

# Communication avoiding LU and QR factorizations

Laura Grigori

ALPINES

INRIA Rocquencourt - LJLL, UPMC

On sabbatical at UC Berkeley

# Motivation - the communication wall

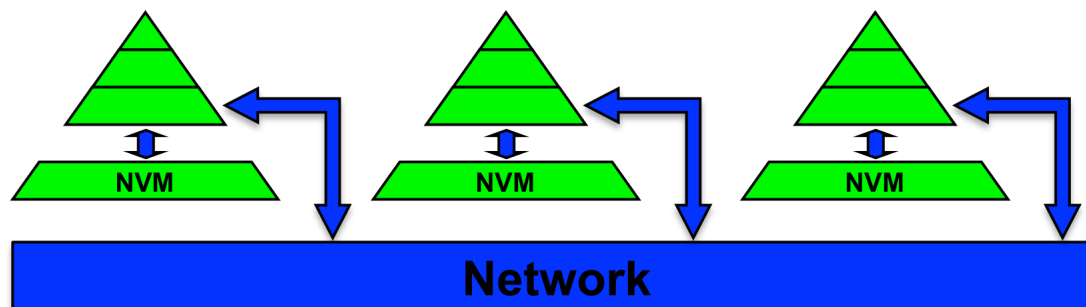
- Time to move data  $\gg$  time per flop
  - Gap steadily and exponentially growing over time

Annual improvements			
Time/flop		Bandwidth	Latency
59%	Network	26%	15%
	DRAM	23%	5.5%

*“Getting up to speed, The future of supercomputing” 2004, data from 1995-2004*

*“We are going to hit the **memory wall**, unless something basic changes”*

[W. Wulf, S. McKee, 95]



# Compelling numbers (1)

## DRAM bandwidth:

- Mid 90's ~ 0.2 bytes/flop – 1 byte/flop
- Past few years ~ 0.02 to 0.05 bytes/flop

## DRAM latency:

- DDR2 (2007) ~ 120 ns 1x
- DDR4 (2014) ~ 45 ns 2.6x in 7 years
- Stacked memory ~ similar to DDR4

## Time/flop

- 2006 Intel Yonah ~ 2GHz x 2 cores (32 GFlops/chip) 1x
- 2015 Intel Haswell ~2.3GHz x 16 cores (588 GFlops/chip) 18x in 9 years

# The role of numerical linear algebra

- Challenging applications often rely on solving linear algebra problems
- Linear systems of equations

Solve  $Ax = b$ , where  $A \in \mathbf{R}^{n \times n}$ ,  $b \in \mathbf{R}^n$ ,  $x \in \mathbf{R}^n$

- Direct methods

$PA = LU$ , then solve  $P^T L U x = b$

LU factorization is backward stable,

$\|PA - \widehat{L} \cdot \widehat{U}\|_{\infty}$  is small, close to machine epsilon in practice

- Iterative methods

- Find a solution  $x_k$  from  $x_0 + K_k(A, r_0)$ , where  $K_k(A, r_0) = \text{span} \{r_0, A r_0, \dots, A^{k-1} r_0\}$  such that the Petrov-Galerkin condition  $b - Ax_k \perp L_k$  is satisfied, where  $L_k$  is a subspace of dimension  $k$  and  $r_0 = Ax_0 - b$ .
- Convergence depends on  $\kappa(A)$  and the eigenvalue distribution (for SPD matrices).

# Least Square (LS) Problems

- Given  $A \in \mathbf{R}^{m \times n}$ ,  $b \in \mathbf{R}^m$ , solve  $\min_x \|Ax - b\|_2$ .
- Any solution of the LS problem satisfies the normal equations:  $A^T Ax = A^T b$
- Given the QR factorization of A

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix} \text{ where } \begin{array}{l} A \text{ is } m \times n \text{ real matrix, } m \geq n \\ R \text{ is } n \times n \text{ upper triangular matrix} \\ Q \text{ is } m \times m \text{ orthogonal matrix} \end{array}$$

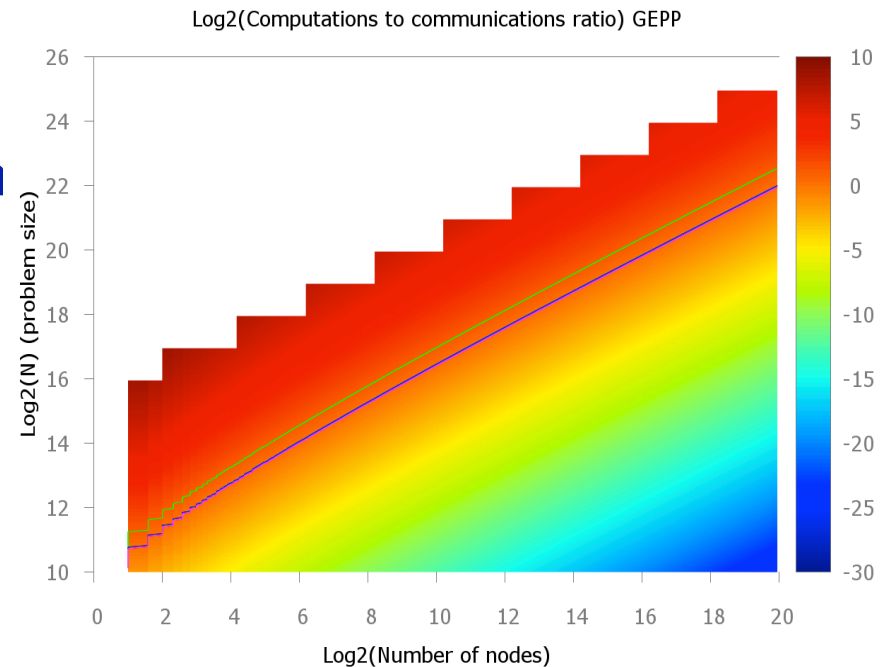
if  $\text{rank}(A) = \text{rank}(R) = n$ , then the LS solution is given by  $Rx = (Q^T b)(1:n)$

- The QR factorization is column-wise backward stable

$$\|A - \hat{Q}\hat{R}\|_2 \text{ is small, close to machine epsilon in practice}$$

# Approaches for reducing communication

- **Tuning**
  - Overlap communication and computation, at most a factor of 2 speedup
- **Same numerical algorithm, different schedule of the computation**
  - Block algorithms for NLA
    - Barron and Swinnerton-Dyer, 1960
    - ScaLAPACK, Blackford et al 97
  - Cache oblivious algorithms for NLA
    - Gustavson 97, Toledo 97, Frens and Wise 03, Ahmed and Pingali 00
- **Same algebraic framework, different numerical algorithm**
  - The approach used in CA algorithms
  - More opportunities for reducing communication, may affect stability



# Motivation

- The communication problem needs to be taken into account higher in the computing stack
- A paradigm shift in the way the numerical algorithms are devised is required
- Communication avoiding algorithms - a novel perspective for numerical linear algebra
  - Minimize volume of communication
  - Minimize number of messages
  - Minimize over multiple levels of memory/parallelism
  - Allow redundant computations (preferably as a low order term)

## Communication Complexity of Dense Linear Algebra

- Matrix multiply, using  $2n^3$  flops (sequential or parallel)
  - Hong-Kung (1981), Irony/Tishkin/Toledo (2004)
  - Lower bound on Bandwidth =  $\Omega(\text{\#flops} / M^{1/2})$
  - Lower bound on Latency =  $\Omega(\text{\#flops} / M^{3/2})$
- Same lower bounds apply to LU using reduction
  - Demmel, LG, Hoemmen, Langou 2008

$$\begin{pmatrix} I & & -B \\ A & I & \\ & & I \end{pmatrix} = \begin{pmatrix} I & & \\ A & I & \\ & & I \end{pmatrix} \begin{pmatrix} I & -B \\ & I & AB \\ & & I \end{pmatrix}$$

- And to almost all direct linear algebra [Ballard, Demmel, Holtz, Schwartz, 09]



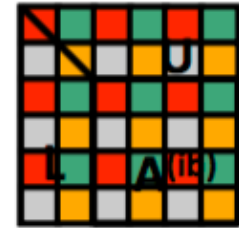
# Sequential algorithms and communication bounds

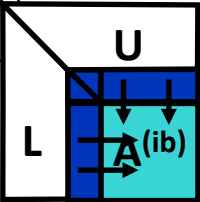
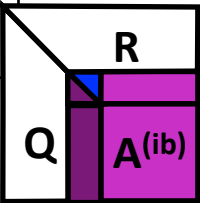
Algorithm	Minimizing #words (not #messages)	Minimizing #words and #messages
Cholesky	LAPACK	[Gustavson, 97] [Ahmed, Pingali, 00]
LU	LAPACK (few cases) [Toledo,97], [Gustavson, 97] both use partial pivoting	[LG, Demmel, Xiang, 08] [Khabou, Demmel, LG, Gu, 12] uses tournament pivoting
QR	LAPACK (few cases) [Elmroth,Gustavson,98]	[Frens, Wise, 03], 3x flops [Demmel, LG, Hoemmen, Langou, 08] [Ballard et al, 14]
RRQR		[Demmel, LG, Gu, Xiang 11] uses tournament pivoting, 3x flops

- Only several references shown for block algorithms (LAPACK), **cache-oblivious algorithms** and **communication avoiding algorithms**
- **CA algorithms** exist also for SVD and eigenvalue computation

# 2D Parallel algorithms and communication bounds

- If memory per processor =  $n^2 / P$ , the lower bounds become  
 $\#words\_moved \geq \Omega ( n^2 / P^{1/2} )$ ,  $\#messages \geq \Omega ( P^{1/2} )$



Algorithm	Minimizing #words (not #messages)	Minimizing #words and #messages
Cholesky	ScaLAPACK	ScaLAPACK
LU	 ScaLAPACK uses partial pivoting	[LG, Demmel, Xiang, 08] [Khabou, Demmel, LG, Gu, 12] uses tournament pivoting
QR	ScaLAPACK	[Demmel, LG, Hoemmen, Langou, 08] [Ballard et al, 14]
RRQR	 ScaLAPACK	[Demmel, LG, Gu, Xiang 13] uses tournament pivoting, 3x flops

- Only several references shown, block algorithms (ScaLAPACK) and communication avoiding algorithms
- CA algorithms exist also for SVD and eigenvalue computation

# Scalability of communication optimal algorithms

- 2D communication optimal algorithms,  $M = 3 \cdot n^2/P$   
(matrix distributed over a  $P^{1/2}$ -by-  $P^{1/2}$  grid of processors)  
 $T_P = O(n^3/P) \gamma + \Omega(n^2/P^{1/2}) \beta + \Omega(P^{1/2}) \alpha$ 
  - Isoefficiency:  $n^3 \propto P^{1.5}$  and  $n^2 \propto P$
  - For GEPP,  $n^3 \propto P^{2.25}$  [Grama et al, 93]
- 3D communication optimal algorithms,  $M = 3 \cdot P^{1/3}(n^2/P)$   
(matrix distributed over a  $P^{1/3}$ -by-  $P^{1/3}$ -by-  $P^{1/3}$  grid of processors)  
 $T_P = O(n^3/P) \gamma + \Omega(n^2/P^{2/3}) \beta + \Omega(\log(P)) \alpha$ 
  - Isoefficiency:  $n^3 \propto P$  and  $n^2 \propto P^{2/3}$
- 2.5D algorithms with  $M = 3 \cdot c \cdot (n^2/P)$ , and 3D algorithms exist for matrix multiplication and LU factorization
  - References: Dekel et al 81, Agarwal et al 90, 95, Johnsson 93, McColl and Tiskin 99, Irony and Toledo 02, Solomonik and Demmel 2011

E - the ratio between execution time on a single processor and total execution time summed over P processors.

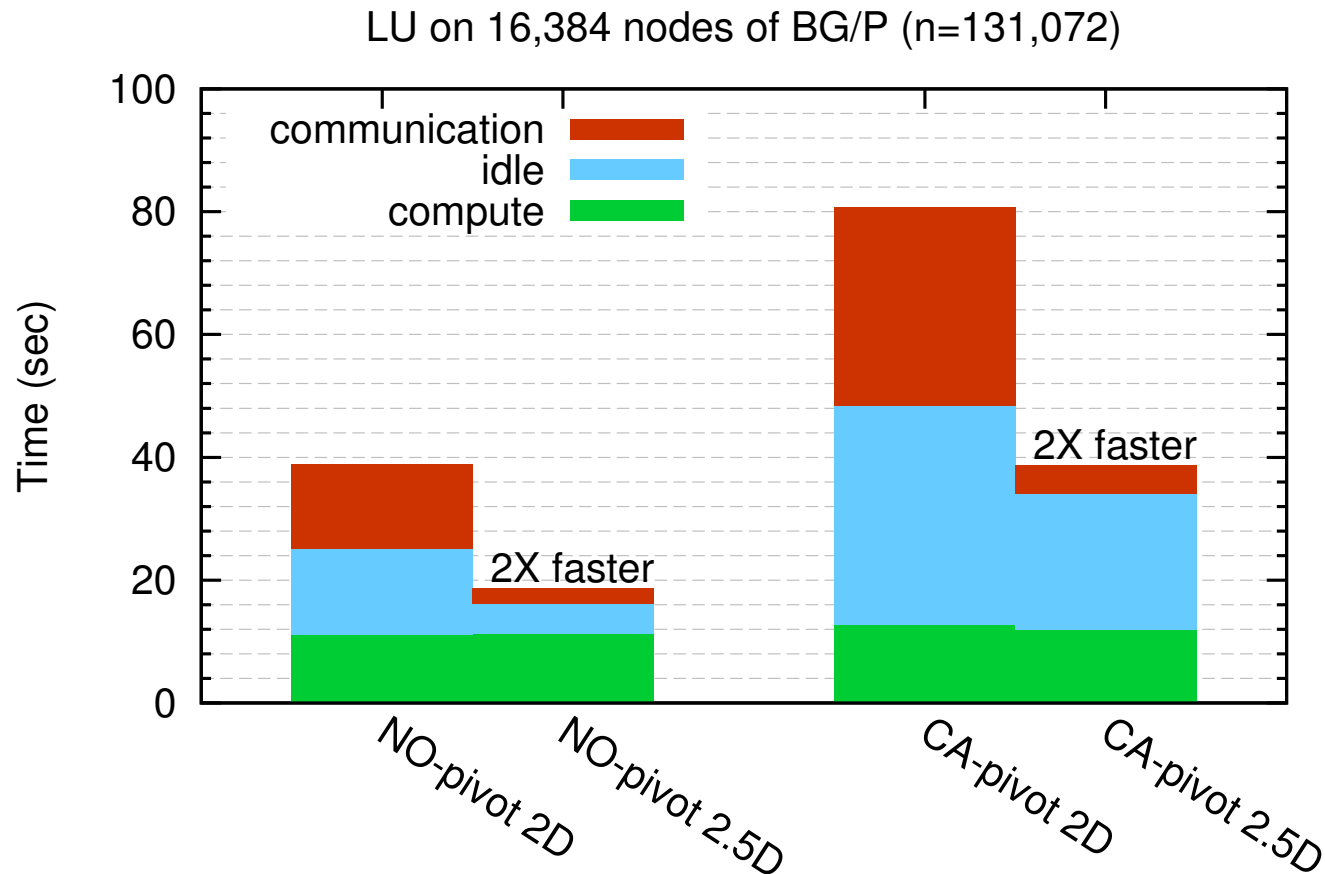
Isoefficiency - how the amount of computation must scale with P to keep E constant.

