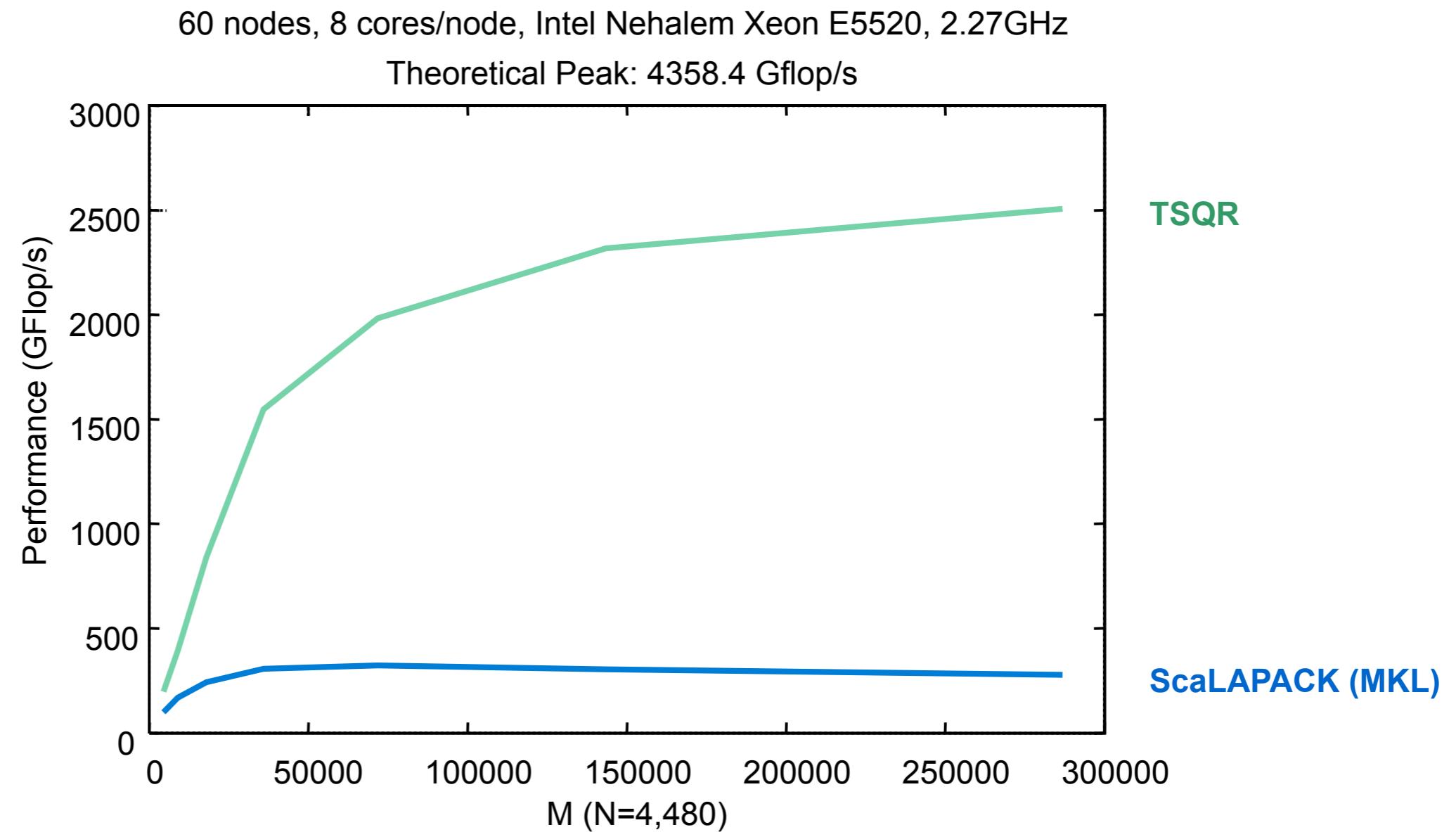
- TSQR / CAQR

- Tall-skinny QR
- Communication avoiding QR
- Tree of pairwise reductions
- Great for tall-skinny matrices (least squares)
- Great for distributed memory
- But triangle-triangle (TT) pairs less efficient than square-square (SS) or triangle-square (TS)



TSQR performance

- Fixed cols, increase rows from square (left) to tall-skinny (right)

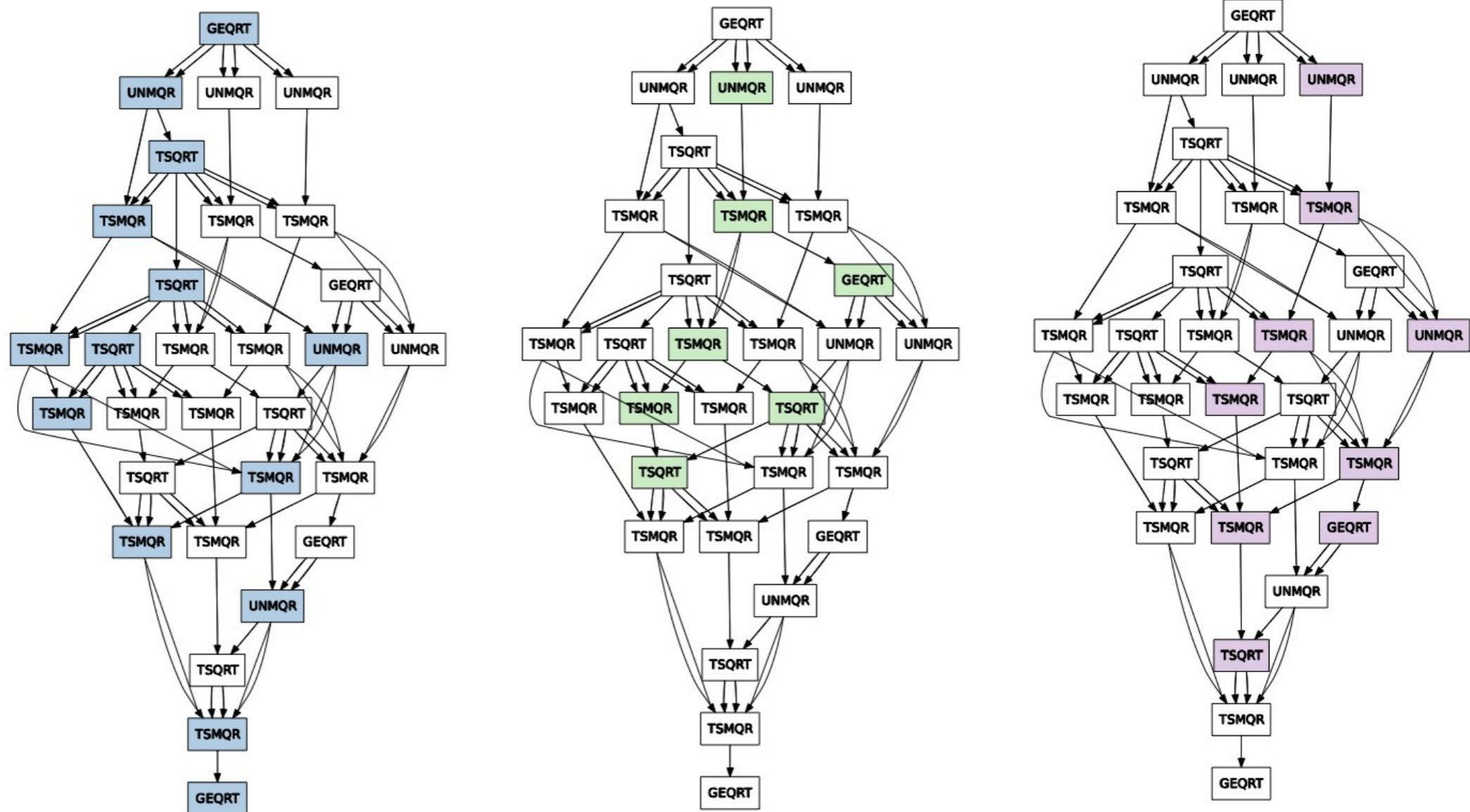


Outline

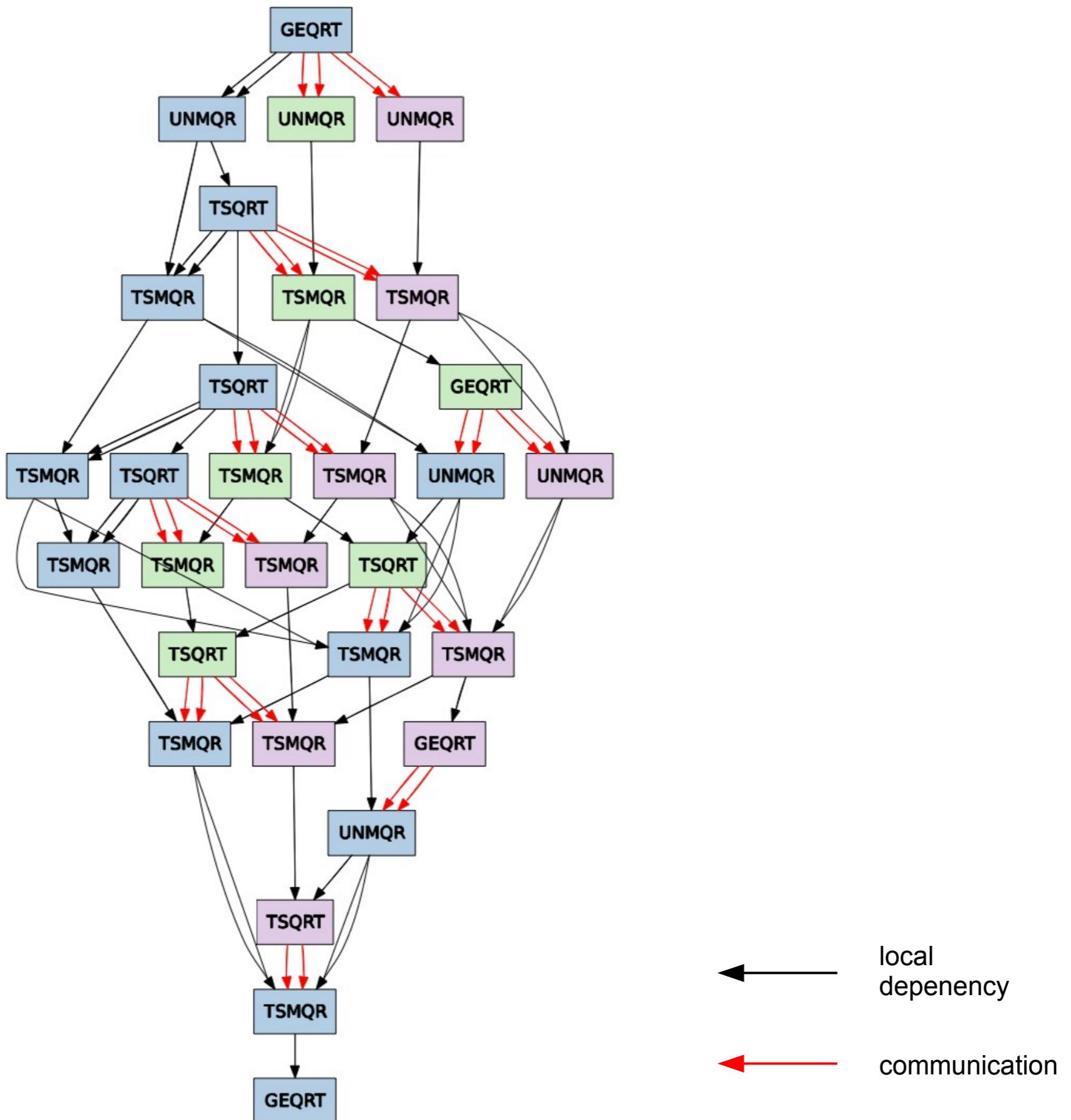
- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - **DPLASMA and PaRSEC**
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

