# Communication avoiding algorithms for LU and QR factorizations

# Laura Grigori *Alpines*

INRIA Paris - LJLL, Sorbonne Université

November 2021

# Plan

- Motivation
- Communication complexity of linear algebra operations
- Communication avoiding for dense linear algebra
  - LU, QR, Rank Revealing QR factorizations
  - Progressively implemented in ScaLAPACK, LAPACK
  - Algorithms for multicore processors
- Conclusions

## 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	Minimizing	
	#words (not #messages)	#words and #messages	
Cholesky	ScaLAPACK	ScaLAPACK	
LU	L ScaLAPACK es 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	Q A <sup>(ib)</sup> 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

#### LU factorization on a $P = P_r \times P_c$ grid of processors For ib = 1 to n-1 step b

LU factorization (as in ScaLAPACK pdgetrf)

#messages

 $O(n\log_2 P_r)$ 

- (1) Compute panel factorization
  - find pivot in each column, swap rows
- (2) Apply all row permutations

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

- broadcast pivot information along the rows
- swap rows at left and right
- (3) Compute block row of U
  - broadcast right diagonal block of L of current panel
- (4) Update trailing matrix
  - broadcast right block column of L
  - broadcast down block row of U





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



A(ib)





 $O(n/b\log_{P_c})$ 

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

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

$$W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix} \xrightarrow{\rightarrow} \begin{bmatrix} R_{00} \\ R_{10} \\ R_{20} \\ R_{30} \end{bmatrix} \xrightarrow{\rightarrow} R_{01} \xrightarrow{} R_{02}$$

J. Demmel, LG, M. Hoemmen, J. Langou, 08

## Parallel TSQR



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

## Algebra of TSQR



Q is represented implicitly as a product Output: { $Q_{00}$ ,  $Q_{10}$ ,  $Q_{00}$ ,  $Q_{20}$ ,  $Q_{30}$ ,  $Q_{01}$ ,  $Q_{11}$ ,  $Q_{02}$ ,  $R_{02}$ }

Page 8

## Flexibility of TSQR and CAQR algorithms

Parallel: 
$$W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix} \xrightarrow{\rightarrow} R_{20} \xrightarrow{R_{01}} R_{11} \xrightarrow{R_{02}} R_{20}$$





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

## Algebra of TSQR





# **QR** for General Matrices

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

 $\Omega(n^2 / P^{1/2}) - just log(P) smaller$ 

• Latency lower bound:

 $\Omega(P^{1/2})$  – just polylog(P) smaller



VS

same

## Performance of TSQR vs Sca/LAPACK

- Parallel
  - Intel Xeon (two socket, quad core machine), 2010
    - Up to **5.3x speedup** (8 cores, 10<sup>5</sup> x 200)
  - Pentium III cluster, Dolphin Interconnect, MPICH, 2008
    - Up to 6.7x speedup (16 procs, 100K x 200)
  - BlueGene/L, 2008
    - Up to **4x speedup** (32 procs, 1M x 50)
  - Tesla C 2050 / Fermi (Anderson et al)
    - Up to **13x** (110,592 x 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.  $\gamma = 2.10^{12}$  s,  $\alpha = 10^5$  s,  $\beta = 2.10^9$  s/word

Page 13

## Impact

### TSQR/CAQR implemented in

- Intel Data analytics library
- GNU Scientific Library
- ScaLAPACK
- Spark for data mining

- CALU implemented in
  - Cray's libsci
  - To be implemented in lapack/scapalack

## 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

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

$$\begin{array}{c|cccc} Q & I & Y & -T & Y_1^T \\ \hline - & \hline & - & \hline \\ \hline & - & \hline & - & \hline & - & \hline \\ \hline & - & \hline & - & \hline \\ \hline \hline \\ \hline \hline \\ \hline \hline \hline \hline \end{array} \end{array}$$

## **Reconstruct Householder vectors TSQR-HR**

- **Perform TSQR** 1.
- Form Q explicitly (tall-skinny orthonormal factor) 2.
- Perform LU decomposition: Q I = LU3.

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

$$I - YTY^{T} = I - \begin{bmatrix} Y_{1} \\ Y_{2} \end{bmatrix} T \begin{bmatrix} Y_{1}^{T} & Y_{2}^{T} \end{bmatrix}$$



Q

-T  $Y_1^T$ 

#### Strong scaling



- 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<sup>2</sup> 2n<sup>3</sup>/3 by measured runtime

Ballard, Demmel, LG, Jacquelin, Knight, Nguyen, and Solomonik, 2015. Page 18

## The LU factorization of a tall skinny matrix

First try the obvious generalization of TSQR.



$$\begin{pmatrix} U_{00} \\ U_{10} \\ U_{20} \\ U_{20} \\ U_{30} \end{pmatrix} = \begin{pmatrix} \prod_{01} & L_{01} \\ \prod_{11} \end{pmatrix} \begin{pmatrix} L_{01} & L_{11} \end{pmatrix} \begin{pmatrix} U_{01} \\ U_{11} \end{pmatrix}$$

Page 19

## Obvious generalization of TSQR to LU

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

$$W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix} \xrightarrow{\rightarrow} U_{00} \xrightarrow{\rightarrow} U_{01} \\ \xrightarrow{\rightarrow} U_{10} \xrightarrow{\rightarrow} U_{02} \\ \xrightarrow{\rightarrow} U_{20} \xrightarrow{\rightarrow} U_{11} \xrightarrow{\rightarrow} U_{02}$$