## 2.5D algorithms for LU, QR

- Assume  $c > 1$  copies of data, memory per processor is  $M \approx c \cdot (n^2/P)$
- For matrix multiplication
  - The bandwidth is reduced by a factor of  $c^{1/2}$
  - The latency is reduced by a factor of  $c^{3/2}$
  - Perfect Strong Scaling regime, given  $P$  such that  $M = 3n^2 / P$   
 $T(cP) = T(P)/c$
- For LU, QR
  - The bandwidth can be reduced by a factor of  $c^{1/2}$
  - But then the latency will increase by a factor of  $c^{1/2}$
  - Thm [Solomonik et al]: Perfect Strong Scaling impossible for LU, because  
 $\text{Latency} \cdot \text{Bandwidth} = \Omega(n^2)$
  - Conjecture: this applies to other factorizations as QR, RRQR, etc.

## 2.5D LU with and without pivoting

- 2.5D algorithms with  $M = 3 \cdot c \cdot (n^2/P)$ , and 3D algorithms exist for matrix multiplication and LU factorization
  - References: Dekel et al 81, Agarwal et al 90, 95, Johnsson 93, McColl and Tiskin 99, Irony and Toledo 02, Solomonik and Demmel 2011 (data presented below)



# The algebra of LU factorization

- Compute the factorization  $PA = LU$
- Given the matrix

$$A = \begin{pmatrix} 3 & 1 & 3 \\ 6 & 7 & 3 \\ 9 & 12 & 3 \end{pmatrix}$$

Let

$$M_1 A = \begin{pmatrix} 1 & & \\ -2 & 1 & \\ -3 & & 1 \end{pmatrix}, \quad M_1 A = \begin{pmatrix} 3 & 1 & 3 \\ 0 & 5 & -3 \\ 0 & 9 & -6 \end{pmatrix}$$

## The need for pivoting

- For stability avoid division by small elements, otherwise  $\|A-LU\|$  can be large
  - Because of roundoff error
- For example

$$A = \begin{pmatrix} 0 & 3 & 3 \\ 3 & 1 & 3 \\ 6 & 2 & 3 \end{pmatrix}$$

has an LU factorization if we permute the rows of A

$$PA = \begin{pmatrix} 6 & 2 & 3 \\ 0 & 3 & 3 \\ 3 & 1 & 3 \end{pmatrix} = \begin{pmatrix} 1 & & \\ & 1 & \\ 0.5 & & 1 \end{pmatrix} \begin{pmatrix} 6 & 2 & 3 \\ 3 & 3 \\ & & 1.5 \end{pmatrix}$$

- Partial pivoting allows to bound all elements of L by 1.

# LU with partial pivoting – BLAS 2 algorithm

```

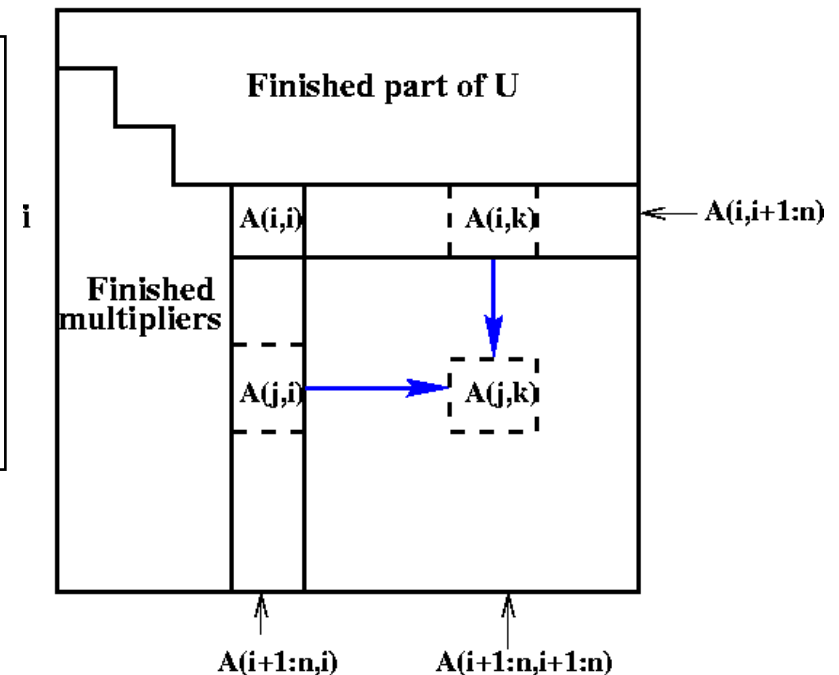
for i = 1 to n-1
  Let A(j,i) be elt. of max magnitude in A(i+1:n,i)
  Permute rows i and j
  for j = i+1 to n
    A(j,i) = A(j,i)/A(i,i)
  for j = i+1 to n
    for k = i+1 to n
      A(j,k) = A(j,k) - A(j,i) * A(i,k)
  
```

- Algorithm using BLAS 1/2 operations

```

for i = 1 to n-1
  Let A(j,i) be of elt. of max magnitude in A(i+1:n,i)
  Permute rows i and j
  A(i+1:n,i) = A(i+1:n,i) * ( 1 / A(i,i) )
  ... BLAS 1 (scale a vector)
  A(i+1:n,i+1:n) = A(i+1:n , i+1:n )
  - A(i+1:n , i) * A(i , i+1:n)
  ... BLAS 2 (rank-1 update)
  
```

Work at step  $i$  of Gaussian Elimination





## Block LU factorization – obtained by delaying updates

- Matrix  $A$  of size  $n \times n$  is partitioned as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \text{ where } A_{11} \text{ is } b \times b$$

- The first step computes LU with partial pivoting of the first block:

$$P_1 \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix} = \begin{pmatrix} L_{11} \\ L_{21} \end{pmatrix} U_{11}$$

- The factorization obtained is:

$$P_1 A = \begin{pmatrix} L_{11} & \\ L_{21} & I_{n-b} \end{pmatrix} \begin{pmatrix} U_{11} & U_{12} \\ & A_{22}^1 \end{pmatrix}, \text{ where } \begin{aligned} U_{12} &= L_{11}^{-1} A_{12} \\ A_{22}^1 &= A_{22} - L_{21} U_{12} \end{aligned}$$

- The algorithm continues recursively on the trailing matrix  $A_{22}^1$

## Block LU factorization – the algorithm

1. Compute LU with partial pivoting of the first panel

$$P_1 \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix} = \begin{pmatrix} L_{11} \\ L_{21} \end{pmatrix} U_{11}$$

2. Pivot by applying the permutation matrix  $P_1$  on the entire matrix

$$P_1 A = \bar{A}$$

3. Solve the triangular system to compute a block row of U

$$U_{12} = L_{12}^{-1} \bar{A}_{12}$$

4. Update the trailing matrix

$$\bar{A}_{22}^1 = \bar{A}_{22} - L_{21} U_{12}$$

5. The algorithm continues recursively on the trailing matrix  $\bar{A}_{22}^1$

# LU factorization (as in ScaLAPACK pdgetrf)

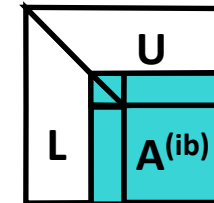


LU factorization on a  $P = P_r \times P_c$  grid of processors

For  $ib = 1$  to  $n-1$  step  $b$

$$A^{(ib)} = A(ib:n, ib:n)$$

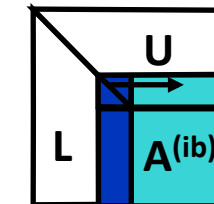
#messages



(1) Compute panel factorization

$$O(n \log_2 P_r)$$

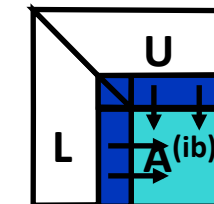
- find pivot in each column, swap rows



(2) Apply all row permutations

$$O(n/b(\log_2 P_c + \log_2 P_r))$$

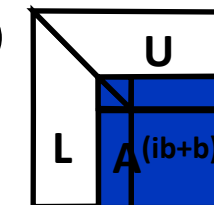
- broadcast pivot information along the rows
- swap rows at left and right



(3) Compute block row of U

$$O(n/b \log_2 P_c)$$

- broadcast right diagonal block of L of current panel



(4) Update trailing matrix

$$O(n/b(\log_2 P_c + \log_2 P_r))$$

- broadcast right block column of L
- broadcast down block row of U

# General scheme for QR factorization by Householder transformations

The Householder matrix

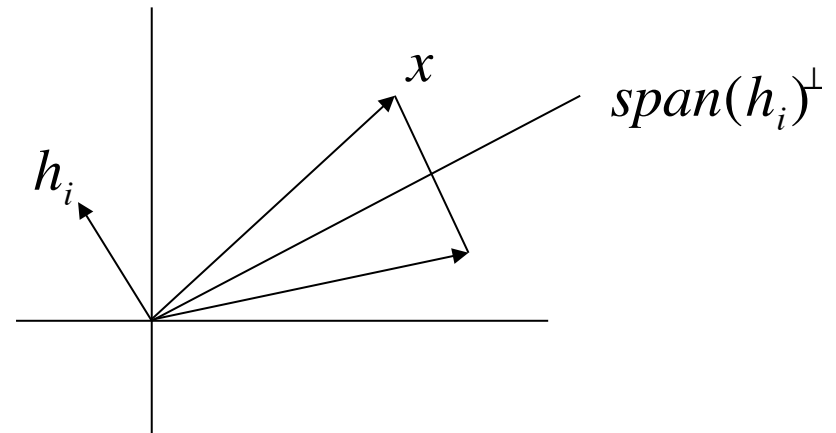
$$H_i = I - \tau_i h_i h_i^T$$

has the following properties:

- is symmetric and orthogonal,

$$H_i^2 = I,$$

- is independent of the scaling of  $h_i$ ,
- it reflects  $x$  about the hyperplane  $\text{span}(h_i)^\perp$



- For QR, we choose a Householder matrix that allows to annihilate the elements of a vector  $x$ , except first one.

## General scheme for QR factorization by Householder transformations

- Apply Householder transformations to annihilate subdiagonal entries

$$\begin{aligned}
 A &= \begin{pmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{pmatrix} = H_1 \begin{pmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & x & x & x \\ 0 & x & x & x \end{pmatrix} = H_1 \begin{pmatrix} 1 & & & \\ & \tilde{H}_2 & & \\ & & & \\ & & & \end{pmatrix} \begin{pmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{pmatrix} = \\
 &= H_1 H_2 \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \tilde{H}_3 & \\ & & & \end{pmatrix} \begin{pmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{pmatrix} = H_1 H_2 H_3 R = QR
 \end{aligned}$$

- For A of size mxn, the factorization can be written as:

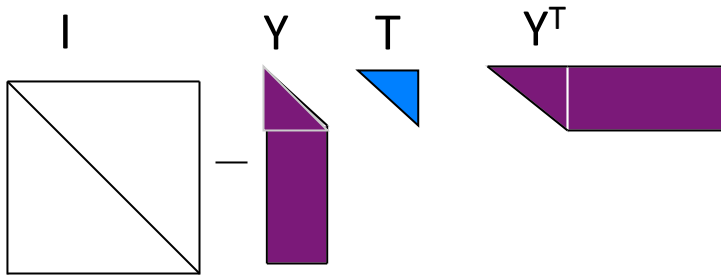
$$\begin{aligned}
 H_n H_{n-1} \dots H_2 H_1 A &= R \rightarrow A = (H_n H_{n-1} \dots H_2 H_1)^T R \\
 Q &= H_1 H_2 \dots H_n
 \end{aligned}$$

## Compact representation for Q

- Orthogonal factor Q can be represented implicitly as

$$Q = H_1 H_2 \dots H_b = (I - \tau_1 h_1 h_1^T) \dots (I - \tau_b h_b h_b^T) = I - YTY^T, \text{ where}$$

$$Y = (h_1 \quad h_2 \quad \dots \quad h_b)$$



- Example for  $b=2$ :

$$Y = (h_1 | h_2), \quad T = \begin{pmatrix} \tau_1 & -\tau_1 h_1^T h_2 \tau_2 \\ & \tau_2 \end{pmatrix}$$

# Algebra of block QR factorization

Matrix  $A$  of size  $n \times n$  is partitioned as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \text{ where } A_{11} \text{ is } b \times b$$

## Block QR algebra

The first step of the block QR factorization algorithm computes:

$$Q_1^T A = \begin{bmatrix} R_{11} & R_{12} \\ & A_{22}^1 \end{bmatrix}$$

The algorithm continues recursively on the trailing matrix  $A_{22}^1$

# Block QR factorization

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} = Q_1 \begin{pmatrix} R_{11} & R_{12} \\ & A_{22}^1 \end{pmatrix}$$

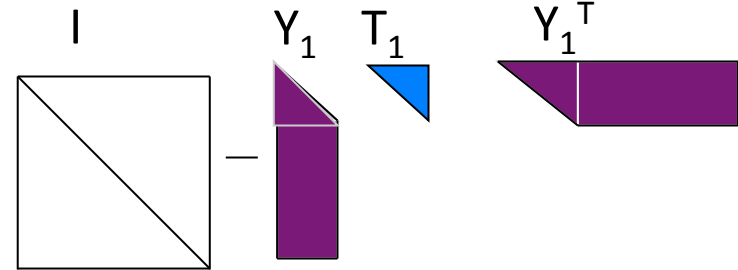
Block QR algebra:

1. Compute panel factorization:

$$\begin{pmatrix} A_{11} \\ A_{12} \end{pmatrix} = Q_1 \begin{pmatrix} R_{11} \\ \end{pmatrix}, \quad Q_1 = H_1 H_2 \dots H_b$$

2. Compute the compact representation:

$$Q_1 = I - Y_1 T_1 Y_1^T$$



3. Update the trailing matrix:

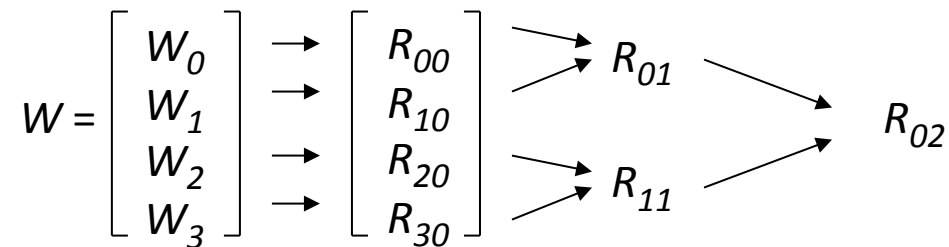
$$\left( I - Y_1 T_1^T Y_1^T \right) \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix} = \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix} - Y_1 \left( T_1^T \left( Y_1^T \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix} \right) \right) = \begin{pmatrix} R_{12} \\ A_{22}^1 \end{pmatrix}$$