Distributed dataflow execution

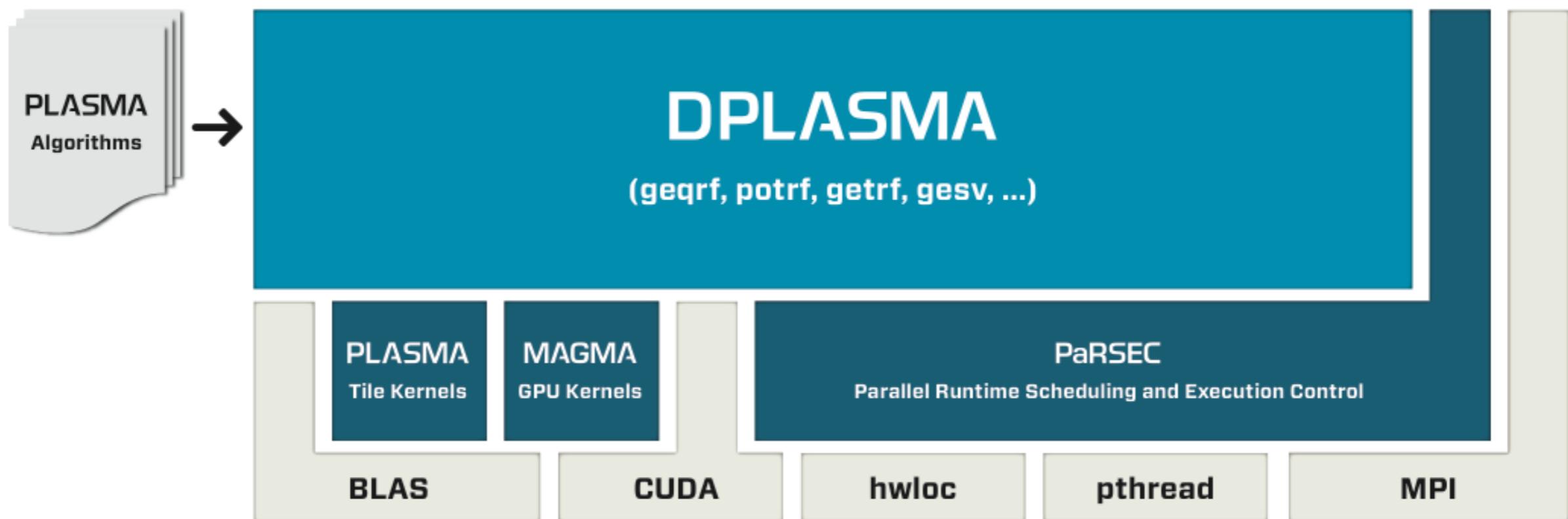
A00	A01	A02	A03
A10	A11	A12	A13
A20	A21	A22	A23
A30	A31	A32	A33



A ₀₀	A ₀₁	A ₀₂	A ₀₃
A ₁₀	A ₁₁	A ₁₂	A ₁₃
A ₂₀	A ₂₁	A ₂₂	A ₂₃
A ₃₀	A ₃₁	A ₃₂	A ₃₃

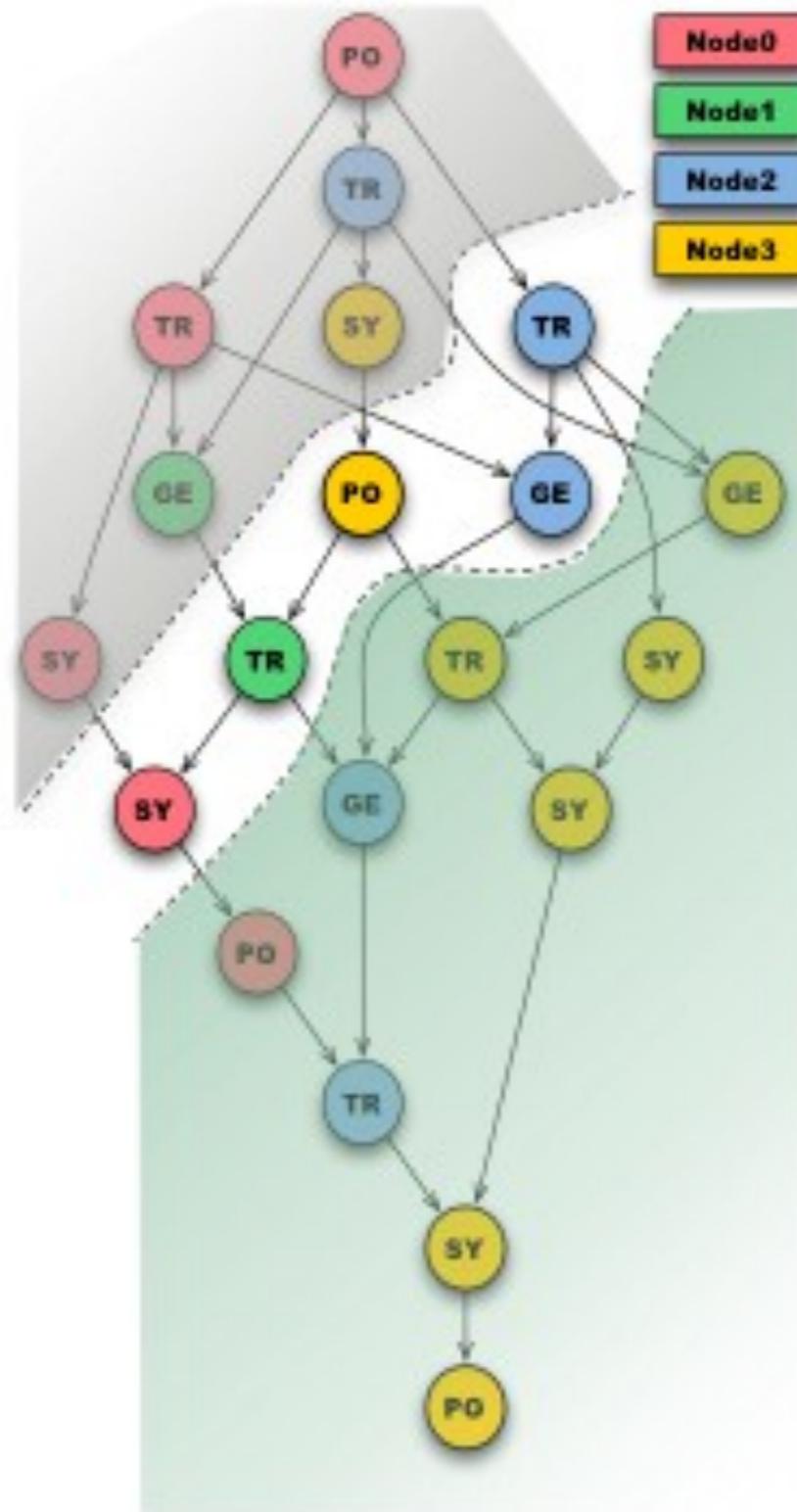


DPLASMA architecture



Parameterized task graph (PTG)

- Symbolic DAG representation
- Problem size independent
- Completely distributed
- DAG is $O(n^3)$,
- PTG avoids explicitly creating it
- Runtime
 - Data-driven execution
 - Locality-aware scheduling
 - Communication overlap



From serial code to PTG

Serial code

```
FOR k=0 TO N-1
    DGEQRT(inoutAkk)
    FOR n=k+1 to N
        DORMQR(inA $\blacksquare_{kk}$ , inoutAkn)
    FOR m=k+1 to N
        DTSMQR(inA $\blacksquare_{kk}$ , inoutAmk)
        FOR n=k+1 to N
            DTSMQR(inAmk, inoutAkn, inoutAmn)
```

PTG representation

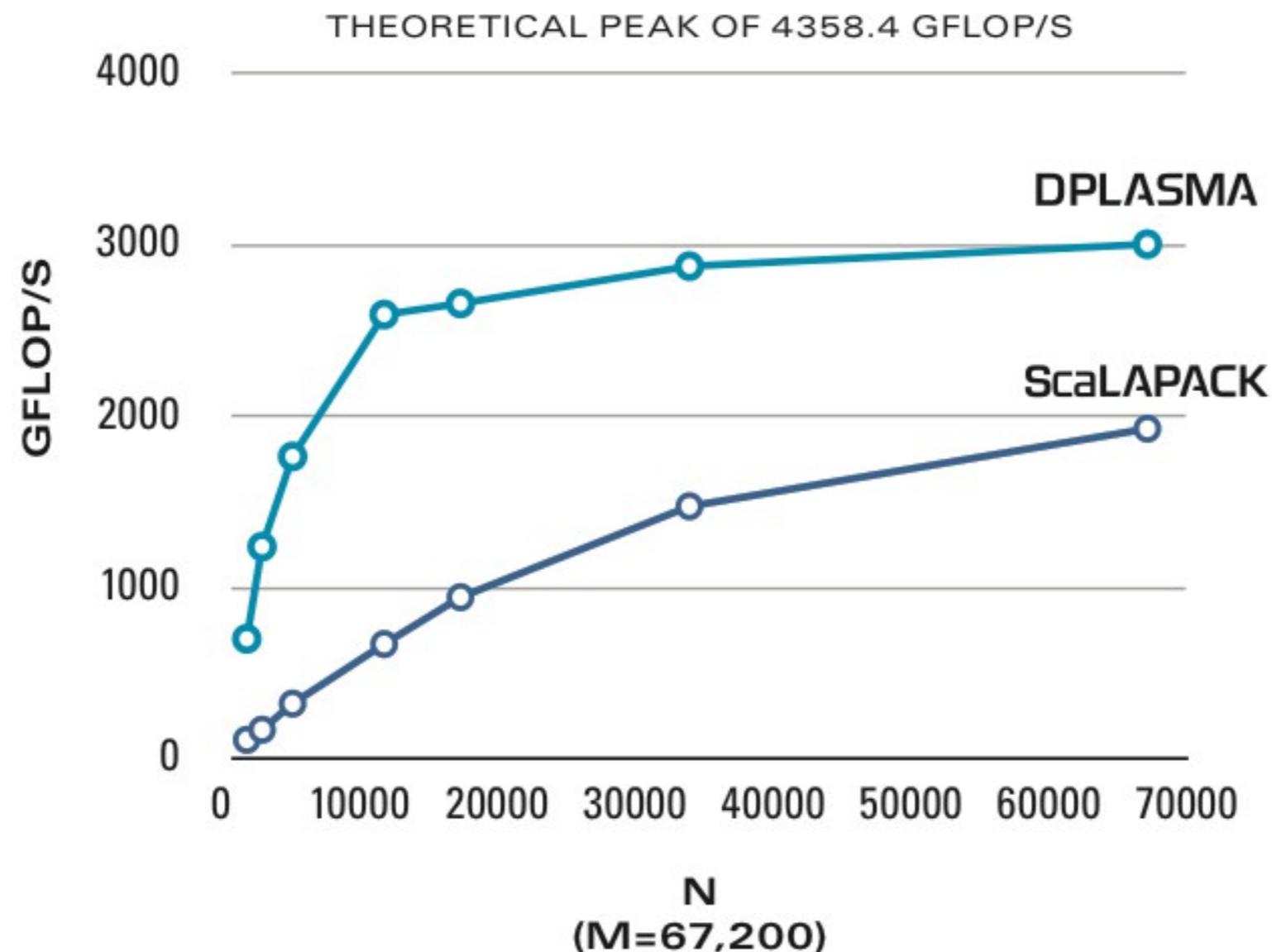
DGEQRT _{kkk}	
1 _{ARG}	$\leftarrow A_{k,k} \mid DTSMQR_{k,k,k-1}$
1 _{ARG}	$\Rightarrow DORMQR_{k,k+1..N,k}(\blacksquare)$
1 _{ARG}	$\Rightarrow DTSMQR_{k+1,k,k}(\blacksquare)$
1 _{ARG}	$\Rightarrow A_{k,k}(\blacksquare)$
DORMQR _{knk}	
1 _{ARG}	$\leftarrow DGEQRT_{k,k,k}(\blacksquare)$
2 _{ARG}	$\leftarrow A_{k,n} \mid DTSMQR_{k,n,k-1}$
2 _{ARG}	$\Rightarrow DTSMQR_{k+1,n,k}$
2 _{ARG}	$\Rightarrow A_{k,n}$
DTSMQR _{mkk}	
1 _{ARG}	$\leftarrow DGEQRT_{m-1,k,k}(\blacksquare) \mid DTSMQR_{m-1,k,k}(\blacksquare)$
1 _{ARG}	$\Rightarrow DTSMQR_{m+1,k,k}(\blacksquare) \mid A_{k,k}(\blacksquare)$
2 _{ARG}	$\leftarrow A_{m,k} \mid DTSMQR_{m,k,k-1}$
2 _{ARG}	$\Rightarrow DTSMQR_{m,k+1..N,k}$
2 _{ARG}	$\Rightarrow A_{m,k}$
DTSMQR _{mnk}	
1 _{ARG}	$\leftarrow DTSMQR_{m,k,k}$
2 _{ARG}	$\leftarrow DORMQR_{m-1,n,k} \mid DTSMQR_{m-1,n,k}$
2 _{ARG}	$\Rightarrow DTSMQR_{m+1,n,k} \mid A_{n,k}$
3 _{ARG}	$\leftarrow A_{m,n} \mid DTSMQR_{m,n,k-1}$
3 _{ARG}	$\Rightarrow DGEQRT_{m,n,k+1} \mid DORMQR_{m,n,k+1} \mid$
3 _{ARG}	$\Rightarrow DTSMQR_{m,n,k+1} \mid DTSMQR_{m,n,k+1} \mid$
3 _{ARG}	$\Rightarrow A_{m,n}$