- 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

$$W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix} \xrightarrow{\rightarrow} U_{01} \xrightarrow{\rightarrow} U_{02} \xrightarrow{\rightarrow} U_{03}$$

## Stability of the LU factorization

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

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

depends on the growth factor

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

- $g_W \le 2^{n-1}$ , attained for Wilkinson matrix but in practice it is on the order of  $n^{2/3} - n^{1/2}$
- 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)

Page 22

## Block pairwise pivoting



Page 23

## Tournament pivoting - the overall idea

• At each iteration of a block algorithm

$$\mathcal{A} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} \begin{cases} b \\ n-b \end{cases}, \text{ where } \quad W = \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$$

- 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.

$$PA = \begin{pmatrix} L_{11} & \\ L_{21} & I_{n-b} \end{pmatrix} \begin{pmatrix} U_{11} & U_{12} \\ & A_{22} - L_2 U_{12} \end{pmatrix}$$

## Tournament pivoting for a tall skinny matrix

1) Compute GEPP factorization of each  $W_{i.}$ , find permutation  $\Pi_0$ 

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

2) Perform  $\log_2(P)$  times GEPP factorizations of 2b-by-b rows, find permutations  $\prod_{1}, \prod_{2}$ 

$$\begin{pmatrix} A_{00} \\ A_{10} \\ \hline A_{20} \\ \hline A_{20} \\ A_{30} \end{pmatrix} = \begin{pmatrix} \prod_{01} L_0 H_{01} \\ \hline \prod_{11} L_1 H_{11} \end{pmatrix}$$
 Pick b pivot rows, form A<sub>01</sub>  
Same for A<sub>11</sub>

3) Compute LU factorization with no pivoting of the permuted matrix:  $\Pi_2^T \Pi_1^T \Pi_0^T W = LU$ 

## **Tournament pivoting**



## Growth factor for binary tree based CALU



- Random matrices from a normal distribution
- Same behaviour for all matrices in our test, and |L| <= 4.2

Page 27

## 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:



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

Page 28

## Outline of the proof of stability for CALU

• Consider 
$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{32} \end{pmatrix}$$
, and the result of TSLU as  $\begin{bmatrix} A_{11} \\ A_{21} \\ A_{31} \end{bmatrix} \xrightarrow{\bullet} A_{11} \xrightarrow{\bullet} A_{11}$ 

• After the factorization of first panel by CALU,  $A_{32}^s$  (the Schur complement of  $A_{32}$ ) is not bounded as in GEPP,

$$\begin{pmatrix} \Pi_{11} & \Pi_{12} \\ \Pi_{21} & \Pi_{22} \\ & & I \end{pmatrix} \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \\ A_{31} & A_{32} \end{pmatrix} = \begin{pmatrix} \overline{A}_{11} & \overline{A}_{12} \\ \overline{A}_{21} & \overline{A}_{22} \\ A_{31} & A_{32} \end{pmatrix} = \begin{pmatrix} \overline{L}_{11} & & \\ \overline{L}_{21} & I \\ \overline{L}_{31} & & I \end{pmatrix} \begin{pmatrix} \overline{U}_{11} & \overline{U}_{12} \\ & A_{22}^{s} \\ & A_{32}^{s} \end{pmatrix}$$

• but A<sup>s</sup><sub>32</sub> can be obtained by GEPP on larger matrix G formed from blocks of A

$$G = \begin{pmatrix} \overline{A}_{11} & \overline{A}_{12} \\ A_{21} & A_{21} \\ & -A_{31} & A_{32} \end{pmatrix} = \begin{pmatrix} \overline{L}_{11} & & \\ A_{21}\overline{U}_{11}^{-1} & L_{21} \\ & & -L_{31} & I \end{pmatrix} \begin{pmatrix} \overline{U}_{11} & & \overline{U}_{12} \\ & U_{21} & -L_{21}^{-1}A_{21}\overline{U}_{11}^{-1}\overline{U}_{12} \\ & & A_{32}^{s} \end{pmatrix}$$

• GEPP on G does not permute and

$$L_{31}L_{21}^{-1}A_{21}\overline{U}_{11}^{-1}\overline{U}_{12} + A_{32}^{s} = L_{31}U_{21}\overline{U}_{11}^{-1}\overline{U}_{12} + A_{32}^{s} = A_{31}\overline{U}_{11}^{-1}\overline{U}_{12} + A_{32}^{s} = \overline{L}_{31}\overline{U}_{12} + A_{32}^{s} = A_{32}$$
  
Page 29

## 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
Upper bound	2 <sup>n(log(P)+1)-1</sup>	2 <sup>n-1</sup>
In practice	n <sup>2/3</sup> n <sup>1/2</sup>	n <sup>2/3</sup> n <sup>1/2</sup>

Better bounds

## CALU – a communication avoiding LU factorization

Consider a 2D grid of P processors P<sub>r</sub>-by-P<sub>c</sub>, using a 2D block cyclic layout with square ٠ blocks of size b. U

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)
- (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

$$\log_2 P_c + \log_2 P_r))$$

 $O(n/b\log_2 P_c)$ 

O(n/b(1))



A(ib)







# LU for General Matrices

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

 $\Omega(P^{1/2})$  – just 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<sup>6</sup> 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 x b.
- The computation of each block is associated with a task.



Page 34

## 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)		$\checkmark$				
Block Cyclic Layout (BCL)		$\checkmark$	$\checkmark$			
2-level Block Layout (2I-BL)		$\checkmark$	$\checkmark$			

S. Donfack, LG, B. Gropp, V. Kale, IPDPS 2012

## 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



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





Page 38