4. The algorithm continues recursively on the trailing matrix.

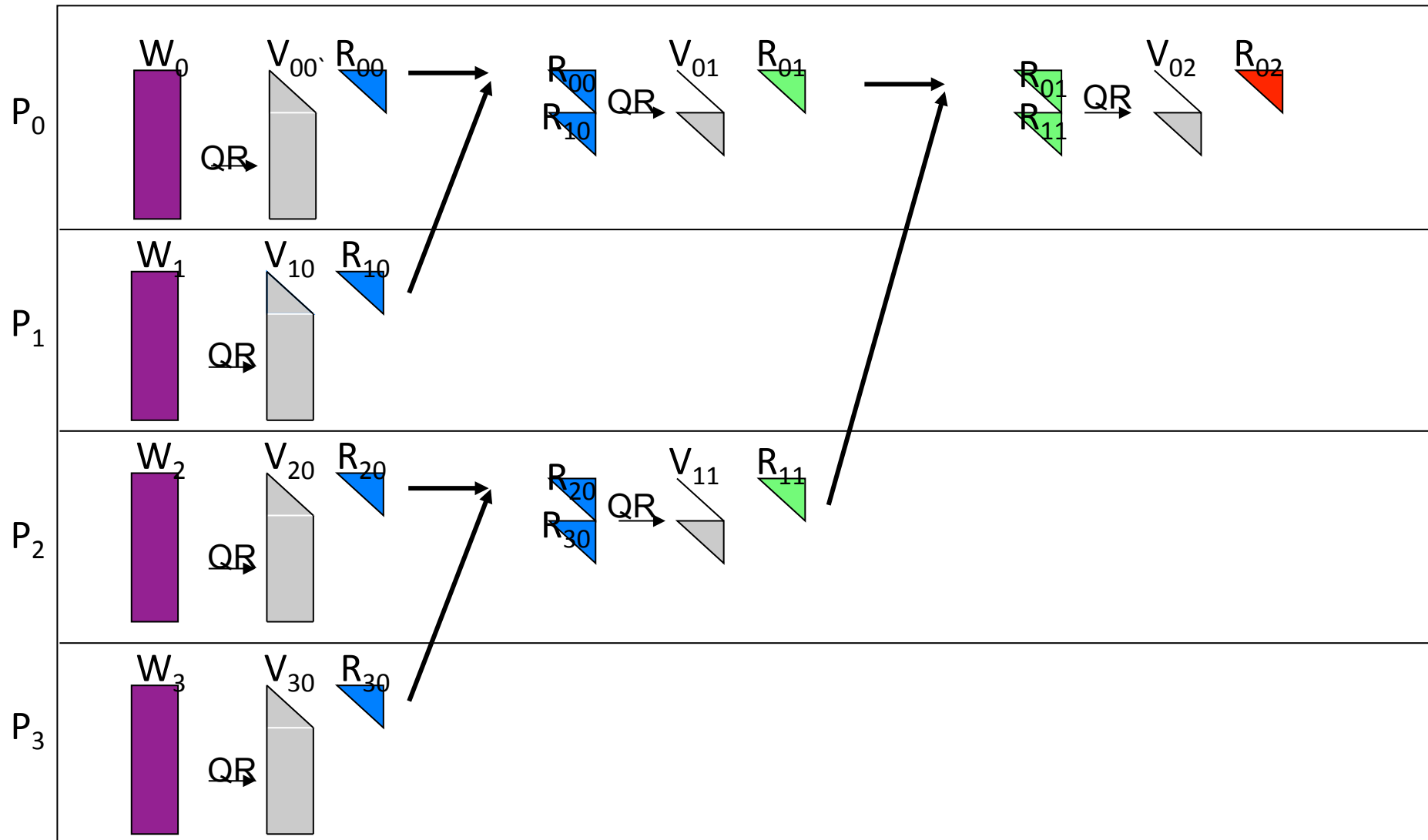


# TSQR: QR factorization of a tall skinny matrix using Householder transformations

- QR decomposition of  $m \times b$  matrix  $W$ ,  $m \gg b$ 
  - $P$  processors, block row layout
- Classic Parallel Algorithm
  - Compute Householder vector for each column
  - Number of messages  $\propto b \log P$
- Communication Avoiding Algorithm
  - Reduction operation, with QR as operator
  - Number of messages  $\propto \log P$

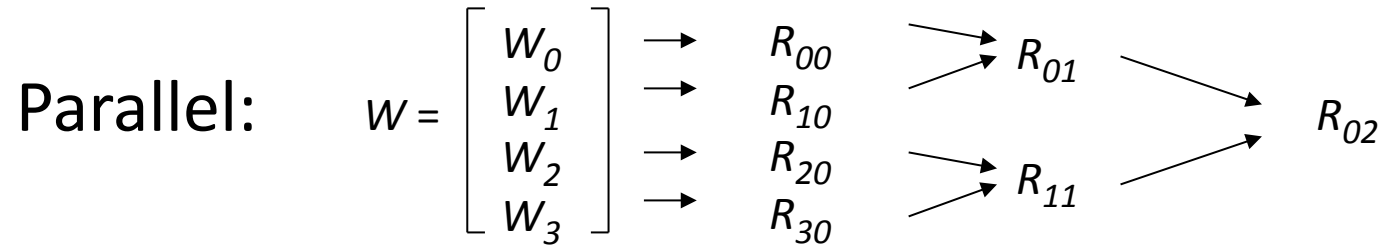


# Parallel TSQR



References: Golub, Plemmons, Sameh 88, Pothen, Raghavan, 89, Da Cunha, Becker, Patterson, 02

# Algebra of TSQR



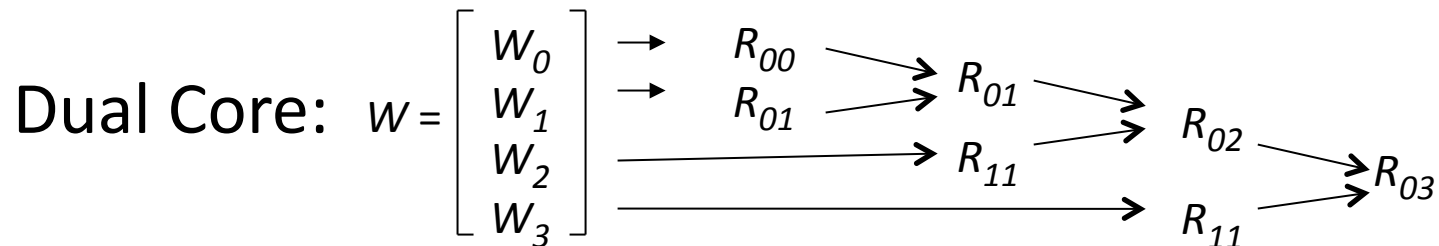
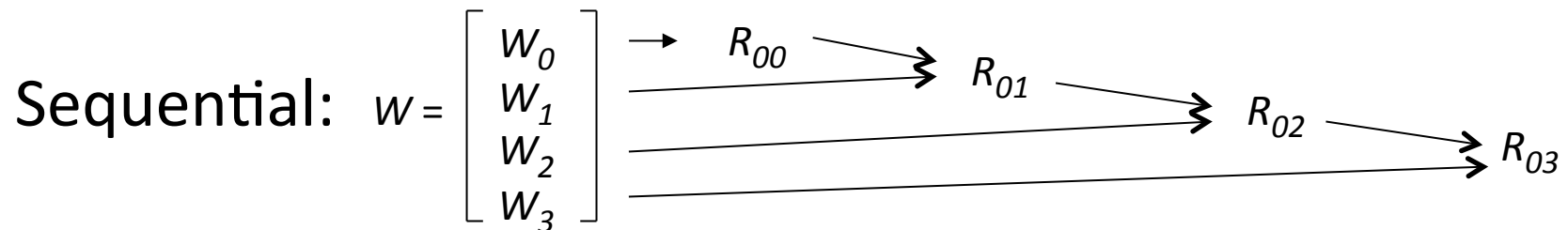
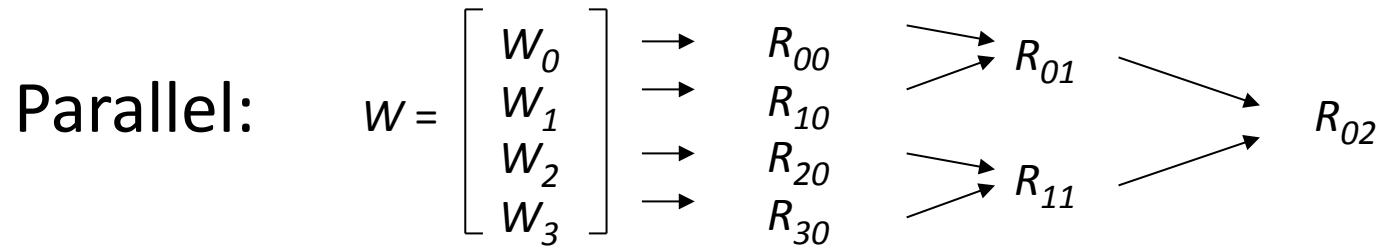
$$W = \begin{pmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{pmatrix} = \begin{pmatrix} \frac{Q_{00}R_{00}}{Q_{10}R_{10}} \\ \frac{Q_{20}R_{20}}{Q_{30}R_{30}} \end{pmatrix} = \begin{pmatrix} Q_{00} \\ \hline Q_{10} \\ \hline Q_{20} \\ \hline Q_{30} \end{pmatrix} \begin{pmatrix} R_{00} \\ R_{10} \\ R_{20} \\ R_{30} \end{pmatrix}$$

$$\begin{pmatrix} R_{00} \\ R_{10} \\ R_{20} \\ R_{30} \end{pmatrix} = \begin{pmatrix} \frac{Q_{01}R_{01}}{Q_{11}R_{11}} \end{pmatrix} = \begin{pmatrix} Q_{01} \\ \hline Q_{11} \end{pmatrix} \begin{pmatrix} R_{01} \\ R_{11} \end{pmatrix} \quad \begin{pmatrix} R_{01} \\ R_{11} \end{pmatrix} = Q_{02}R_{02}$$

Q is represented implicitly as a product

Output:  $\{Q_{00}, Q_{10}, Q_{00}, Q_{20}, Q_{30}, Q_{01}, Q_{11}, Q_{02}, R_{02}\}$

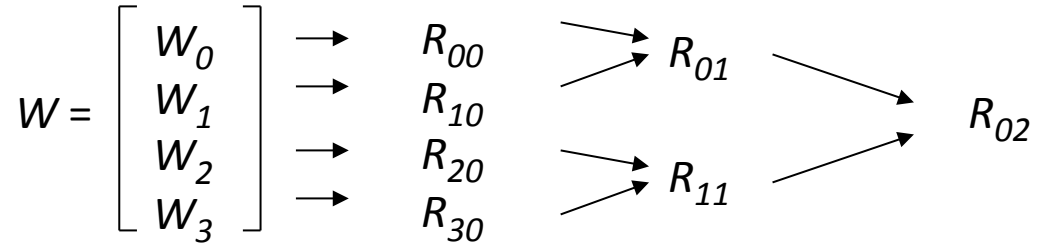
## Flexibility of TSQR and CAQR algorithms



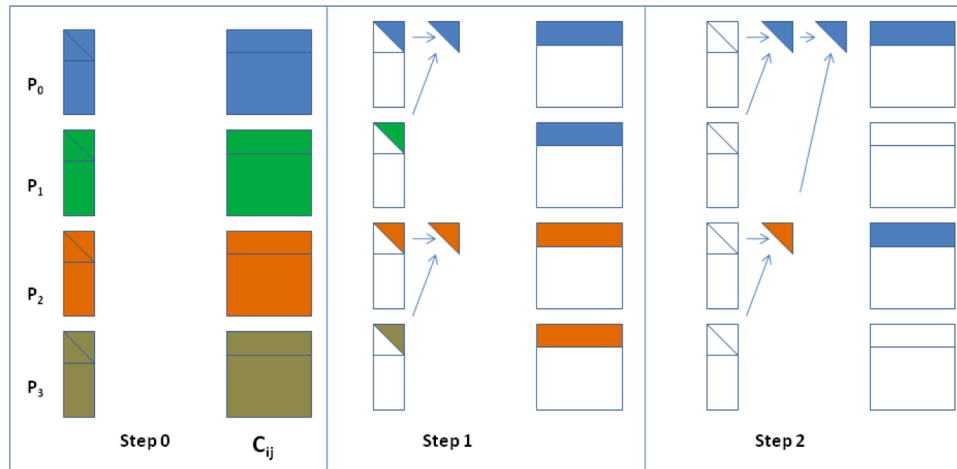
Reduction tree will depend on the underlying architecture,  
could be chosen dynamically

# Algebra of TSQR

Parallel:

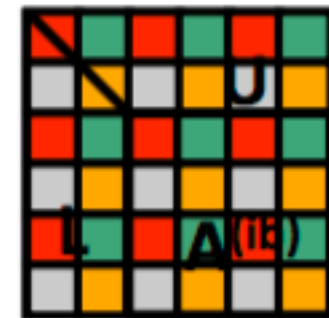


## CAQR



# QR for General Matrices

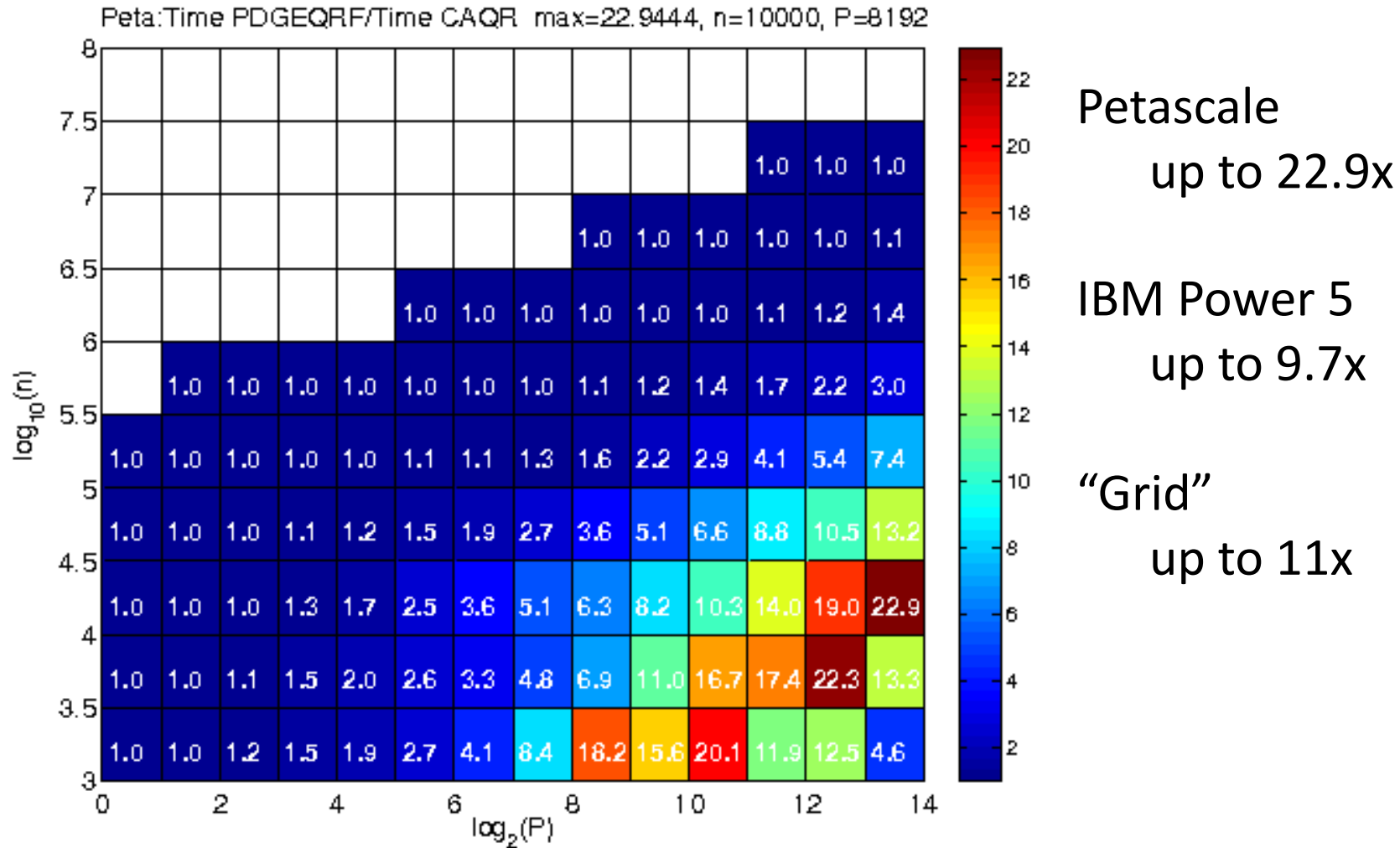
- Cost of **CAQR** vs **ScaLAPACK's PDGEQRF**
  - $n \times n$  matrix on  $P^{1/2} \times P^{1/2}$  processor grid, block size  $b$
  - Flops:  $(4/3)n^3/P + (3/4)n^2b \log P/P^{1/2}$  vs  $(4/3)n^3/P$
  - Bandwidth:  $(3/4)n^2 \log P/P^{1/2}$  vs **same**
  - Latency:  $2.5 n \log P / b$  vs  $1.5 n \log P$
- Close to optimal (modulo  $\log P$  factors)
  - Assume:  $O(n^2/P)$  memory/processor,  $O(n^3)$  algorithm,
  - Choose  $b$  near  $n / P^{1/2}$  (its upper bound)
  - Bandwidth lower bound:
    - $\Omega(n^2 / P^{1/2})$  – just  $\log(P)$  smaller
  - Latency lower bound:
    - $\Omega(P^{1/2})$  – just  $\text{polylog}(P)$  smaller



# Performance of TSQR vs Sca/LAPACK

- Parallel
  - Intel Xeon (two socket, quad core machine), 2010
    - Up to **5.3x speedup** (8 cores,  $10^5 \times 200$ )
  - Pentium III cluster, Dolphin Interconnect, MPICH, 2008
    - Up to **6.7x speedup** (16 procs,  $100K \times 200$ )
  - BlueGene/L, 2008
    - Up to **4x speedup** (32 procs,  $1M \times 50$ )
  - Tesla C 2050 / Fermi (Anderson et al)
    - Up to **13x** ( $110,592 \times 100$ )
  - Grid – **4x** on 4 cities vs 1 city (Dongarra, Langou et al)
  - QR computed locally using recursive algorithm (Elmroth-Gustavson) – enabled by TSQR
- Results from many papers, for some see [Demmel, LG, Hoemmen, Langou, SISC 12], [Donfack, LG, IPDPS 10].

# Modeled Speedups of CAQR vs ScaLAPACK

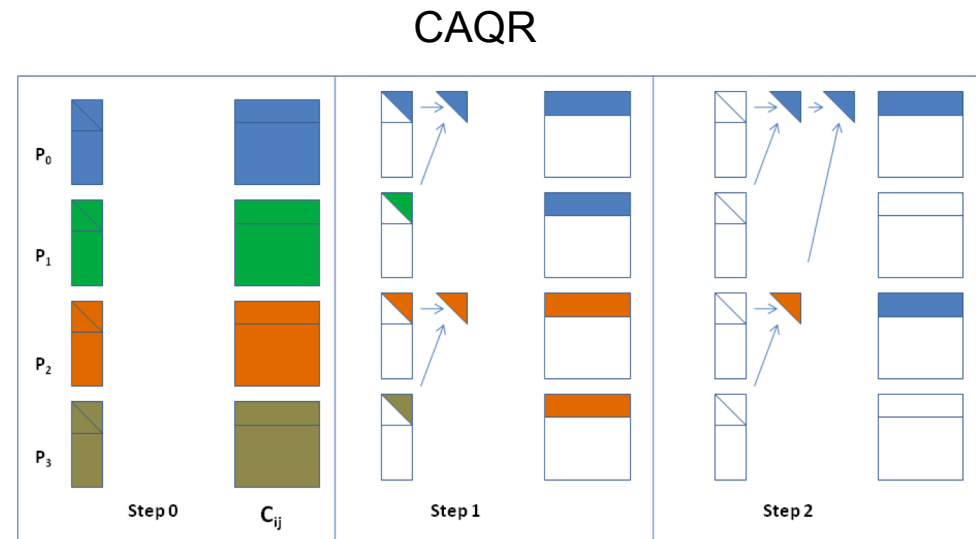
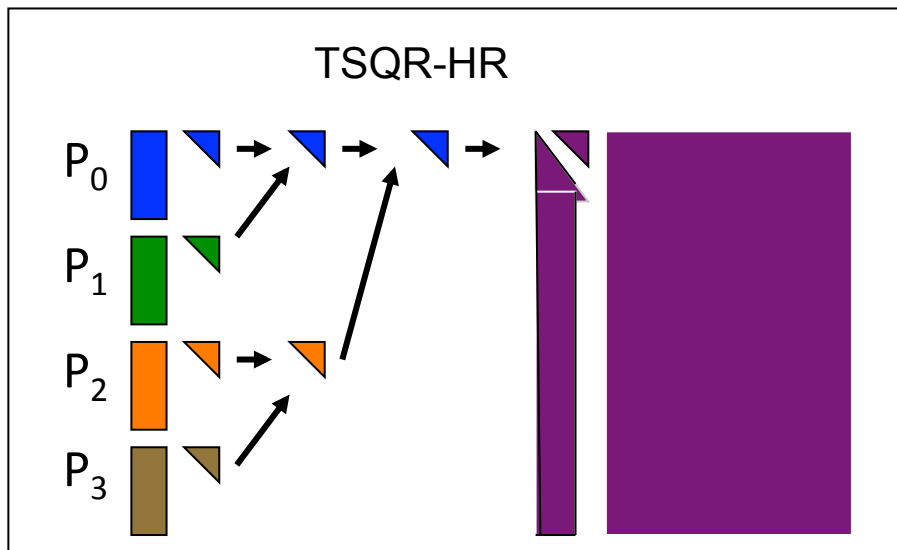
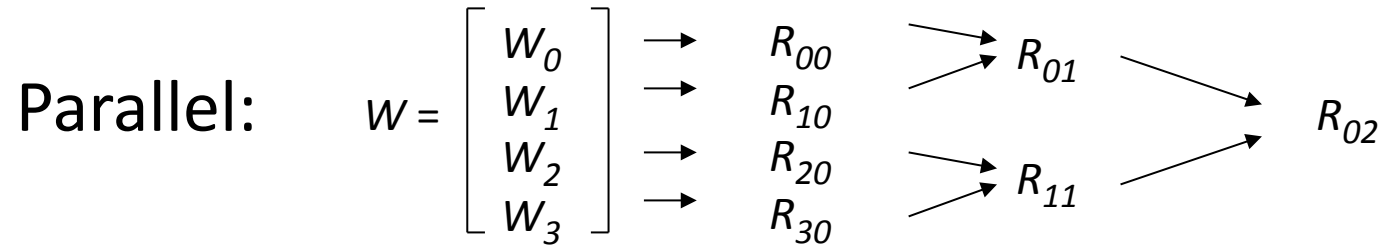


Petascale machine with 8192 procs, each at 500 GFlops/s, a bandwidth of 4 GB/s.

$$\gamma = 2 \cdot 10^{-12} s, \alpha = 10^{-5} s, \beta = 2 \cdot 10^{-9} s / \text{word}.$$



# Algebra of TSQR



# Reconstruct Householder vectors from TSQR

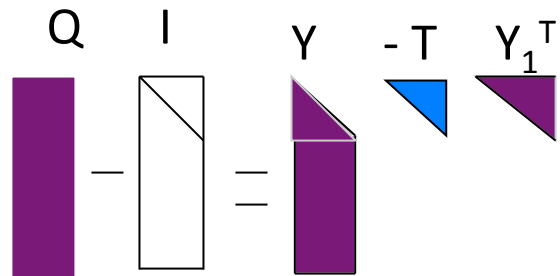
The QR factorization using Householder vectors

$$W = QR = (I - YTY_1^T)R$$

can be re-written as an LU factorization

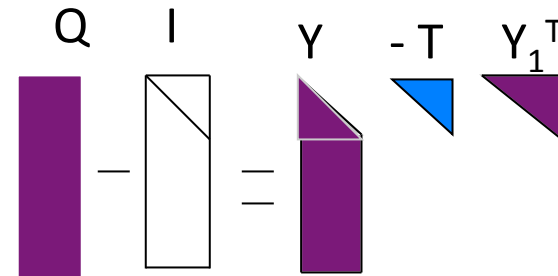
$$W - R = Y(-TY_1^T)R$$

$$Q - I = Y(-TY_1^T)$$



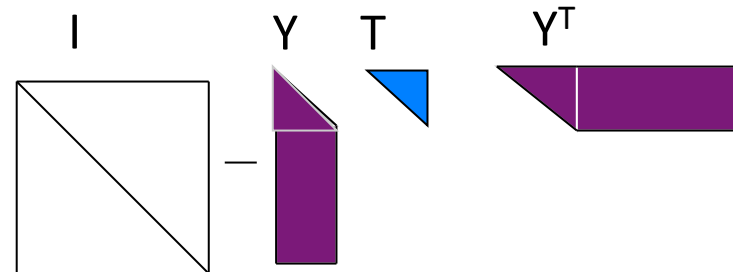
# Reconstruct Householder vectors TSQR-HR

1. Perform TSQR
2. Form Q explicitly (tall-skinny orthonormal factor)
3. Perform LU decomposition:  $Q - I = LU$



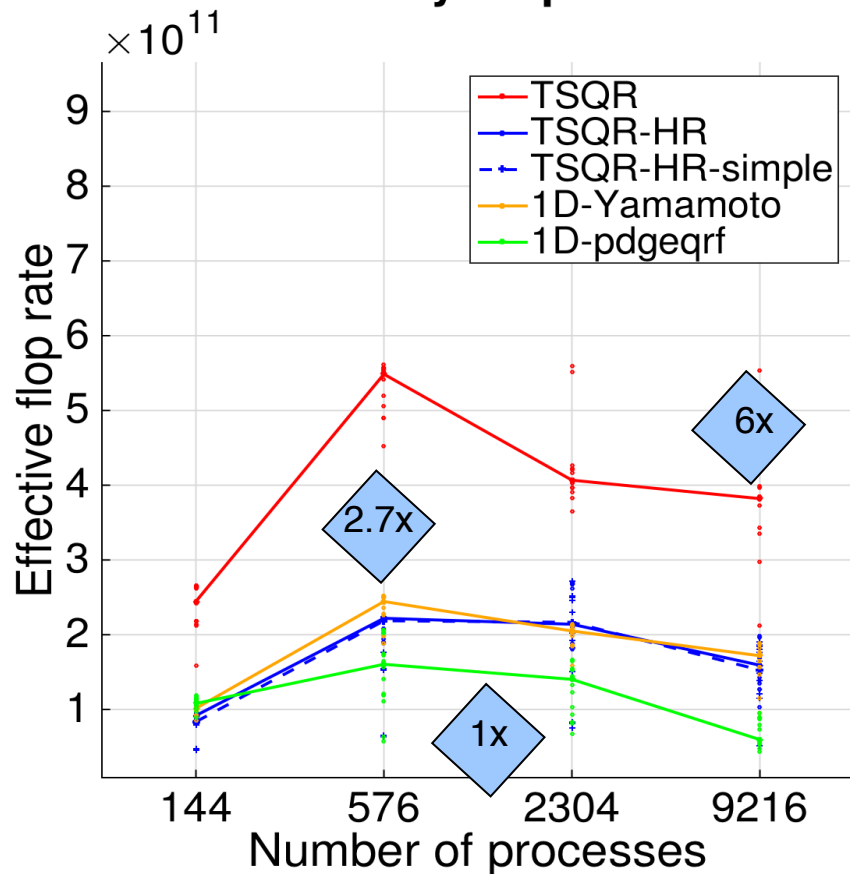
4. Set  $Y = L$
5. Set  $T = -U Y_1^{-T}$

$$I - YTY^T = I - \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} T \begin{bmatrix} Y_1^T & Y_2^T \end{bmatrix}$$