QR performance

- Fixed rows, increase cols, from tall-skinny (left) to square (right)

Solving Linear Least Square Problem (DGEQRF)

60-node, 480-core, 2.27GHz Intel Xeon Nehalem, IB 20G System

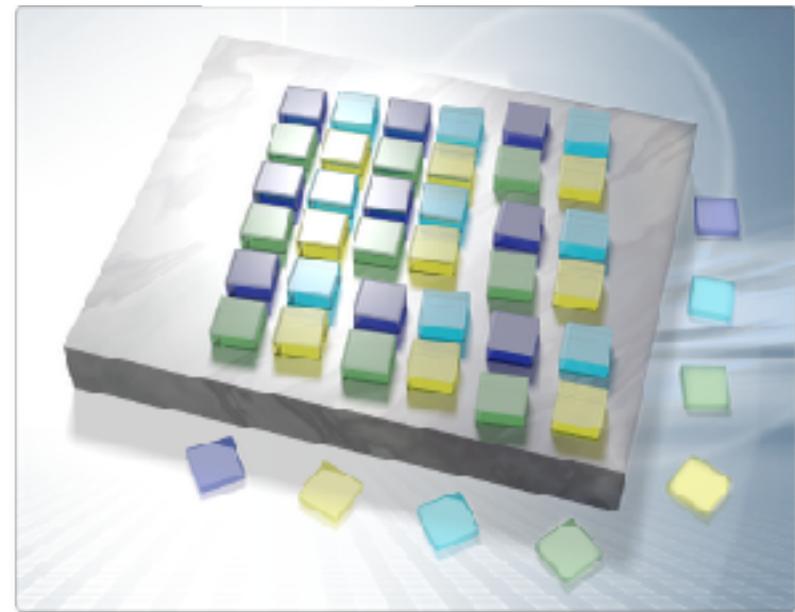


Outline

- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - **SLATE**
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - Case study: 2-stage SVD

SLATE

- C++11, MPI, OpenMP, cuBLAS
- Map of tiles
 - `A.tiles_[{ i, j, device }]`
where `{ i, j, device }` is `std::tuple< int, int, int >`,
`i` and `j` are always global tile indices
 - Supports both contiguous (PLASMA) and non-contiguous (ScaLAPACK) tiles
 - Tiles can be allocated independently, unlike PLASMA
- Trailing matrix updates
 - CPU: OpenMP tasks or Batched BLAS
 - GPU: Batched BLAS
 - Batched BLAS is many small, independent gemms in one call, in parallel
- Explicit look-ahead
 - Factor next panel(s) while trailing matrix is updating
 - Hides communication



SLATE components

- BLAS++
- LAPACK++
 - C++ wrappers around BLAS and LAPACK routines
 - Enables C++ template code
 - Removes 1st letter for precision – determined by arguments
 - $\text{gemm}(A, B, C) \Rightarrow \text{sgemm}, \text{dgemm}, \text{cgemm}, \text{or zgemm}$ depending on type of A, B, and C.
 - Names and arguments consistent across all 4 precisions
 - herk of real matrix $\Rightarrow \text{syrk}$
 - geev always takes complex vector of eigenvalues (w), instead of split-complex as in LAPACK's dgeev (wr, wi)
 - Allocates optimal workspace internally
 - Throws C++ exceptions on hard errors (invalid dimensions, ...); returns numerical errors (singular matrix, ...)

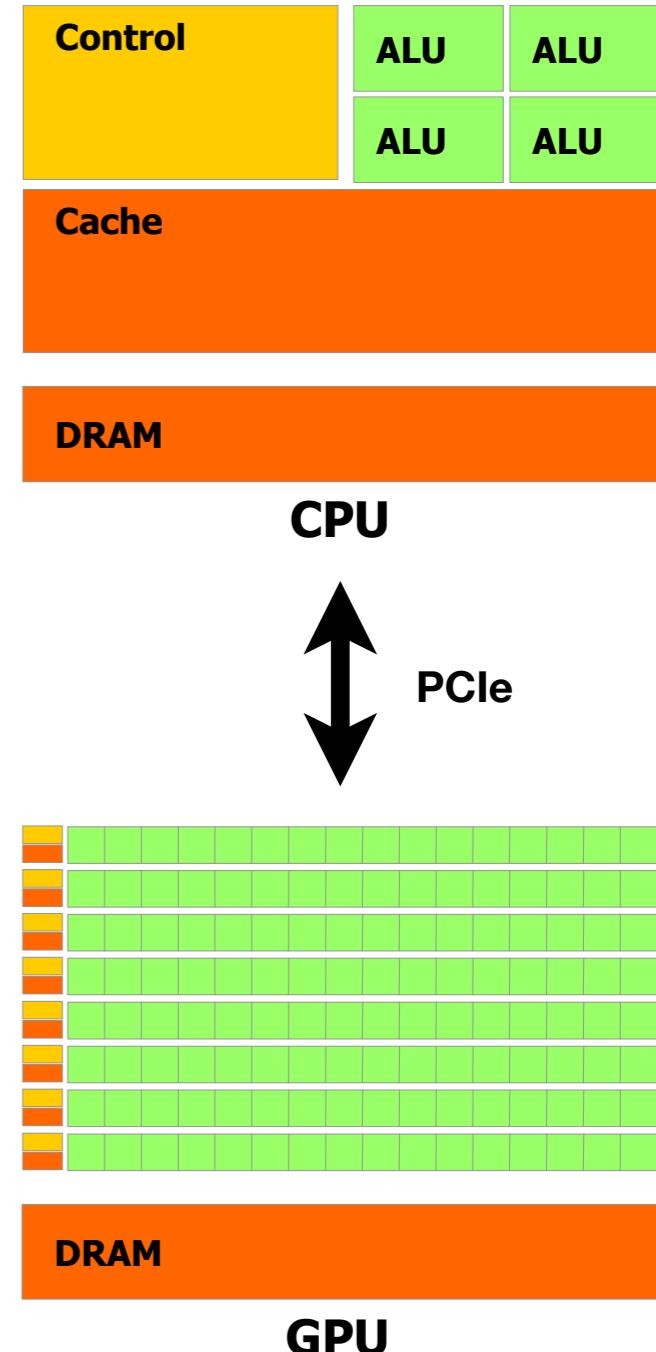


Outline

- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - **MAGMA for CUDA, OpenCL, or Xeon Phi**
 - Case study: 2-stage SVD

Challenges of GPUs

- High levels of parallelism
- Expensive branches (if-then-else)
- Computation vs. communication gap growing
 - NVIDIA Pascal P100 has 4670 Gflop/s, 732 GB/s memory, 16 GB/s PCIe (80 GB/s NVLINK)
- Use hybrid approach
 - Small, non-parallelizable tasks on CPU
 - Large, parallel tasks on GPU
 - Overlap communication & computation