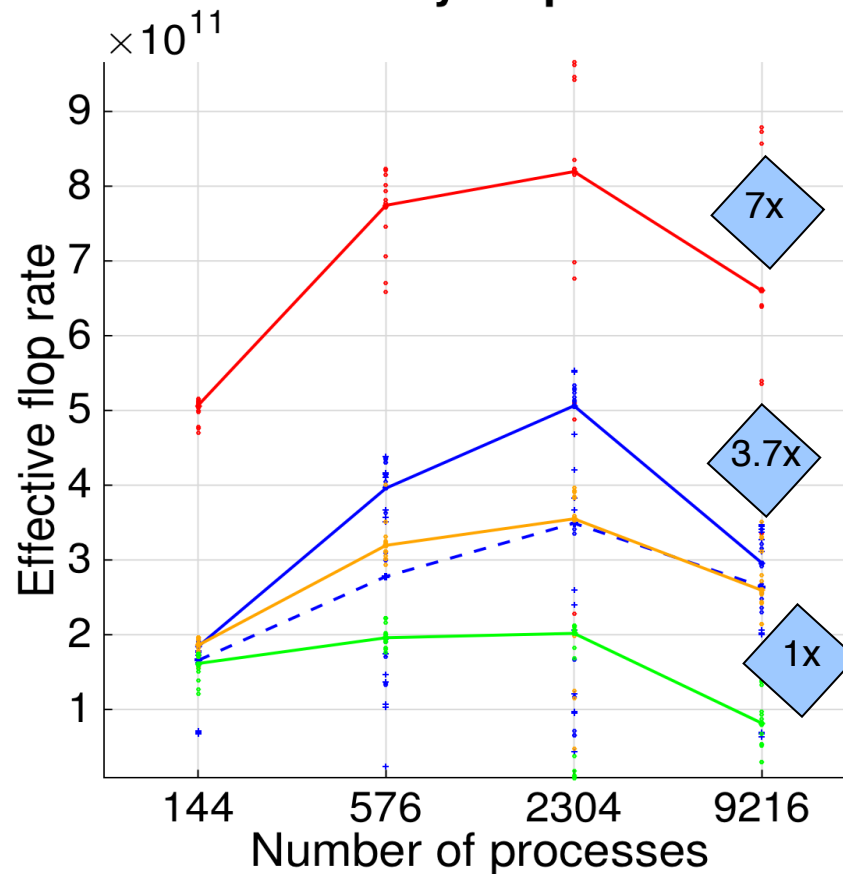


## Strong scaling

**Strong Scaling, Hopper (MKL)  
294912-by-32 problem**



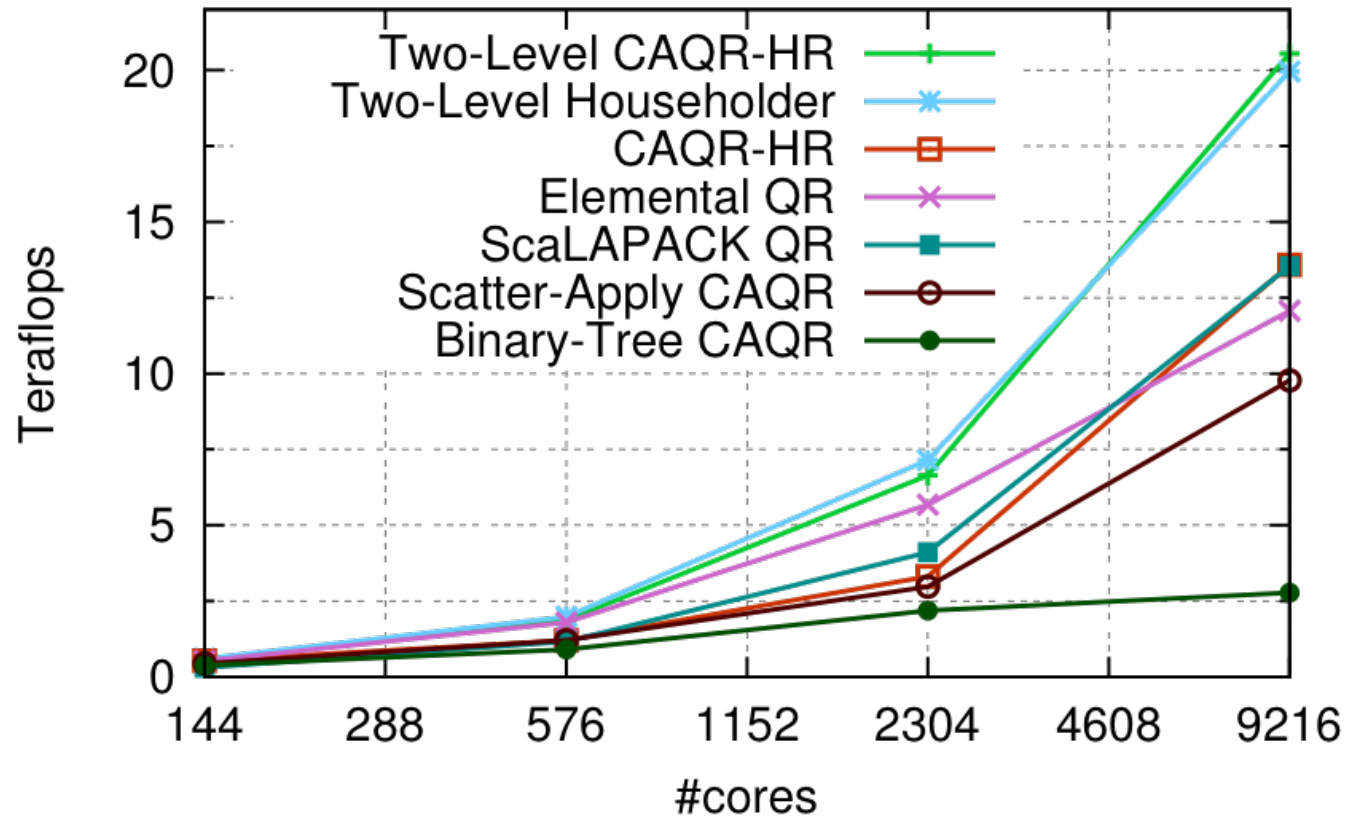
**Strong Scaling, Edison (MKL)  
294912-by-32 problem**



- Hopper: Cray XE6 (NERSC) – 2 x 12-core AMD Magny-Cours (2.1 GHz)
- Edison: Cray CX30 (NERSC) – 2 x 12-core Intel Ivy Bridge (2.4 GHz)
- Effective flop rate, computed by dividing  $2mn^2 - 2n^3/3$  by measured runtime

## Weak scaling QR on Hopper

QR weak scaling on Hopper (15K-by-15K to 131K-by-131K)



- Matrix of size 15K-by-15K to 131K-by-131K
- Hopper: Cray XE6 supercomputer (NERSC) – dual socket 12-core Magny-Cours Opteron (2.1 GHz)

# The LU factorization of a tall skinny matrix

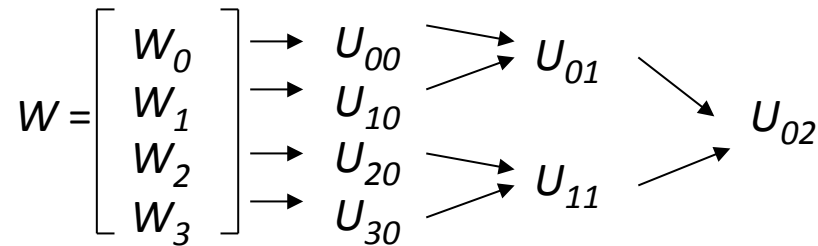
First try the obvious generalization of TSQR.

$$W = \begin{pmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{pmatrix} = \underbrace{\begin{pmatrix} \Pi_{00} & & & \\ & \Pi_{10} & & \\ & & \Pi_{20} & \\ & & & \Pi_{30} \end{pmatrix}}_{\Pi_0} \cdot \begin{pmatrix} L_{00} & & & \\ & L_{10} & & \\ & & L_{20} & \\ & & & L_{30} \end{pmatrix} \cdot \begin{pmatrix} U_{00} \\ U_{10} \\ U_{20} \\ U_{30} \end{pmatrix}$$

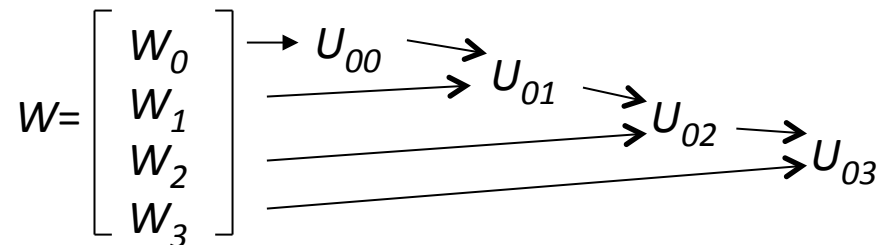
$$\begin{pmatrix} U_{00} \\ U_{10} \\ U_{20} \\ U_{30} \end{pmatrix} = \underbrace{\begin{pmatrix} \Pi_{01} & & \\ & \Pi_{11} & \\ & & \Pi_{21} \\ & & & \Pi_{31} \end{pmatrix}}_{\Pi_1} \cdot \begin{pmatrix} L_{01} & & \\ & L_{11} & \\ & & L_{21} \\ & & & L_{31} \end{pmatrix} \cdot \begin{pmatrix} U_{01} \\ U_{11} \\ U_{21} \\ U_{31} \end{pmatrix} \quad \begin{pmatrix} U_{01} \\ U_{11} \end{pmatrix} = \underbrace{\Pi_{02}}_{\Pi_2} L_{02} U_{02}$$

# Obvious generalization of TSQR to LU

- Block parallel pivoting:
  - uses a binary tree and is optimal in the parallel case



- Block pairwise pivoting:
  - uses a flat tree and is optimal in the sequential case
  - introduced by Barron and Swinnerton-Dyer, 1960: block LU factorization used to solve a system with 100 equations on EDSAC 2 computer using an auxiliary magnetic-tape
  - used in PLASMA for multicore architectures and FLAME for out-of-core algorithms and for multicore architectures



# Stability of the LU factorization

- The backward stability of the LU factorization of a matrix A of size n-by-n

$$\left\| \hat{L} \cdot \hat{U} \right\|_{\infty} \leq (1 + 2(n^2 - n)g_w) \|A\|_{\infty}$$

depends on the growth factor

$$g_w = \frac{\max_{i,j,k} |a_{ij}^k|}{\max_{i,j} |a_{ij}|} \quad \text{where } a_{ij}^k \text{ are the values at the } k\text{-th step.}$$

$$A = \text{diag}(\pm 1) \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 & 1 \\ -1 & 1 & & \cdots & 0 & 1 \\ -1 & -1 & 1 & \ddots & 0 & 1 \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ -1 & -1 & \cdots & -1 & 1 & 1 \\ -1 & -1 & \cdots & -1 & -1 & 1 \end{pmatrix}$$

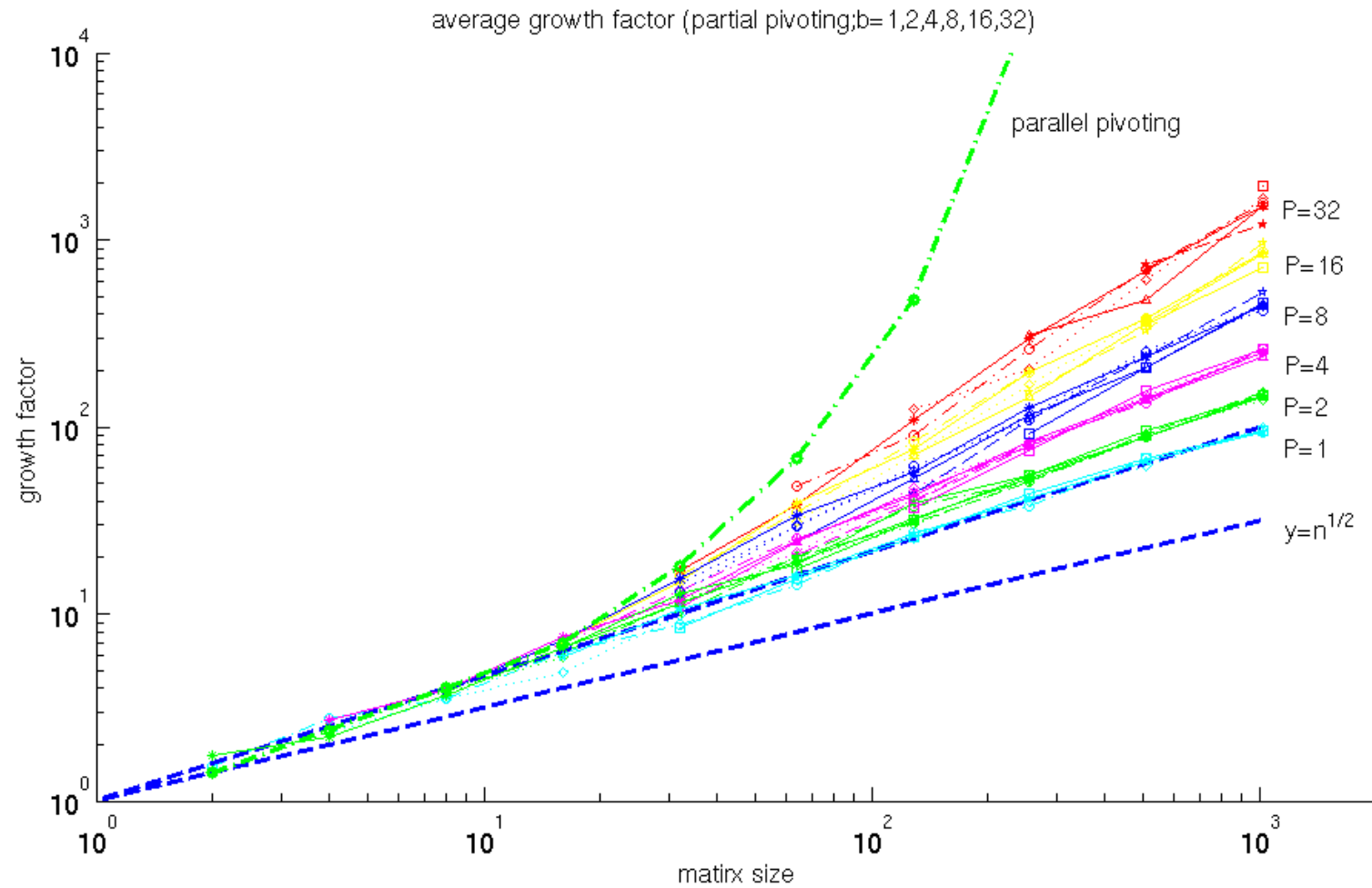
- $g_w \leq 2^{n-1}$ , attained for Wilkinson matrix

but in practice it is on the order of  $n^{2/3} \sim n^{1/2}$

- Two reasons considered to be important for the average case stability [Trefethen and Schreiber, 90]:
  - the multipliers in L are small,
  - the correction introduced at each elimination step is of rank 1.