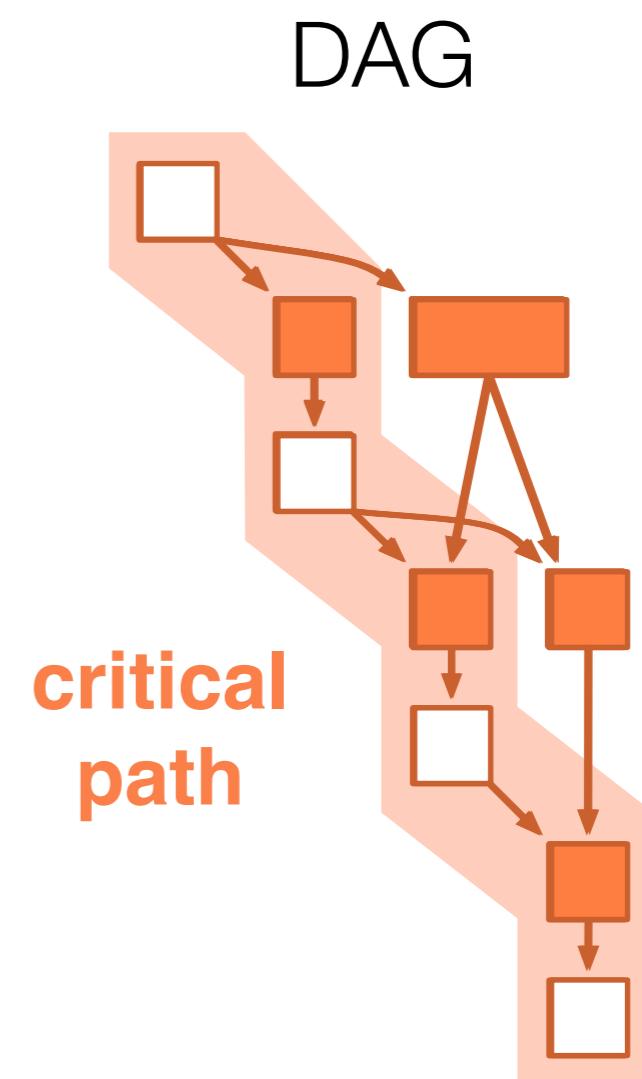
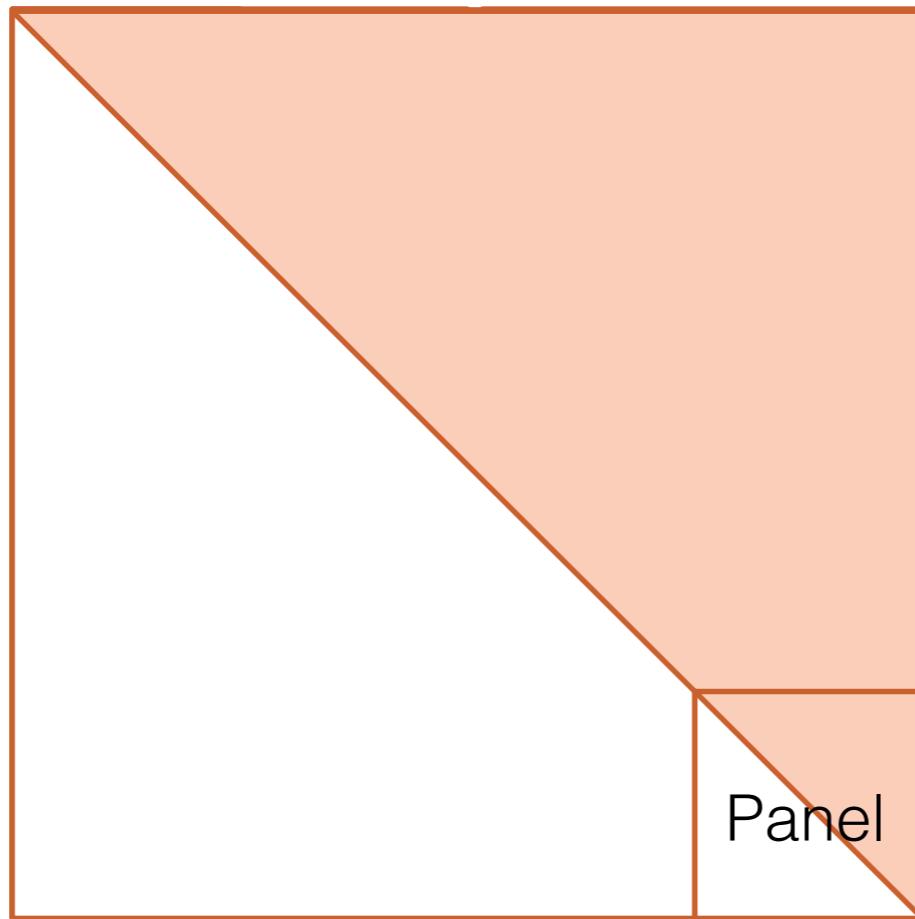


One-sided factorization

- LU, Cholesky, QR factorizations for solving linear systems

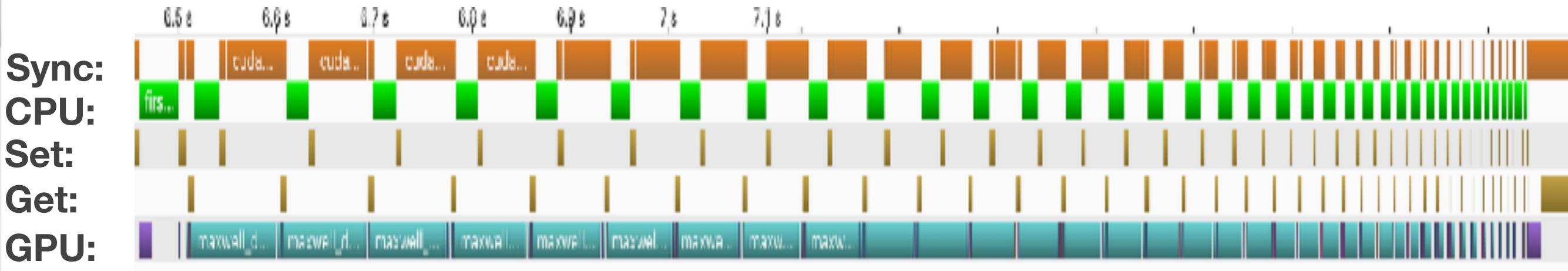
Level 2
BLAS on
CPU

Level 3
BLAS on
GPU



Execution trace

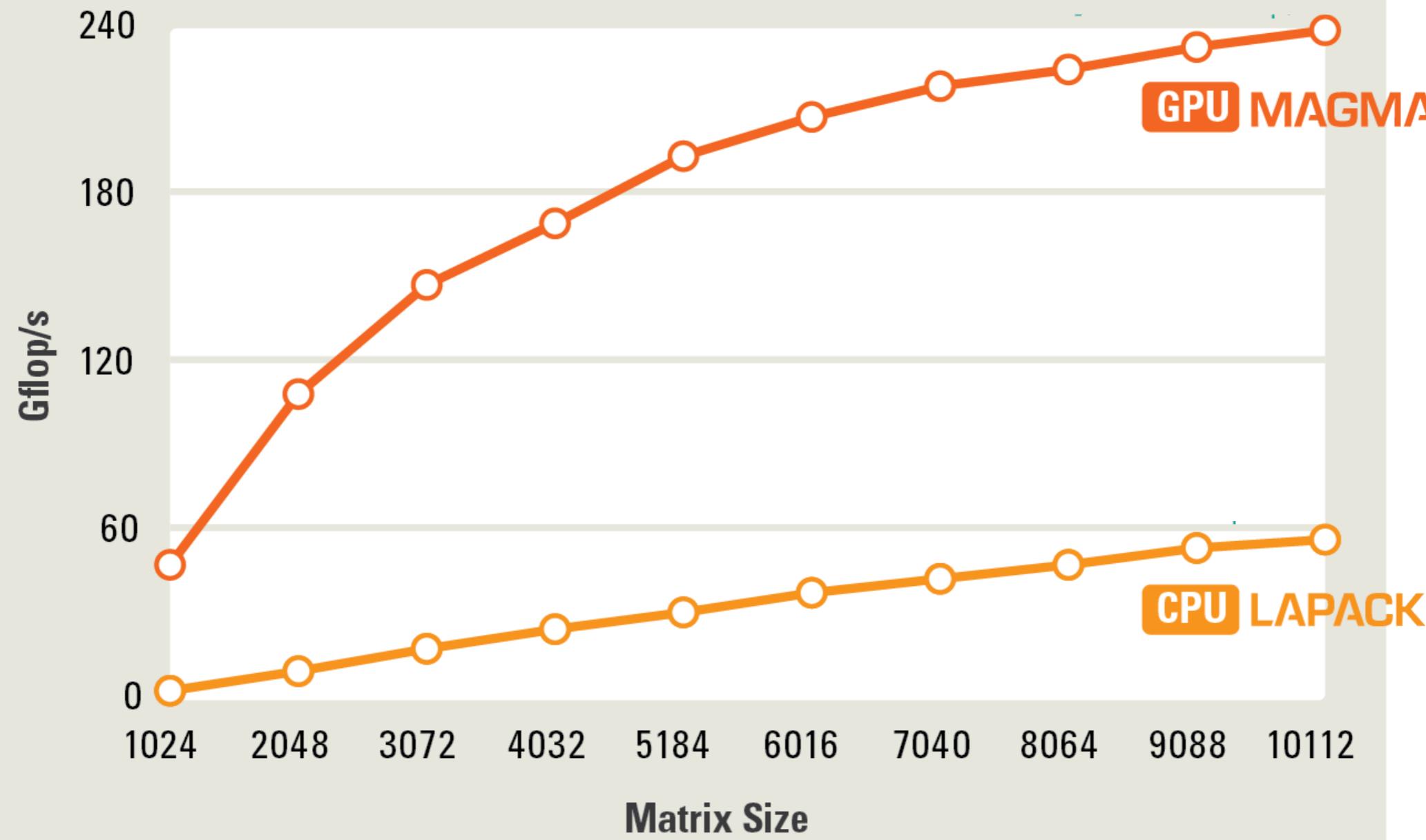
- **Panels** on CPU (**green**) and **set/get communication** (**brown**) overlapped with **trailing matrix updates** (**teal**) on GPU
- Goal to keep GPU busy all the time; CPU may idle



LU factorization (dgetrf), $n = 20000$
P100 GPU, $2 \times 10\text{-core } 2.3 \text{ GHz Haswell}$

- Optimization: for LU, we transpose matrix on GPU so row-swaps are fast

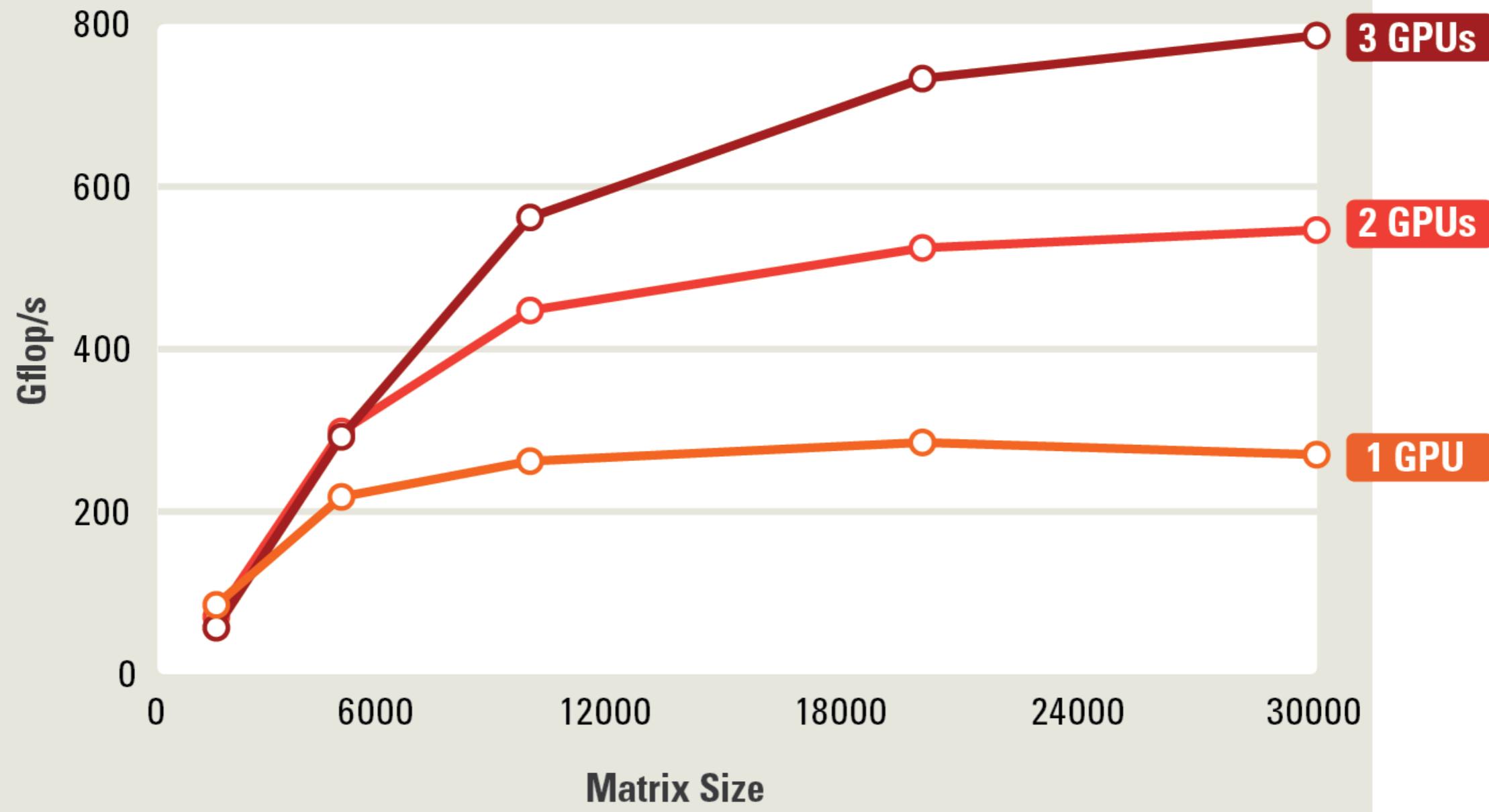
MAGMA LU in double precision on single GPU (C2050)



GPU Fermi C2050 (448 CUDA Cores @ 1.15 GHz)
+ Intel Q9300 (4 cores @ 2.50 GHz)
DP peak **515 + 40** GFlop/s
Power * ~**220** W

CPU AMD Istanbul
[8 sockets x 6 cores (48 cores) @ 2.8GHz]
DP peak **538** GFlop/s
Power * ~**1,022** W

MAGMA LU in double precision on multi-GPUs (Fermi C2070)



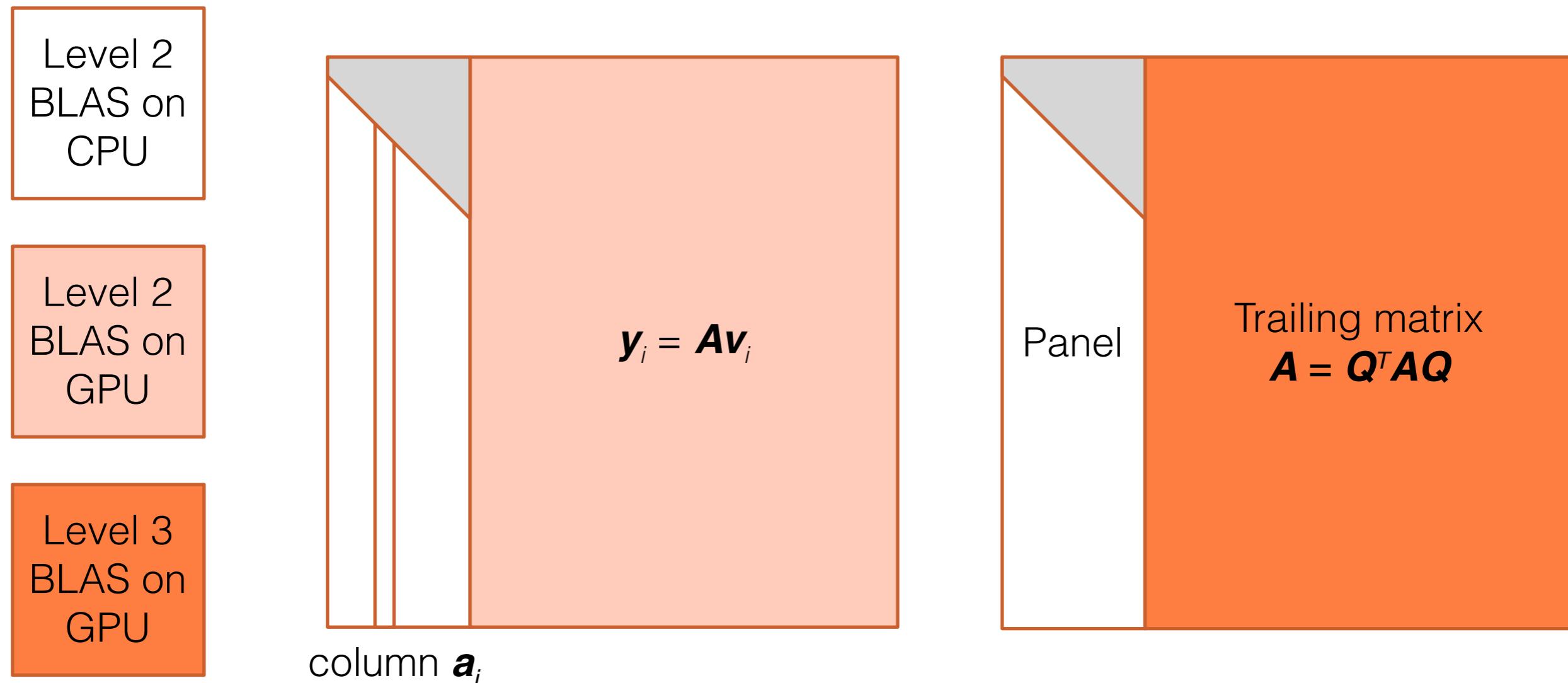
Keeneland system, using one node

3 Nvidia GPUs (M2070 @ 1.1 GHz, 5.4 GB)

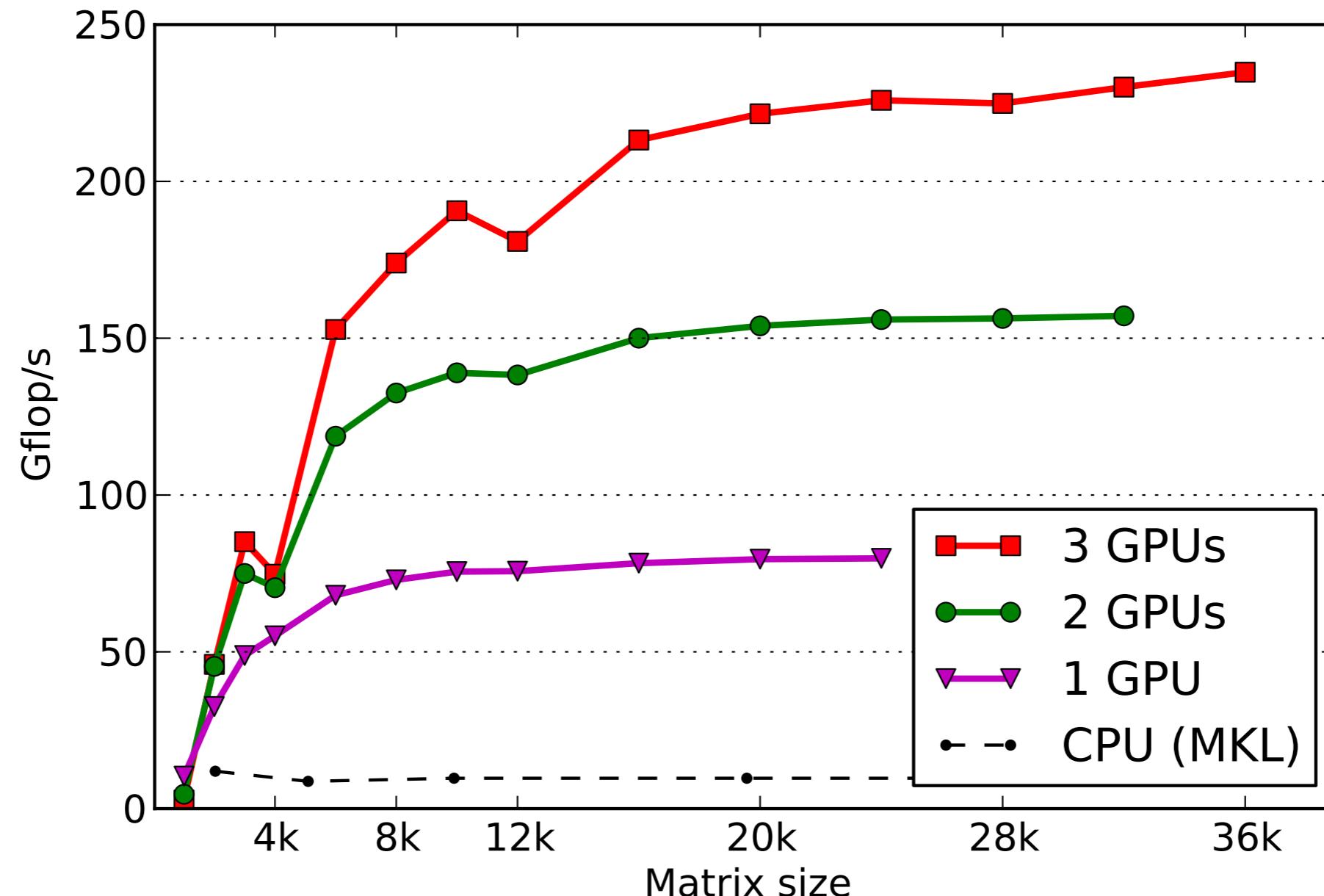
2 x 6 Intel Cores (X5660 @ 2.8 GHz, 23 GB)

Two-sided factorization

- Hessenberg, tridiagonal factorizations for eigenvalue problems



MAGMA Hessenberg in double precision



Keeneland system, using one node

3 Nvidia GPUs (M2070 @ 1.1 GHz, 5.4 GB)

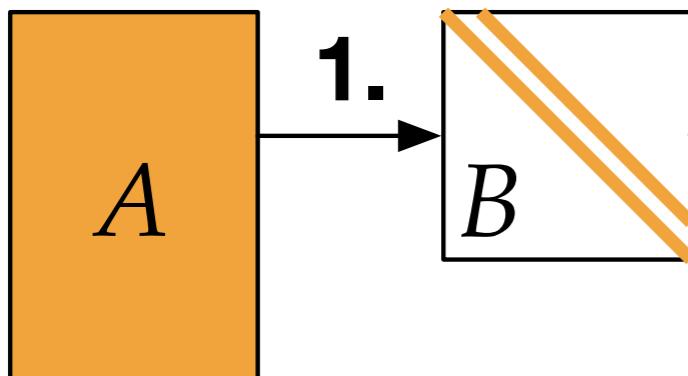
2 x 6 Intel Cores (X5660 @ 2.8 GHz, 23 GB)

Outline

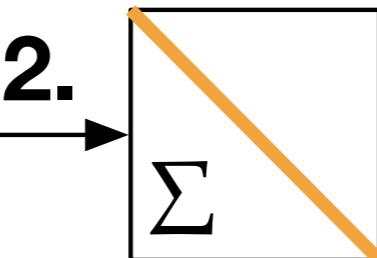
- Legacy Software
 - BLAS
 - LINPACK
 - LAPACK
 - ScaLAPACK
- New Software
 - PLASMA using OpenMP
 - DPLASMA and PaRSEC
 - SLATE
 - SLATE
 - MAGMA for CUDA, OpenCL, or Xeon Phi
 - **Case study: 2-stage SVD**

Computing SVD in 3 Phases

$$A = U_1 B V_1^T$$



$$B = U_2 \Sigma V_2^T$$



3.