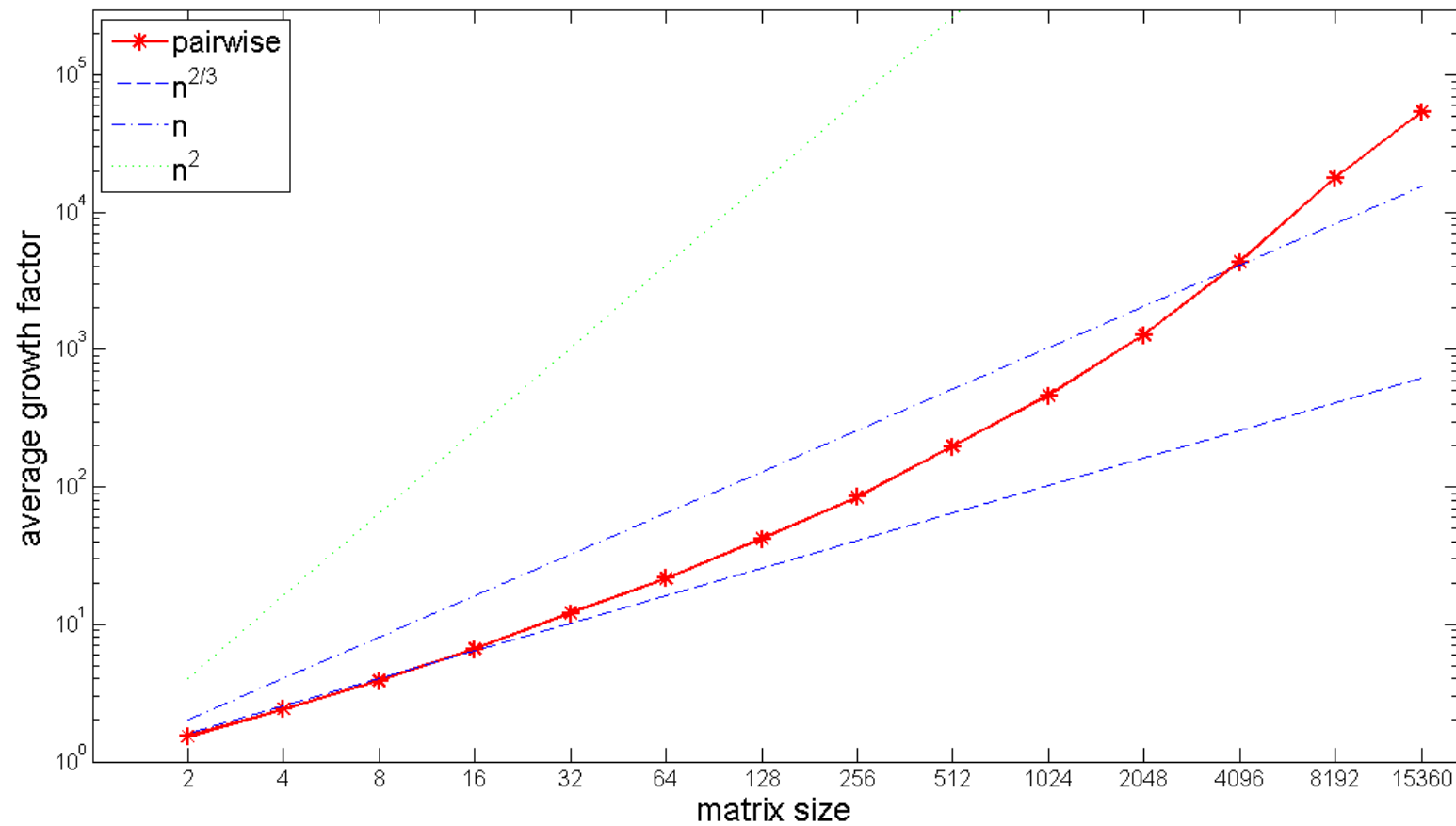
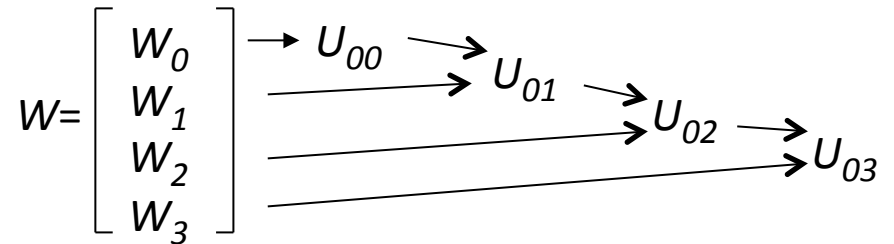
# Block parallel pivoting



- Unstable for large number of processors  $P$
- When  $P$ =number rows, it corresponds to parallel pivoting, known to be unstable (Trefethen and Schreiber, 90)

# Block pairwise pivoting

- Results shown for random matrices
- Will become unstable for large matrices



## Tournament pivoting - the overall idea

- At each iteration of a block algorithm

$$A = \left( \begin{array}{cc} \tilde{A}_{11} & \tilde{A}_{21} \\ A_{21} & A_{22} \end{array} \right) \left. \begin{array}{l} \} b \\ \} n-b \end{array} \right\} , \text{ where } W = \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$$

- Preprocess  $W$  to find at low communication cost good pivots for the LU factorization of  $W$ , return a permutation matrix  $P$ .
- Permute the pivots to top, ie compute  $PA$ .
- Compute LU with no pivoting of  $W$ , update trailing matrix.

$$PA = \begin{pmatrix} L_{11} & \\ L_{21} & I_{n-b} \end{pmatrix} \begin{pmatrix} U_{11} & U_{12} \\ & A_{22} - L_{21}U_{12} \end{pmatrix}$$

# Tournament pivoting for a tall skinny matrix

- 1) Compute GEPP factorization of each  $W_i$ , find permutation  $\Pi_0$

$$W = \begin{pmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{pmatrix} = \begin{pmatrix} \Pi_{00} L_{00} U_{00} \\ \Pi_{10} L_{10} U_{10} \\ \Pi_{20} L_{20} U_{20} \\ \Pi_{30} L_{30} U_{30} \end{pmatrix}, \begin{array}{l} \text{Pick } b \text{ pivot rows, form } A_{00} \\ \text{Same for } A_{10} \\ \text{Same for } A_{20} \\ \text{Same for } A_{30} \end{array}$$

- 2) Perform  $\log_2(P)$  times GEPP factorizations of  $2b$ -by- $b$  rows, find permutations  $\Pi_1, \Pi_2$

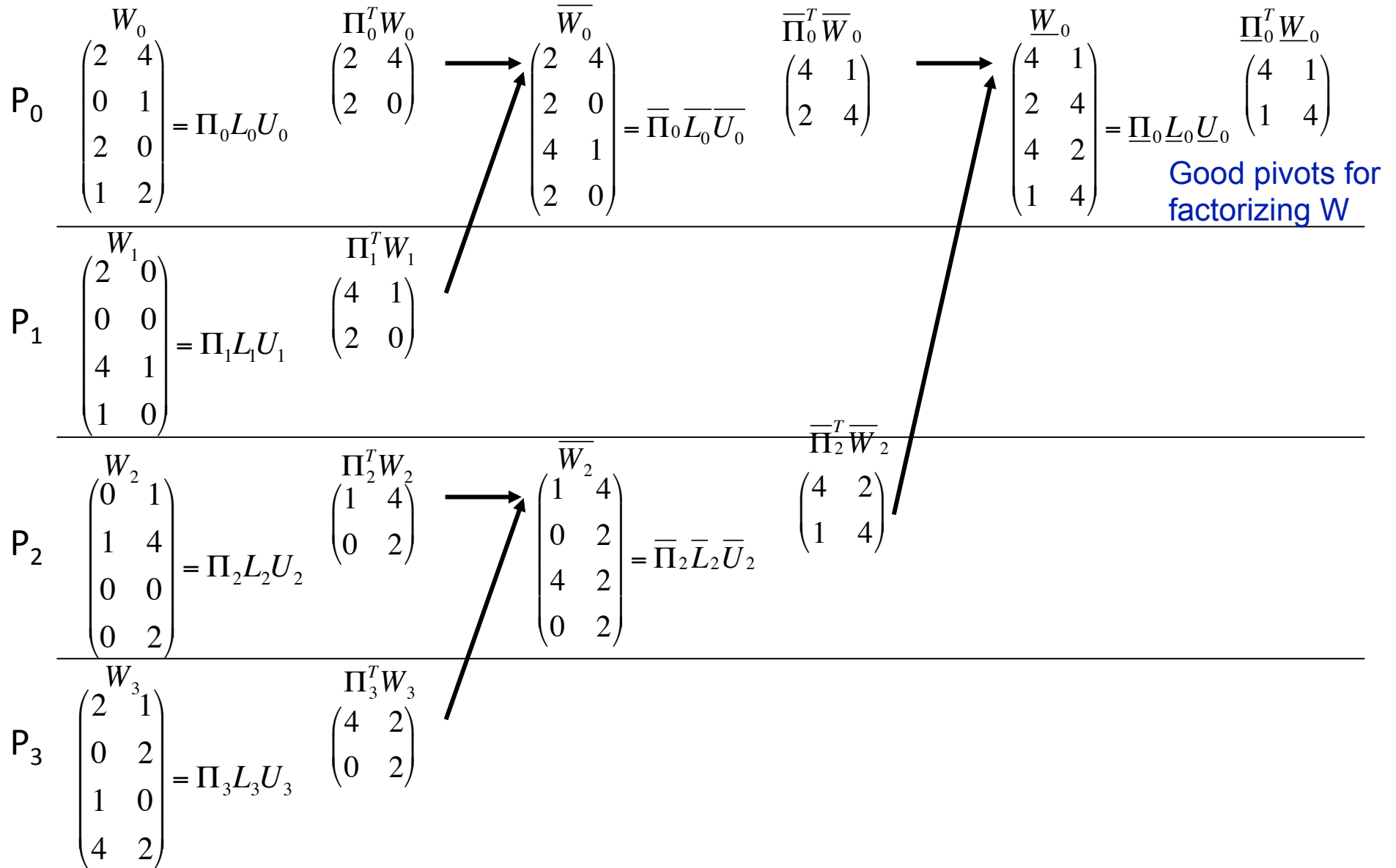
$$\begin{pmatrix} A_{00} \\ A_{10} \\ A_{20} \\ A_{30} \end{pmatrix} = \begin{pmatrix} \Pi_{01} L_{01} U_{01} \\ \Pi_{11} L_{11} U_{11} \end{pmatrix}, \begin{array}{l} \text{Pick } b \text{ pivot rows, form } A_{01} \\ \text{Same for } A_{11} \end{array}$$

$$\begin{pmatrix} A_{01} \\ A_{11} \end{pmatrix} = \underbrace{\Pi_{02}}_{\Pi_2} L_{02} U_{02}$$

- 3) Compute LU factorization with no pivoting of the permuted matrix:

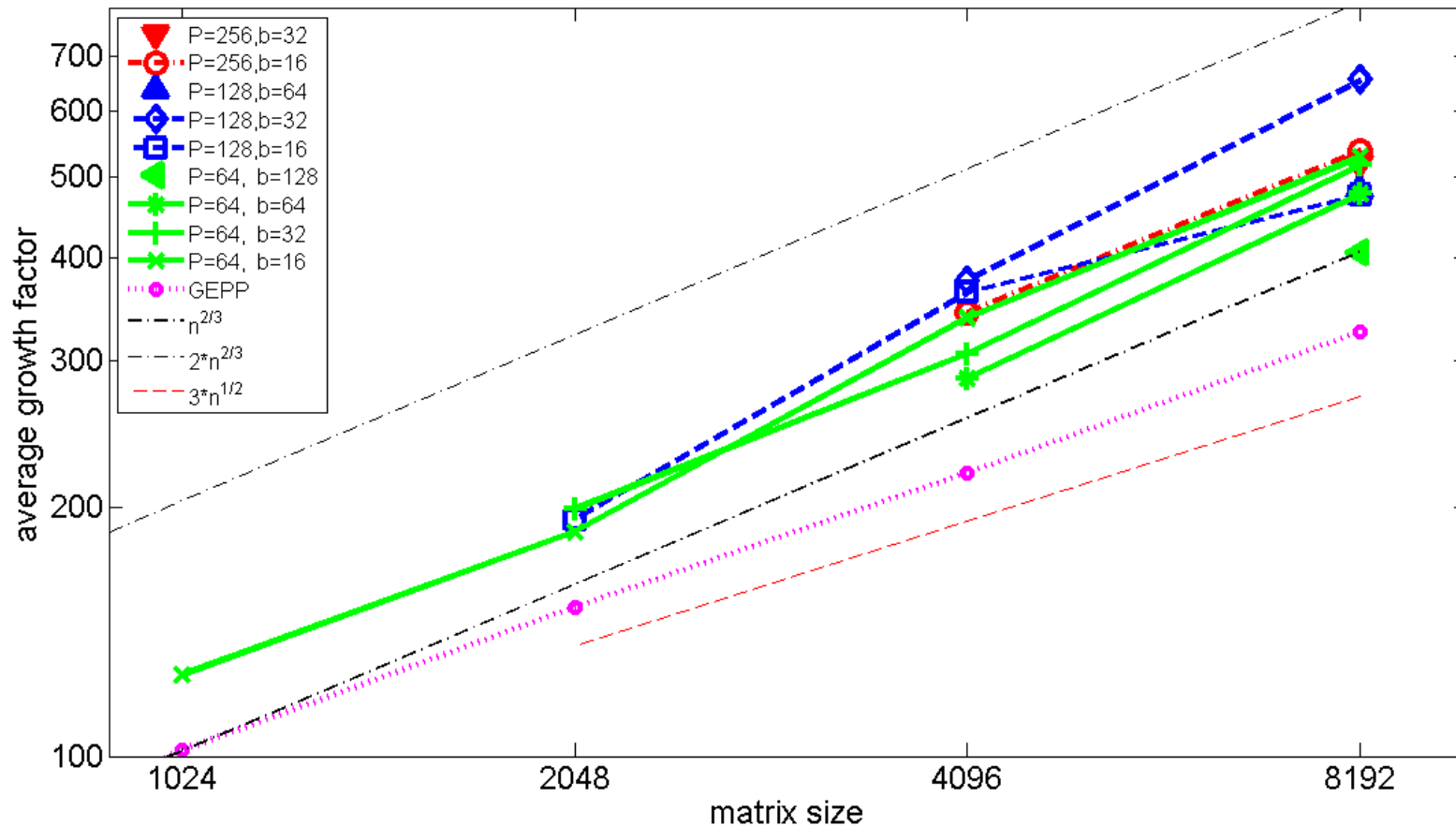
$$\Pi_2^T \Pi_1^T \Pi_0^T W = LU$$

# Tournament pivoting



time  $\longrightarrow$

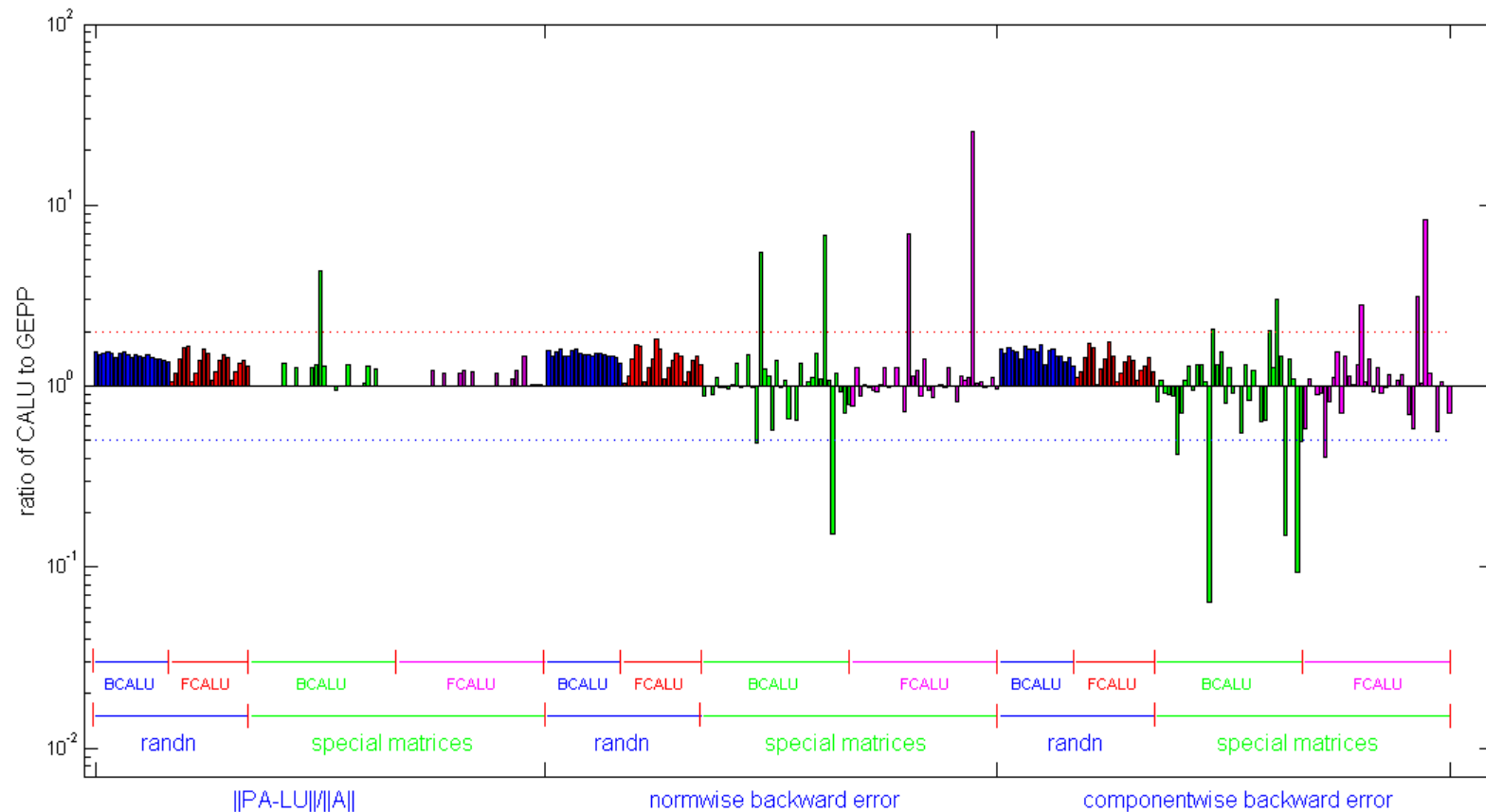
# Growth factor for binary tree based CALU



- Random matrices from a normal distribution
- Same behaviour for all matrices in our test, and  $|L| \leq 4.2$

# Stability of CALU (experimental results)

- Results show  $\|PA-LU\|/\|A\|$ , normwise and componentwise backward errors, for random matrices and special ones
  - See [LG, Demmel, Xiang, SIMAX 2011] for details
  - BCALU denotes binary tree based CALU and FCALU denotes flat tree based CALU

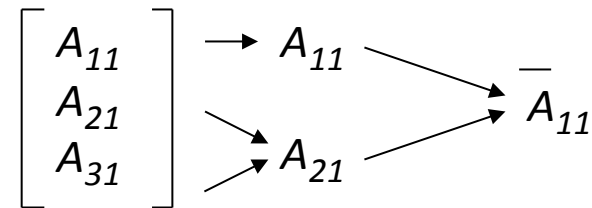


## Our “proof of stability” for CALU

- CALU as stable as GEPP in following sense:  
 In exact arithmetic, CALU process on a matrix  $A$  is equivalent to GEPP process on a larger matrix  $G$  whose entries are blocks of  $A$  and zeros.
- Example of one step of tournament pivoting:

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{32} \end{pmatrix}$$

tournament pivoting:



$$G = \begin{pmatrix} \bar{A}_{11} & & \bar{A}_{12} \\ A_{21} & A_{21} & \\ & -A_{31} & A_{32} \end{pmatrix}$$

- Proof possible by using original rows of  $A$  during tournament pivoting (not the computed rows of  $U$ ).



# LU factorization and low rank matrices

- For low rank matrices, the factorization of  $A_1$  computed as following might not be stable

Compute  $PA=LU$  by using GEPP

$$L(k+1:end,k) = A(k+1:end,k)/A(k,k)$$

Permute the matrix  $A_1=PA$

Compute LU with no pivoting  $A_1=L_1U_1$

$$L(k+1:end,k) = L(k+1:end,k)^* (1/A(k,k))$$

- Example  $A = \text{randn}(6,3)*\text{randn}(3,5)$ ,  $\max(\text{abs}(L)) = 1$ ,  $\max(\text{abs}(L_1)) = 10^{15}$

After 4 steps of factorization of  $PA$  we obtain :

$$PA^4 = \begin{pmatrix} 1.0000 & & & & & \\ 0.1729 & 1.0000 & & & & \\ 0.6061 & 0.8608 & 1.0000 & & & \\ 0.5776 & 0.0543 & 0.3264 & 1.0000 & & \\ 0.4789 & -0.2877 & -0.1545 & 2.3333 & 2.3e-16 & \\ -0.3264 & -0.7514 & -0.4597 & 1.7778 & 8.3e-17 & \end{pmatrix} \cdot \begin{pmatrix} 4.4766 & 3.0163 & -4.7390 & 4.2180 & -0.8164 & \\ & -1.5439 & -0.4703 & 1.9267 & 1.0925 & \\ & & 1.6149 & 2.3623 & 0.3167 & \\ & & & 9.9e-16 & 1.6e-16 & \\ & & & & 1 & \end{pmatrix}$$

After 4 steps of factorization of  $A_1$  we obtain :

$$A_1^4 = \begin{pmatrix} 1.0000 & & & & & \\ 0.1729 & 1.0000 & & & & \\ 0.6061 & 0.8608 & 1.0000 & & & \\ 0.5776 & 0.0543 & 0.3264 & 1.0000 & & \\ 0.4789 & -0.2877 & -0.1545 & 2.3333 & 4.9e-32 & \\ -0.3264 & -0.7514 & -0.4597 & 1.7778 & -7.4e-17 & \end{pmatrix} \cdot \begin{pmatrix} 4.4766 & 3.0163 & -4.7390 & 4.2180 & -0.8164 & \\ & -1.5439 & -0.4703 & 1.9267 & 1.0925 & \\ & & 1.6149 & 2.3623 & 0.3167 & \\ & & & 9.9e-16 & 1.6e-16 & \\ & & & & 1 & \end{pmatrix}$$

Schur complement after 4 elimination steps

## LU\_PRRP: LU with panel rank revealing pivoting

- Pivots are selected by using strong rank revealing QR on each panel
- The factorization after one panel elimination is written as

$$PA = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} I_b & \\ A_{21}A_{11}^{-1} & I_{n-b} \end{pmatrix} \begin{pmatrix} A_{11} & A_{12} \\ & A_{22} - A_{21}A_{11}^{-1}A_{12} \end{pmatrix}$$

$A_{21} A_{11}^{-1}$  is computed through strong rank revealing QR  
and  $\max(|A_{21} A_{11}^{-1}|)_{ij} \leq f$

- LU\_PRRP and CALU\_PRRP stable for pathological cases (Wilkinson matrix) and matrices from two real applications (Volterra integral equation - Foster, a boundary value problem - Wright) on which GEPP fails.

## Growth factor in exact arithmetic

- Matrix of size m-by-n, reduction tree of height  $H=\log(P)$ .
- (CA)LU\_PRRP select pivots using strong rank revealing QR (A. Khabou, J. Demmel, LG, M. Gu, SIMAX 2013)
- “In practice” means observed/expected/conjectured values.

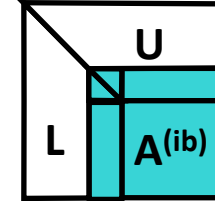
	CALU	GEPP	CALU_PRRP	LU_PRRP
Upper bound	$2^{n(\log(P)+1)-1}$	$2^{n-1}$	$(1+2b)^{(n/b)\log(P)}$	$(1+2b)^{(n/b)}$
In practice	$n^{2/3} \text{ -- } n^{1/2}$	$n^{2/3} \text{ -- } n^{1/2}$	$(n/b)^{2/3} \text{ -- } (n/b)^{1/2}$	$(n/b)^{2/3} \text{ -- } (n/b)^{1/2}$


  
**Better bounds**

- For a matrix of size  $10^7$ -by- $10^7$  (using petabytes of memory)
  - $n^{1/2} = 10^{3.5}$
- When will Linpack have to use the QR factorization for solving linear systems ?

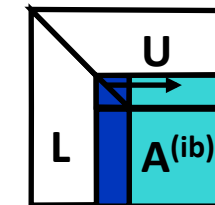
# CALU – a communication avoiding LU factorization

- Consider a 2D grid of  $P$  processors  $P_r$ -by- $P_c$ , using a 2D block cyclic layout with square blocks of size  $b$ .

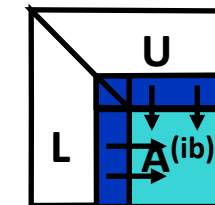


For  $ib = 1$  to  $n-1$  step  $b$   
 $A^{(ib)} = A(ib:n, ib:n)$

(1) Find permutation for current panel using TSLU  $O(n/b \log_2 P_r)$



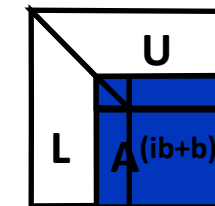
(2) Apply all row permutations (**pdlaswp**)  $O(n/b(\log_2 P_c + \log_2 P_r))$   
 - broadcast pivot information along the rows of the grid



(3) Compute panel factorization (**dtrsm**)

$$O(n/b \log_2 P_c)$$

(4) Compute block row of  $U$  (**pdtrsm**)  
 - broadcast right diagonal part of  $L$  of current panel



(5) Update trailing matrix (**pdgemm**)  
 - broadcast right block column of  $L$   
 - broadcast down block row of  $U$

$$O(n/b(\log_2 P_c + \log_2 P_r))$$

# LU for General Matrices

- Cost of **CALU** vs **ScaLAPACK's PDGETRF**
  - $n \times n$  matrix on  $P^{1/2} \times P^{1/2}$  processor grid, block size  $b$
  - Flops:  $(2/3)n^3/P + (3/2)n^2b / P^{1/2}$  vs  $(2/3)n^3/P + n^2b/P^{1/2}$
  - Bandwidth:  $n^2 \log P/P^{1/2}$  vs same
  - Latency:  $3 n \log P / b$  vs  $1.5 n \log P + 3.5n \log P / b$
- Close to optimal (modulo  $\log P$  factors)
  - Assume:  $O(n^2/P)$  memory/processor,  $O(n^3)$  algorithm,
  - Choose  $b$  near  $n / P^{1/2}$  (its upper bound)
  - Bandwidth lower bound:
    - $\Omega(n^2 / P^{1/2})$  – just  $\log(P)$  smaller
  - Latency lower bound:
    - $\Omega(P^{1/2})$  – just  $\text{polylog}(P)$  smaller

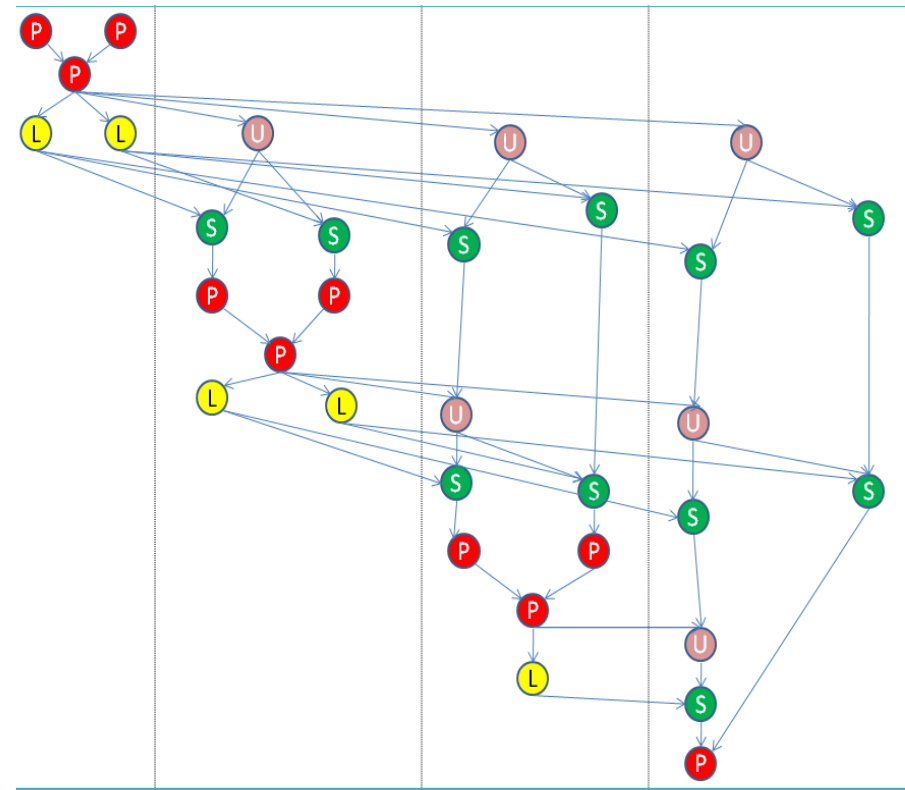


# Performance vs ScaLAPACK

- Parallel TSLU (LU on tall-skinny matrix)
  - IBM Power 5
    - Up to **4.37x** faster (16 procs, 1M x 150)
  - Cray XT4
    - Up to **5.52x** faster (8 procs, 1M x 150)
- Parallel CALU (LU on general matrices)
  - Intel Xeon (two socket, quad core)
    - Up to **2.3x** faster (8 cores,  $10^6$  x 500)
  - IBM Power 5
    - Up to **2.29x** faster (64 procs, 1000 x 1000)
  - Cray XT4
    - Up to **1.81x** faster (64 procs, 1000 x 1000)
- Details in SC08 (LG, Demmel, Xiang), IPDPS'10 (S. Donfack, LG).