$$U = U_1$$

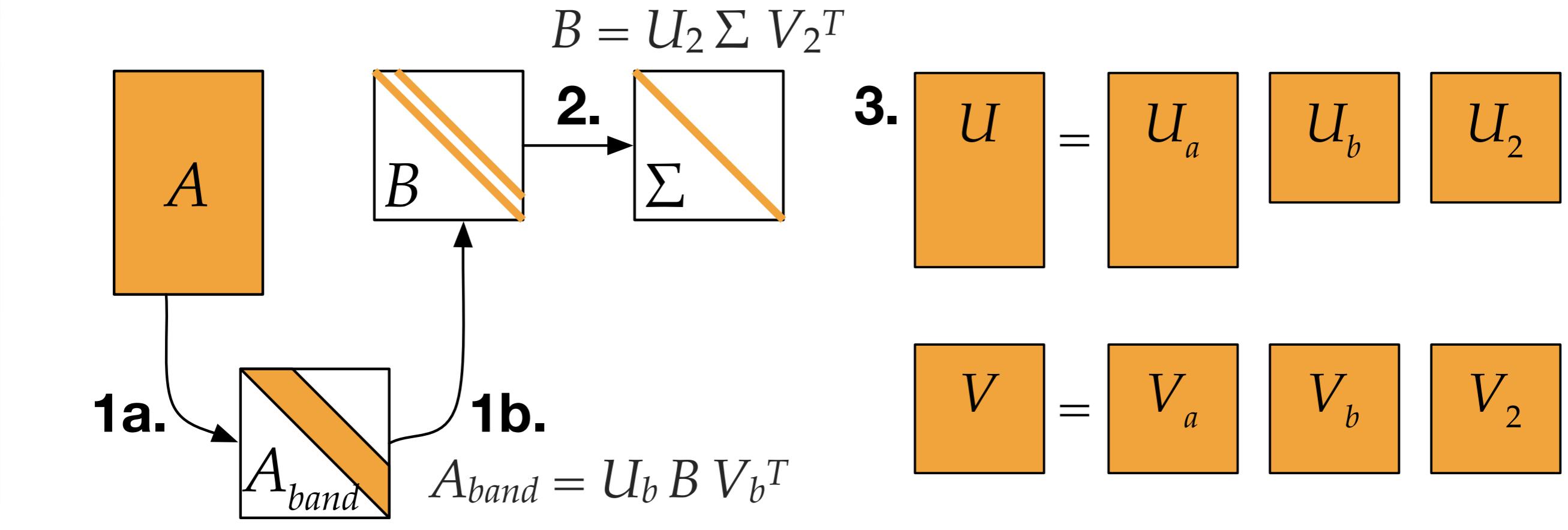
$$U_2$$

$$V = V_1$$

$$V_2$$

1. Reduction to bidiagonal form
2. Bidiagonal SVD
3. Back transform singular vectors

Computing SVD in 3 Phases



1. Reduction to bidiagonal form (2 stage)

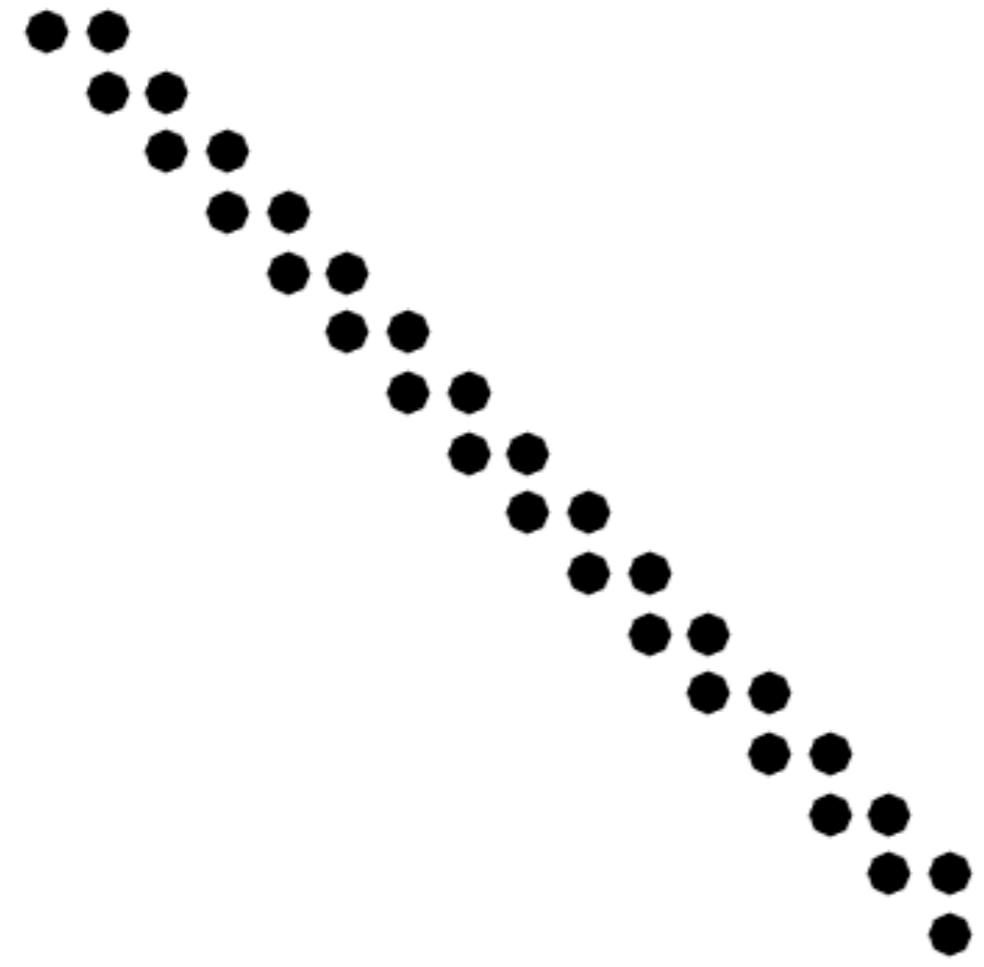
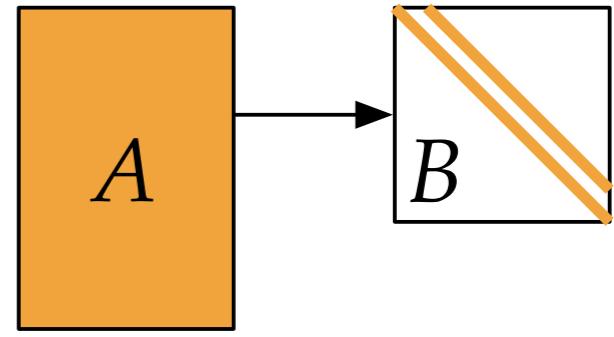
1a. Full to band

1b. Band to bidiagonal

2. Bidiagonal SVD

3. Back transform singular vectors

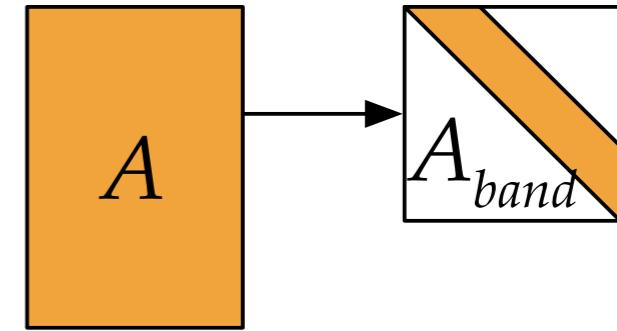
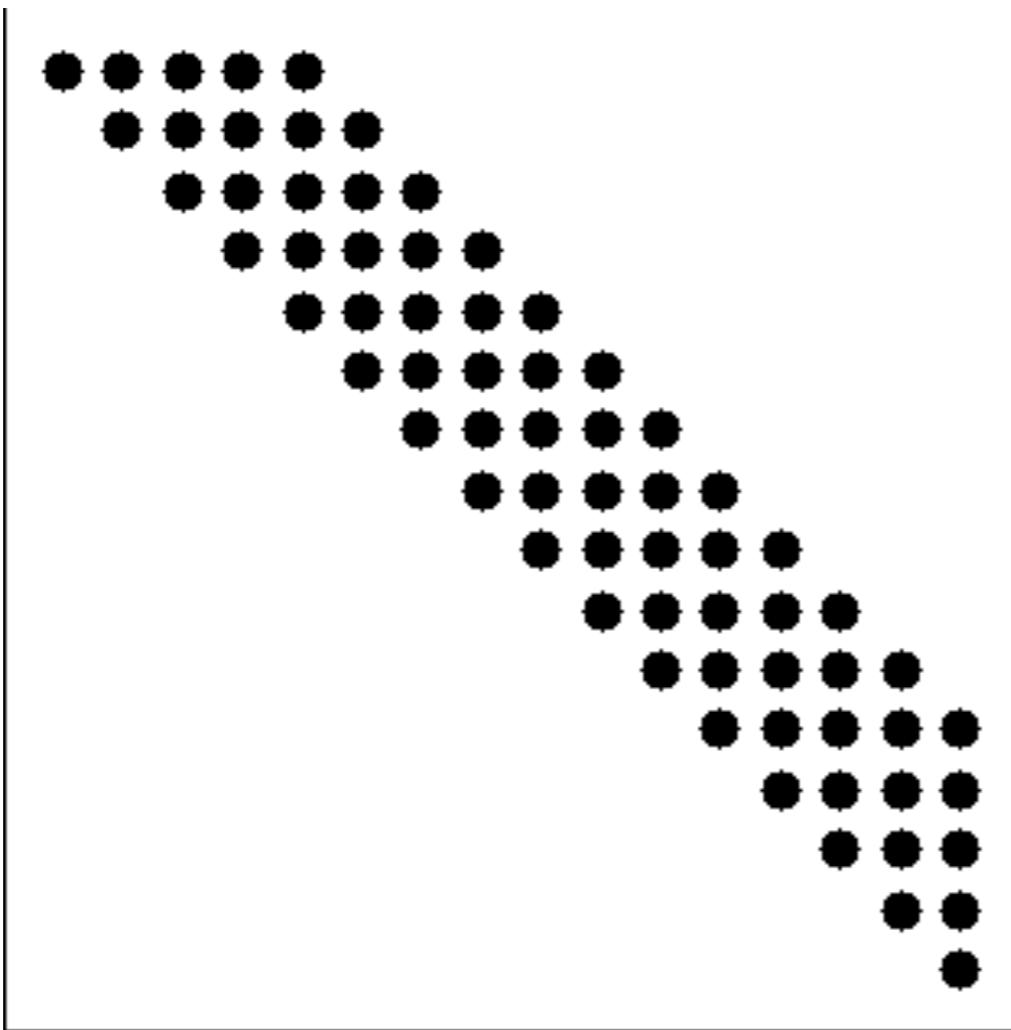
One stage reduction



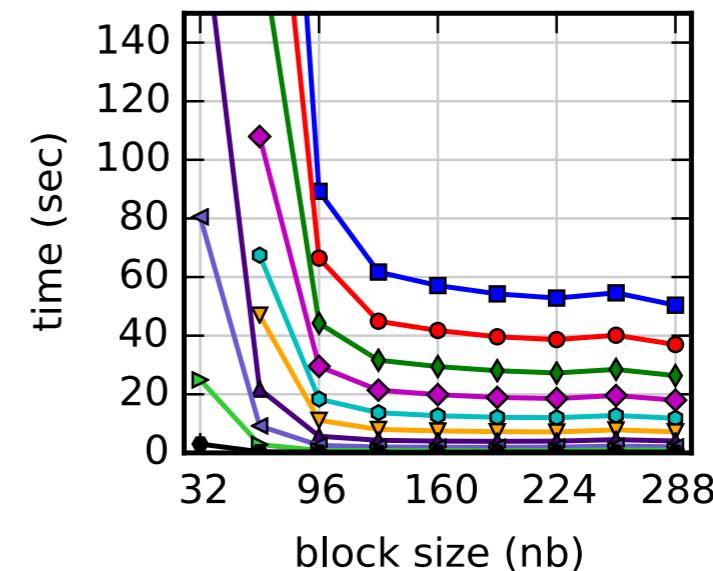
- $\frac{4}{3}n^3$ flops in BLAS 2 gemv
- $\frac{4}{3}n^3$ flops in BLAS 3 gemm
- Performance limited to twice gemv speed