# CALU and its task dependency graph

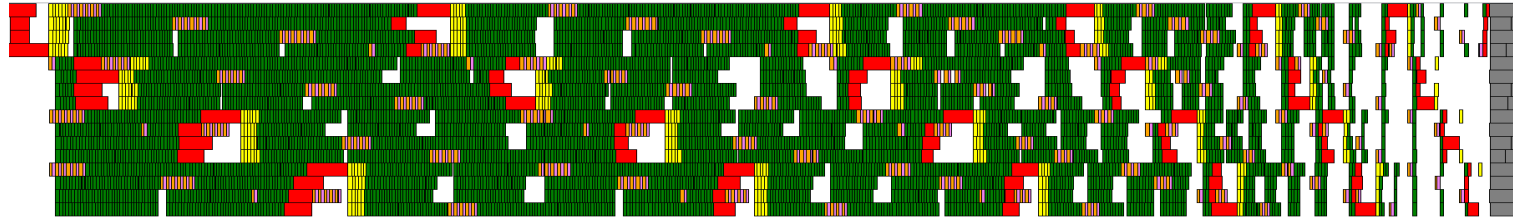
- The matrix is partitioned into blocks of size  $T \times b$ .
- The computation of each block is associated with a task.



# Scheduling CALU's Task Dependency Graph

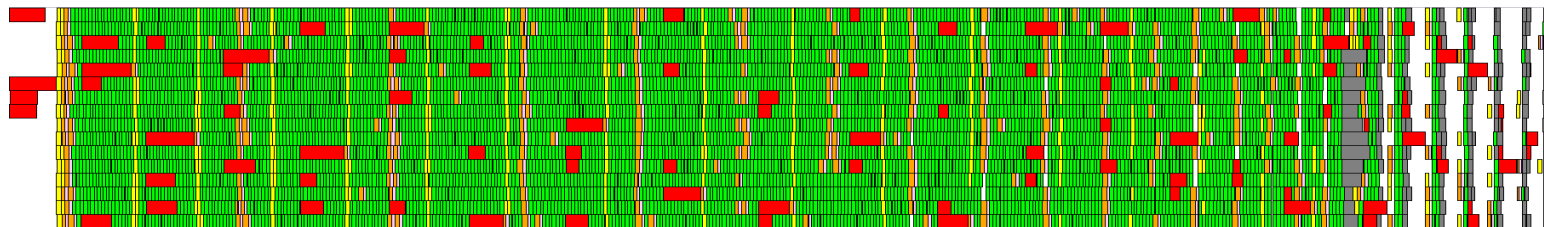
- Static scheduling

- + Good locality of data
- Ignores noise



- Dynamic scheduling

- + Keeps cores busy
- Poor usage of data locality
- Can have large dequeue overhead





# Lightweight scheduling

- Emerging complexities of multi- and mani-core processors suggest a need for self-adaptive strategies
  - One example is work stealing
- Goal:
  - Design a tunable strategy that is able to provide a good trade-off between load balance, data locality, and dequeue overhead.
  - Provide performance consistency
- Approach: combine static and dynamic scheduling
  - Shown to be efficient for regular mesh computation [B. Gropp and V. Kale]

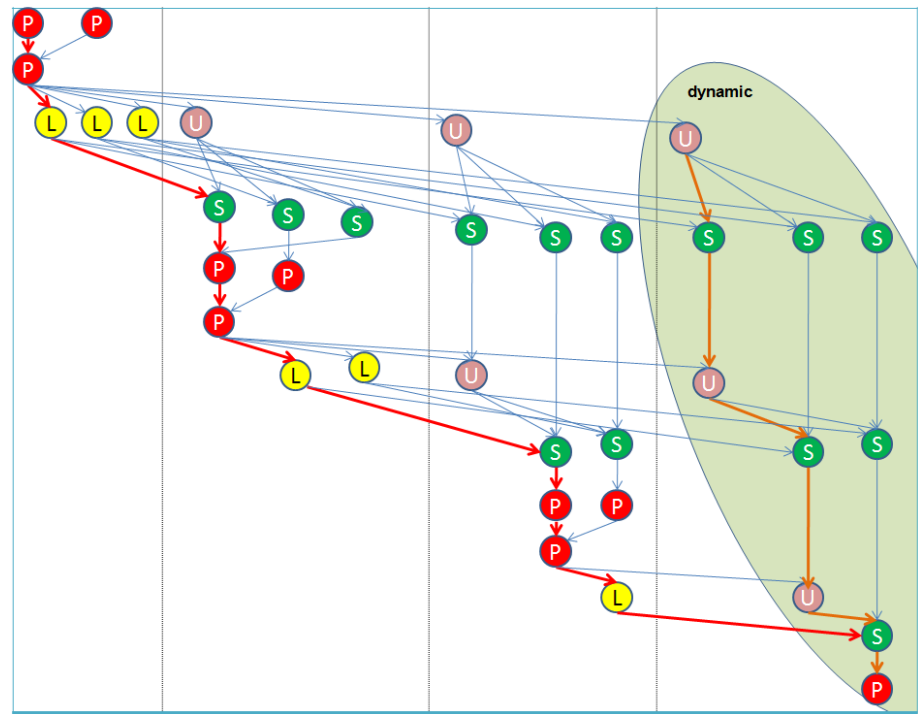
Design space			
Data layout/scheduling	Static	Dynamic	Static/(%dynamic)
Column Major Layout (CM)		√	
Block Cyclic Layout (BCL)	√	√	√
2-level Block Layout (2l-BL)	√	√	√

# Lightweight scheduling

- A self-adaptive strategy to provide
  - A good trade-off between load balance, data locality, and dequeue overhead.
  - Performance consistency
  - Shown to be efficient for regular mesh computation [B. Gropp and V. Kale]

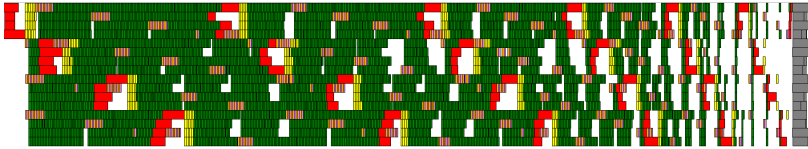
Combined static/dynamic scheduling:

- A thread executes in priority its statically assigned tasks
- When no task ready, it picks a ready task from the dynamic part
- The size of the dynamic part is guided by a performance model

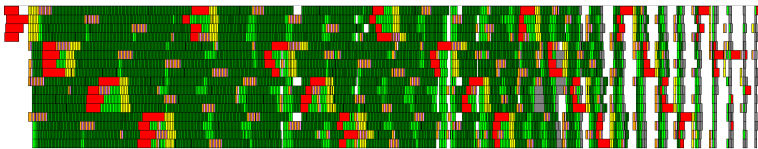


# Best performance of CALU on multicore architectures

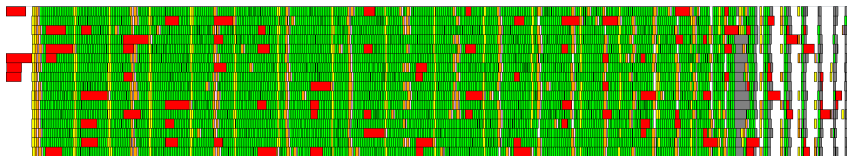
Static scheduling



Static + 10% dynamic scheduling

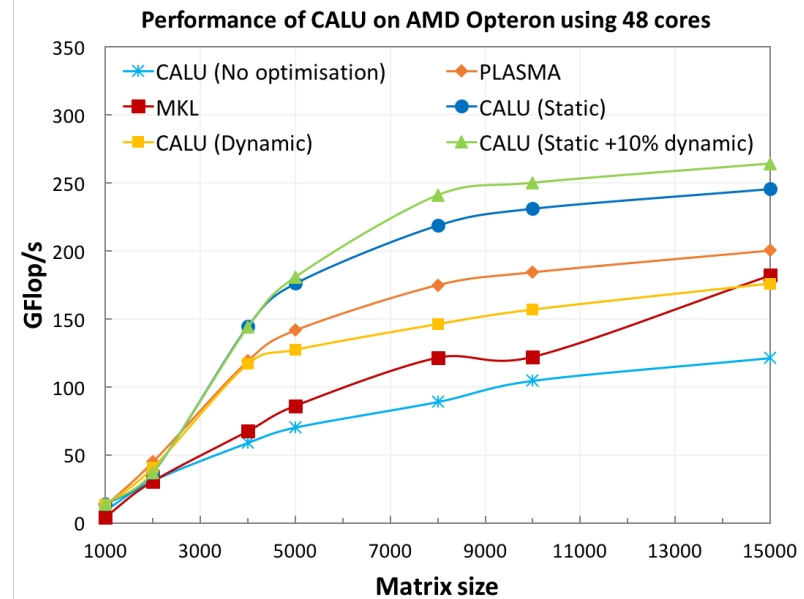
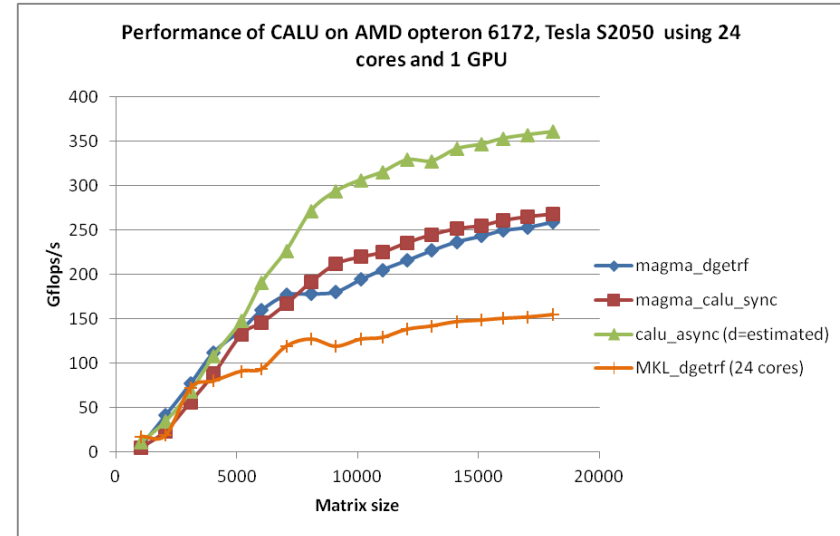


100% dynamic scheduling



time →

- Reported performance for PLASMA uses LU with block pairwise pivoting.
- GPU data courtesy of S. Donfack



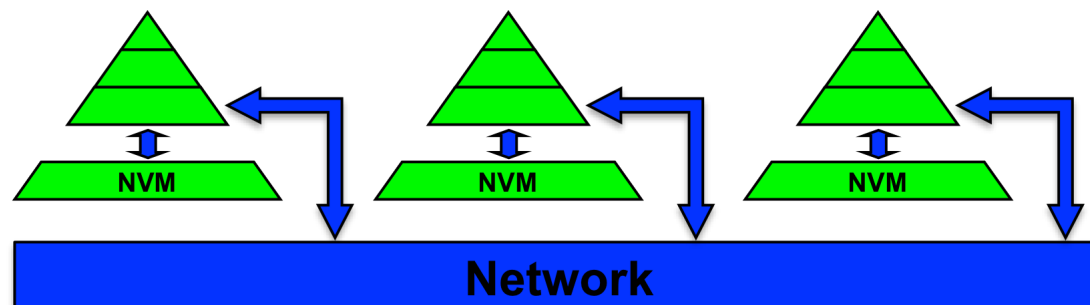
## Parallel write avoiding algorithms

Need to avoid writing suggested by emerging memory technologies, as NVMs:

- Writes more expensive (in time and energy) than reads
- Writes are less reliable than reads

Some examples:

- Phase Change Memory: Reads 25 us latency  
Writes: 15x slower than reads (latency and bandwidth)  
consume 10x more energy
- Conductive Bridging RAM - CBRAM  
Writes: use more energy (1pJ) than reads (50 fJ)
- Gap improving by new technologies such as XPoint and other FLASH alternatives, but not eliminated



## Parallel write-avoiding algorithms

- Matrix A does not fit in DRAM (of size M), need to use NVM (of size  $n^2 / P$ )
- Two lower bounds on volume of communication
  - Interprocessor communication:  $\Omega (n^2 / P^{1/2})$
  - Writes to NVM:  $n^2 / P$
- Result: any three-nested loop algorithm (matrix multiplication, LU,..), must asymptotically exceed at least one of these lower bounds
  - If  $\Omega (n^2 / P^{1/2})$  words are transferred over the network, then  $\Omega (n^2 / P^{2/3})$  words must be written to NVM !
- Parallel LU: choice of best algorithm depends on hardware parameters

	#words interprocessor comm.	#writes NVM
Left-looking	$O((n^3 \log^2 P) / (P M^{1/2}))$	$O(n^2 / P)$
Right-looking	$O((n^2 \log P) / P^{1/2})$	$O((n^2 \log^2 P) / P^{1/2})$

# Conclusions

- Many previous results
  - Only several cited, many references given in the papers
  - Flat trees algorithms for QR factorization, called tiled algorithms used in the context of
    - Out of core - Gunter, van de Geijn 2005
    - Multicore, Cell processors - Buttari, Langou, Kurzak and Dongarra (2007, 2008), Quintana-Orti, Quintana-Orti, Chan, van Zee, van de Geijn (2007, 2008)

# References

Results presented from:

- J. Demmel, L. Grigori, M. F. Hoemmen, and J. Langou, *Communication-optimal parallel and sequential QR and LU factorizations*, UCB-EECS-2008-89, 2008, SIAM journal on Scientific Computing, Vol. 34, No 1, 2012.
- L. Grigori, J. Demmel, and H. Xiang, *Communication avoiding Gaussian elimination*, Proceedings of the IEEE/ACM SuperComputing SC08 Conference, November 2008.
- L. Grigori, J. Demmel, and H. Xiang, *CALU: a communication optimal LU factorization algorithm*, SIAM. J. Matrix Anal. & Appl., 32, pp. 1317-1350, 2011.
- M. Hoemmen's Phd thesis, *Communication avoiding Krylov subspace methods*, 2010.
- L. Grigori, P.-Y. David, J. Demmel, and S. Peyronnet, *Brief announcement: Lower bounds on communication for sparse Cholesky factorization of a model problem*, ACM SPAA 2010.
- S. Donfack, L. Grigori, and A. Kumar Gupta, *Adapting communication-avoiding LU and QR factorizations to multicore architectures*, Proceedings of IEEE International Parallel & Distributed Processing Symposium IPDPS, April 2010.
- S. Donfack, L. Grigori, W. Gropp, and V. Kale, *Hybrid static/dynamic scheduling for already optimized dense matrix factorization*, Proceedings of IEEE International Parallel & Distributed Processing Symposium IPDPS, 2012.
- A. Khabou, J. Demmel, L. Grigori, and M. Gu, *LU factorization with panel rank revealing pivoting and its communication avoiding version*, LAWN 263, SIAM Journal on Matrix Analysis, in revision, 2012.
- L. Grigori, S. Moufawad, *Communication avoiding ILU0 preconditioner*, Inria TR 8266, 2013.
- J. Demmel, L. Grigori, M. Gu, H. Xiang, *Communication avoiding rank revealing QR factorization with column pivoting*, 2013.