- For $i = 1$ to n by n_b
 - Eliminate cols $i : i + n_b$ and rows $i : i + n_b$
 - Update trailing matrix using block Householder (BLAS 3 matrix multiplies)

Two stage reduction

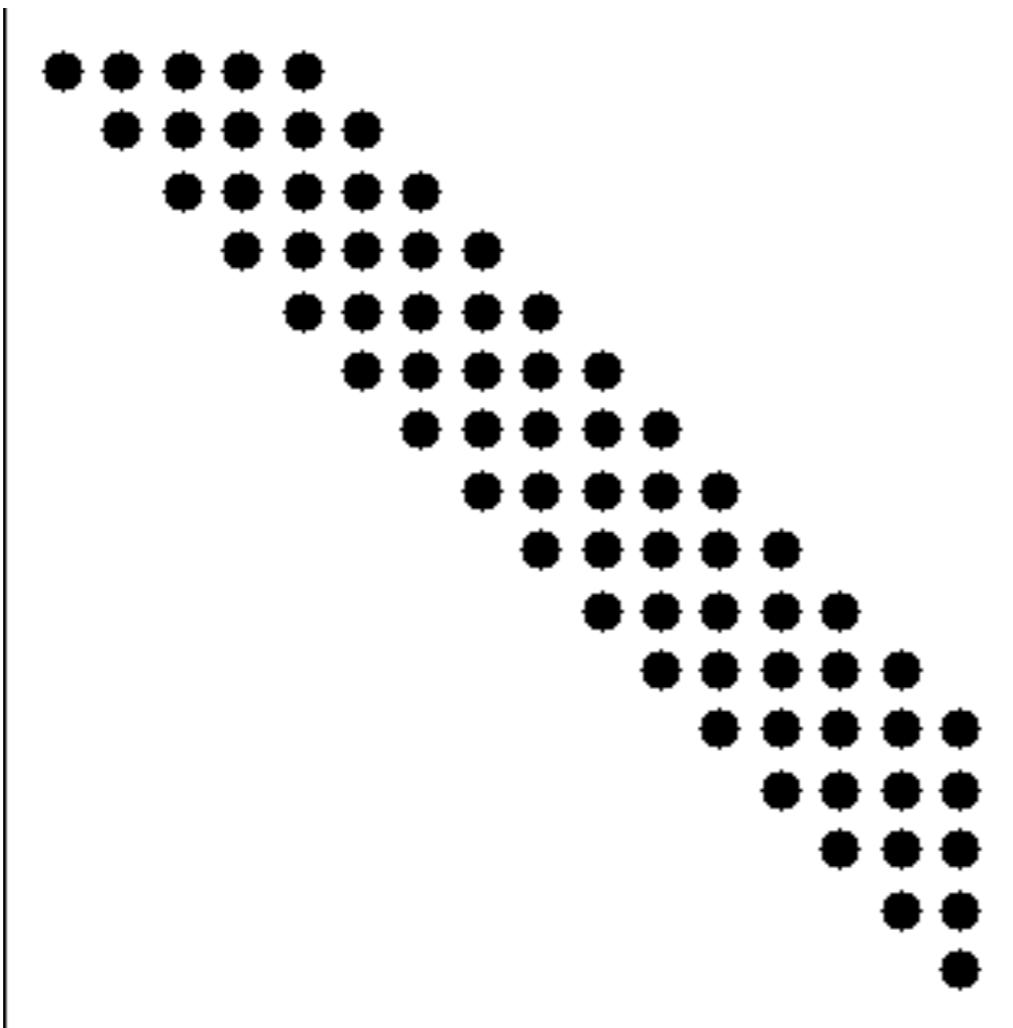
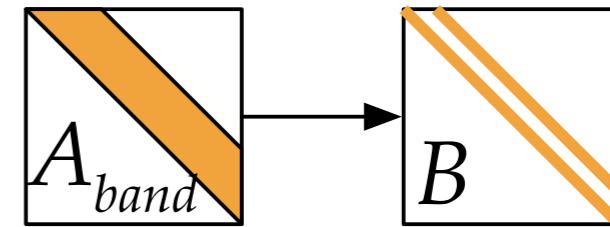


- $\frac{8}{3}n^3$ flops in BLAS 3 gemm
- More efficient with **large bandwidth n_b**

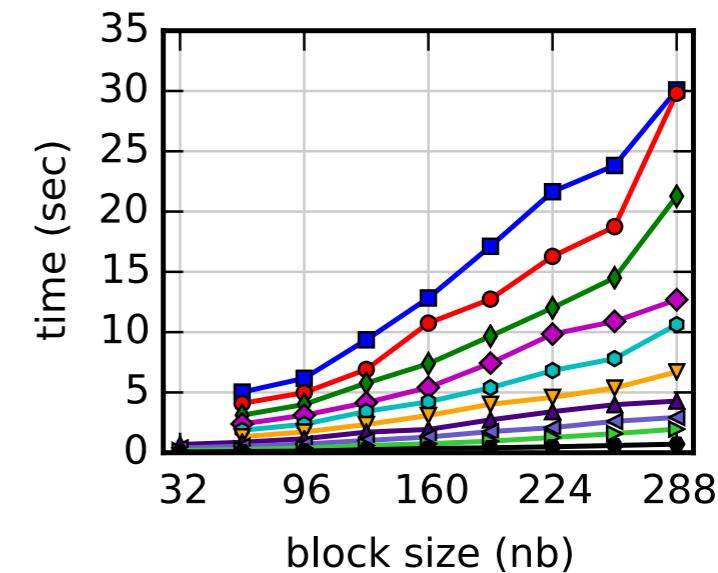


- **1st stage: full to band**
 - QR factorize panel, update trailing matrix
 - LQ factorize panel, update trailing matrix

Two stage reduction



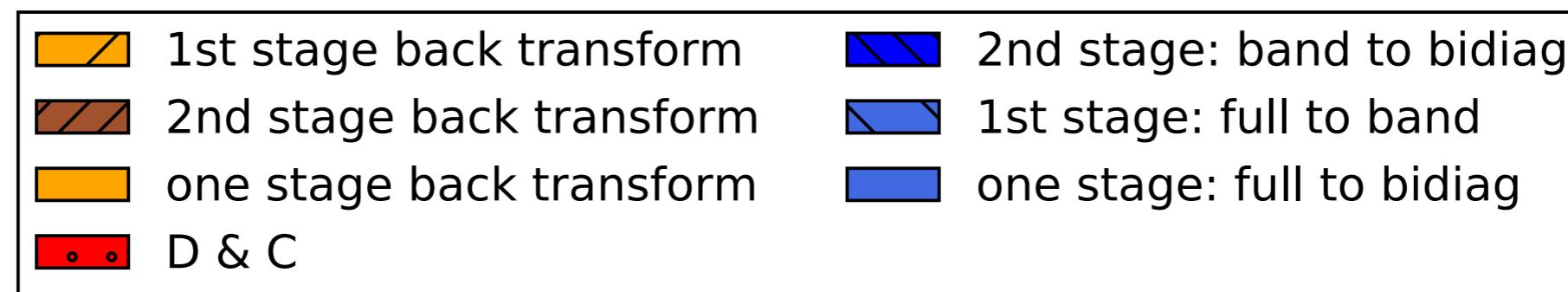
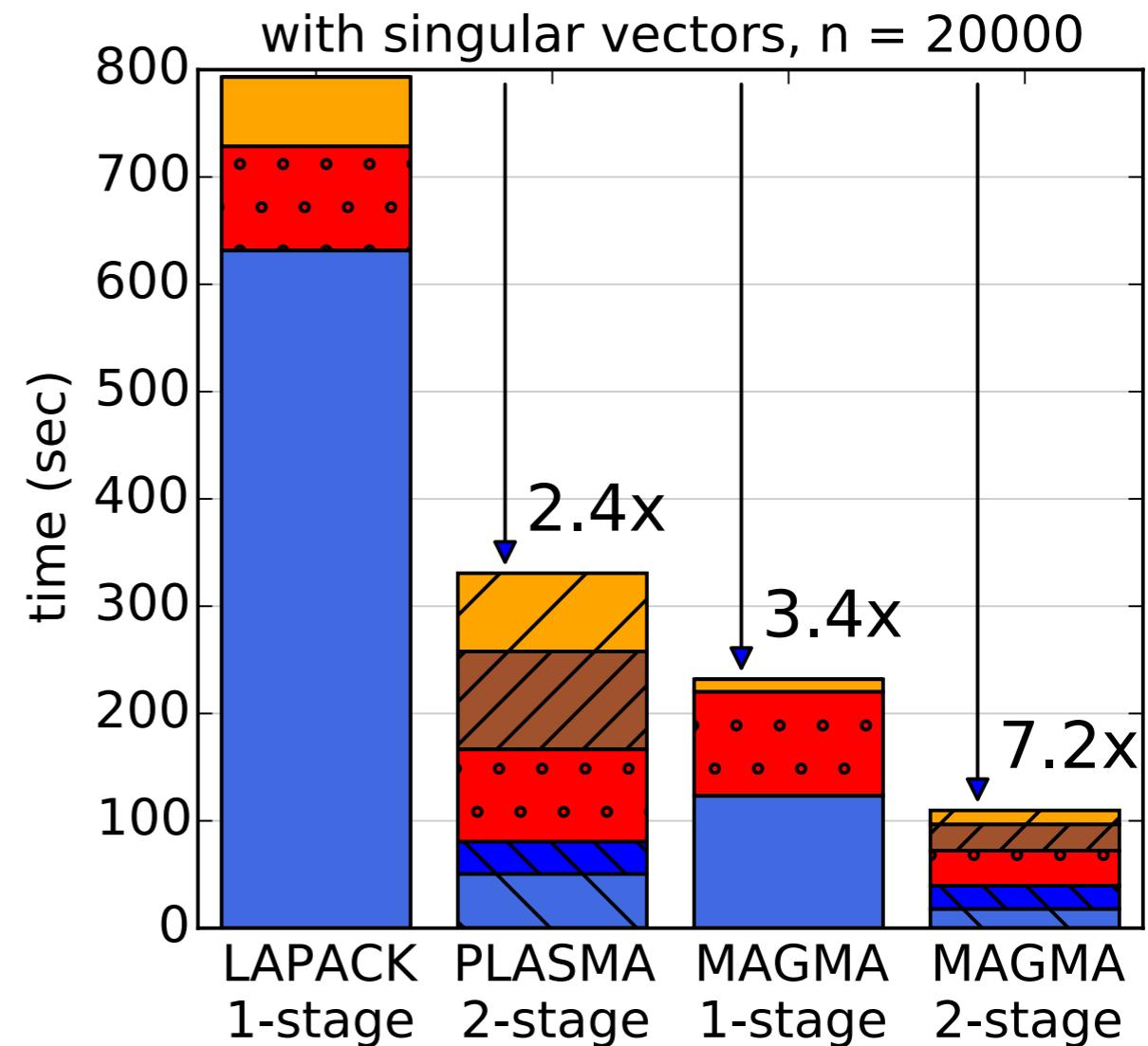
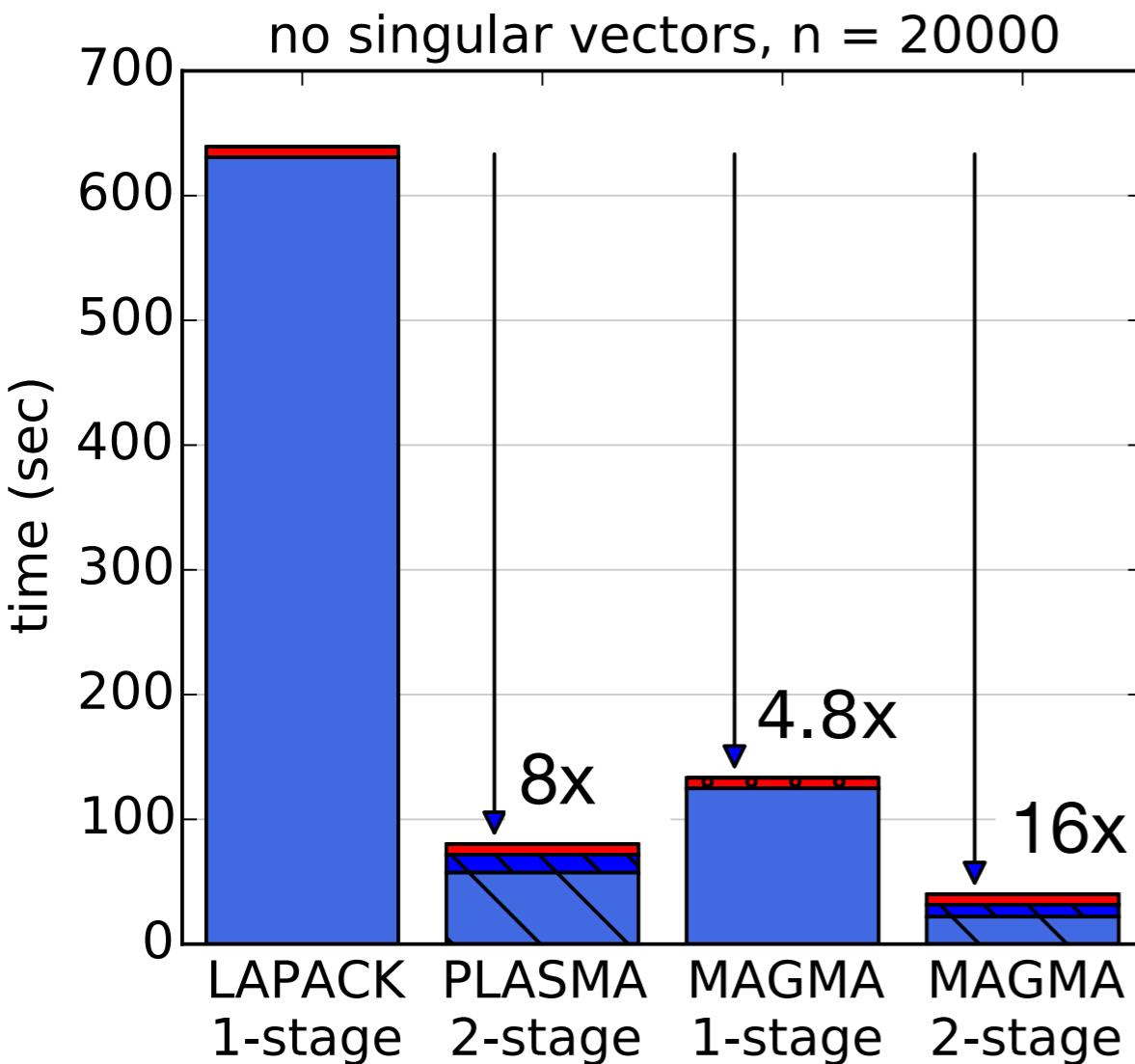
- $O(n_b n^2)$ flops in BLAS 2
- More efficient with **small bandwidth** n_b
- Tune optimal n_b



- **2nd stage: band to bidiagonal**
 - Cache friendly kernels
 - Pipeline multiple sweeps in parallel

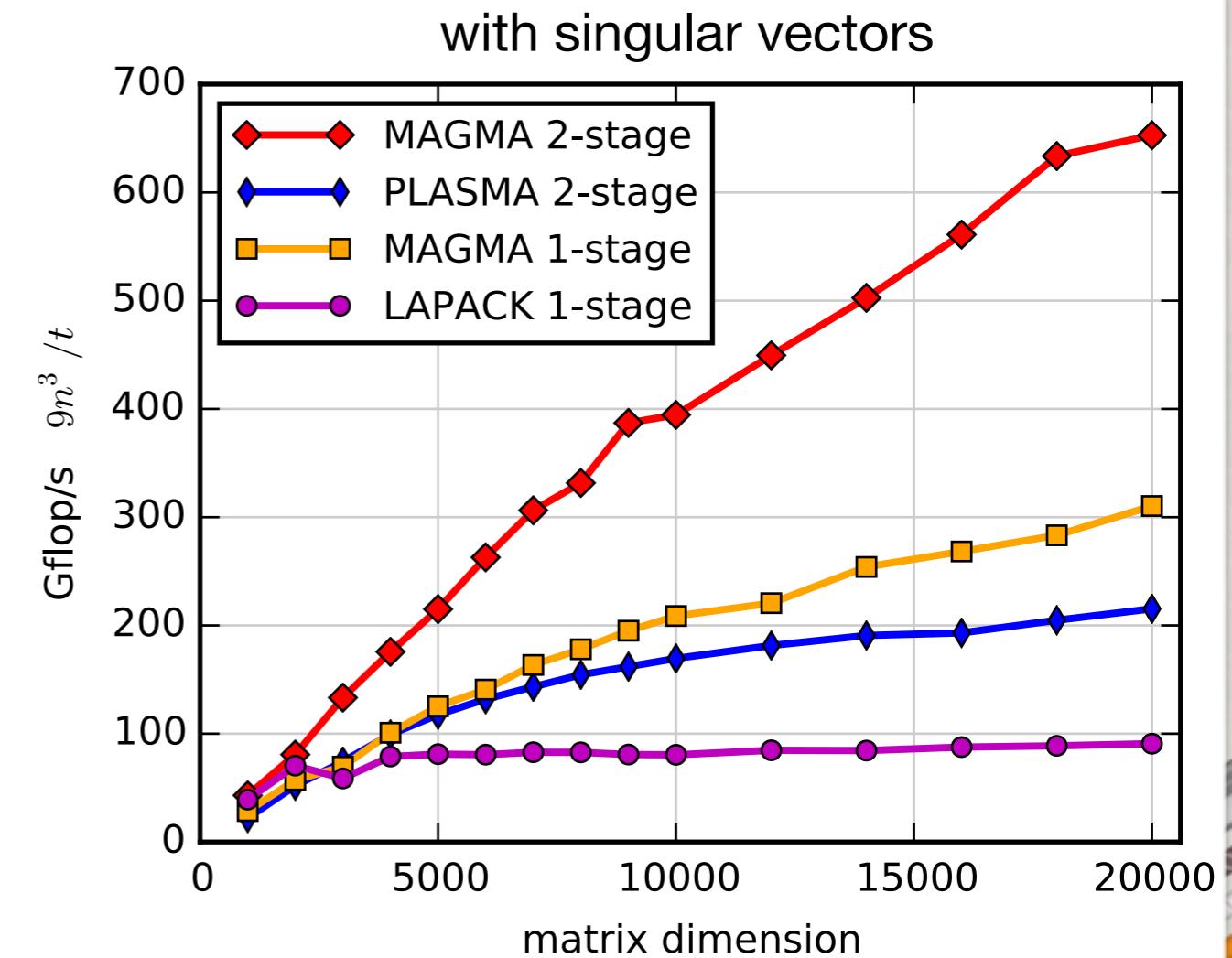
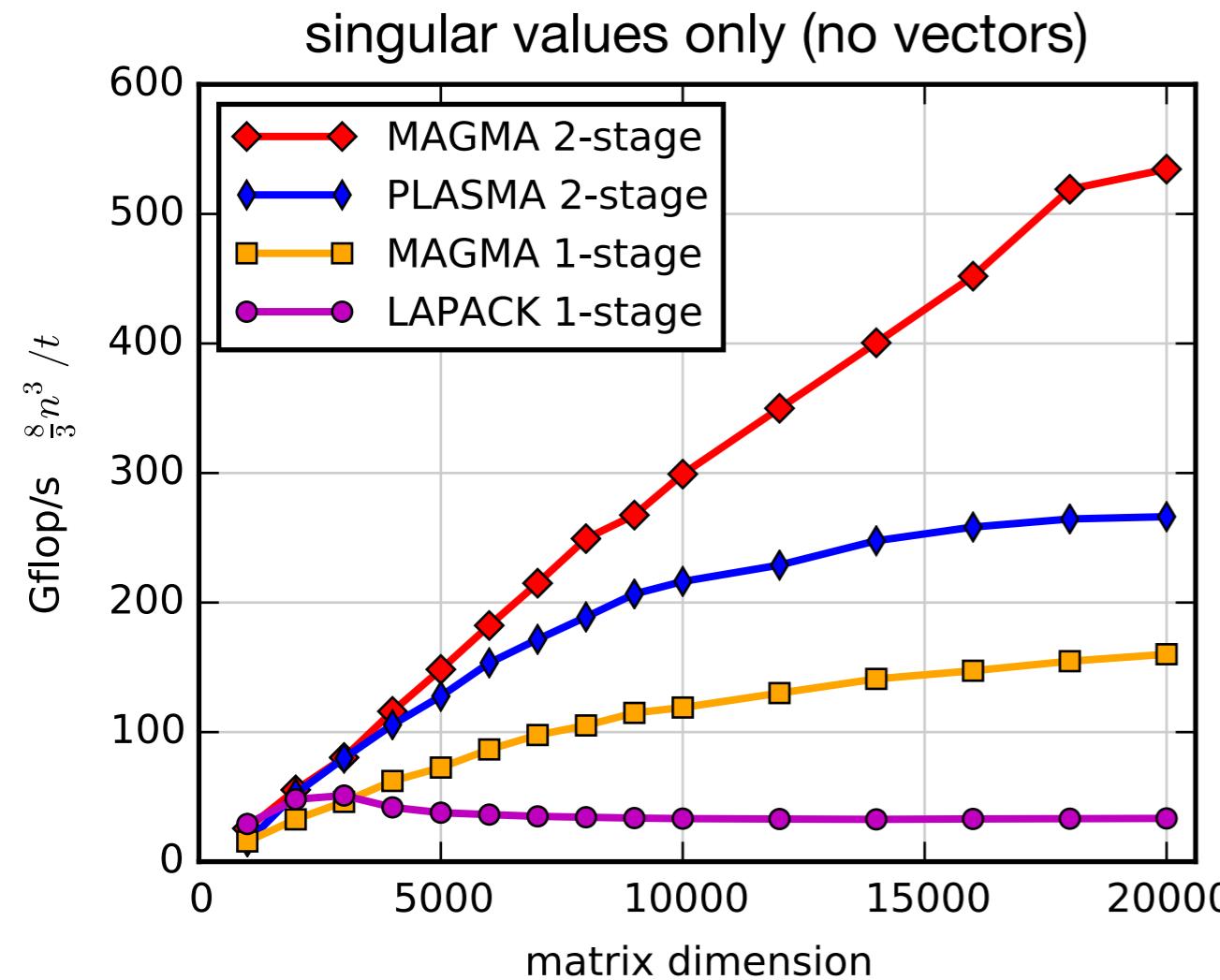
Accelerated phases

↓ lower is better



Overall performance

↑ higher is better





Questions?



THE UNIVERSITY OF
TENNESSEE
KNOXVILLE