# Case One: Communication-avoiding Factorizations

Kirk M. Soodhalter

School of Mathematics
Trinity College Dublin,
the University of Dublin
http://math.soodhalter.com

January 20, 2025

**Trinity College Dublin**
Coláiste na Tríonóide, Baile Átha Cliath
The University of Dublin

# Overview

- Solving dense linear systems
- LU- and QR-factorizations
- Stability of LU
- Parallel LU
- Communication-avoiding LU

**Linear systems of equations**

- Solve $\boldsymbol{Ax} = \boldsymbol{b}$ with $\boldsymbol{A} \in \mathbb{R}^{n \times n}$, $\boldsymbol{b} \in \mathbb{R}^n$ and $\boldsymbol{A}$ dense
- Iterative methods for *sparse matrices* or *matrix-free* problems to follow
- We consider direct methods (LU-factorization)
  - Factorize $\boldsymbol{PA} = \boldsymbol{LU}$: $\boldsymbol{L}$, $\boldsymbol{U}$ lower/upper triangular, $\boldsymbol{P}$ is a row exchange matrix
  - Solve $\boldsymbol{P}^T \boldsymbol{LUx} = \boldsymbol{b}$
  - Serial LU-factorization with partial pivoting is <span style="color:red">backwards stable</span>
  - $\|\boldsymbol{PA} - \boldsymbol{LU}\|_\infty$ small (close to $\varepsilon_{mach}$) in practice

- Given $\boldsymbol{A} \in \mathbb{R}^{m \times n}$, $\boldsymbol{b} \in: \mathbb{R}^m$, solve $\min_{\boldsymbol{x} \in \mathbb{R}^n} \|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|_2$  *Ax ≈ b*
- LS solutions satisfy normal equations: $\boldsymbol{A}^T \boldsymbol{A}\boldsymbol{x} = \boldsymbol{A}^T \boldsymbol{b}$  *b ∉ R(A)*
- For $m \geq n$, we denote as the "~~skinny~~" *narrow* QR-factorization of $\boldsymbol{A}$

*different than Gram-Schmidt version.*

$$\boldsymbol{A} = \boldsymbol{Q} \begin{bmatrix} \boldsymbol{R} \\ \underline{0} \end{bmatrix}$$

with $\boldsymbol{Q} \in \mathbb{R}^{m \times m}$ orthogonal and $\boldsymbol{R} \in \mathbb{R}^{n \times n}$ upper-triangular
- If rank $\boldsymbol{A} = $ rank $\boldsymbol{R} = n$, then LS solution given by

$$\boldsymbol{R}\boldsymbol{x} = \left(\boldsymbol{Q}^T \boldsymbol{b}\right)_{1:n}$$

- Serial QR-factorization (with appropriate orthogonalization routine) is <span style="color:red">backwards stable</span>
- $\left\| \boldsymbol{A} - \boldsymbol{Q} \begin{bmatrix} \boldsymbol{R} \\ \underline{0} \end{bmatrix} \right\|_\infty$ small (close to $\varepsilon_{mach}$) in practice

# Communication lower bounds

The same communication lower bounds shown for serial and parallel matrix-matrix multiplication have been proven to extend to other dense/direct linear algebra algorithms and more generally to many procedures involving three nested loops.

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

- Scale first column: Let $S_1 = \begin{bmatrix} 1/3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$;

$$S_1 A = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 6 & 7 & 3 \\ 9 & 12 & 3 \end{bmatrix}$$

- Eliminating first column: let $L_1 = \begin{bmatrix} 1 & & \\ -6 & 1 & \\ -9 & & 1 \end{bmatrix}$, then

$$L_1 S_1 A = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 5 & -3 \\ 0 & 9 & -6 \end{bmatrix}$$

- $$\boldsymbol{L_1 S_1 A} = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 5 & -3 \\ 0 & 9 & -6 \end{bmatrix}$$

- Scale second column: Let $\boldsymbol{S_2} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1/5 & 0 \\ 0 & 0 & 1 \end{bmatrix}$;

$$\boldsymbol{S_2 L_1 S_1 A} = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 1 & -\frac{3}{5} \\ 0 & 9 & -6 \end{bmatrix}$$

- Eliminating second column: let $\boldsymbol{L_2} = \begin{bmatrix} 1 & & \\ & 1 & \\ & -9 & 1 \end{bmatrix}$, then

$$\boldsymbol{L_2 S_2 L_1 S_1 A} = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 1 & -\frac{3}{5} \\ 0 & 0 & -\frac{3}{5} \end{bmatrix}$$

- $L_2 S_2 L_1 S_1 A = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 1 & -\frac{3}{5} \\ 0 & 0 & -\frac{3}{5} \end{bmatrix}$

- Scale third column: Let $S_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & -5/3 \end{bmatrix}$;

$$S_3 L_2 S_2 L_1 S_1 A = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 1 & -\frac{3}{5} \\ 0 & 0 & 1 \end{bmatrix}$$

- We set $L = S_3 L_2 S_2 L_1 S_1$ and $U = \begin{bmatrix} 1 & \frac{1}{3} & 1 \\ 0 & 1 & -\frac{3}{5} \\ 0 & 0 & 1 \end{bmatrix}$, and

$A = LU$

- Stability: avoid division by zero or **small elements**
  - $\|\boldsymbol{A} - \boldsymbol{LU}\|$ can be large due to roundoff

- Example: $\boldsymbol{A}_1 = \begin{bmatrix} 0 & 3 & 3 \\ 3 & 1 & 3 \\ 6 & 2 & 3 \end{bmatrix}$ has no LU-factorization

- Example: $\boldsymbol{A}_2 = \begin{bmatrix} 10^{-20} & 3 & 3 \\ 3 & 1 & 3 \\ 6 & 2 & 3 \end{bmatrix}$ has an unstable

  LU-factorization, with $\boldsymbol{L}_1 = \begin{bmatrix} 1 & & \\ -3 \times 10^{20} & 1 & \\ -6 \times 10^{20} & & 1 \end{bmatrix}$

- Partial pivoting to move col-max to diagonal:

$$\widehat{\boldsymbol{P}}\boldsymbol{A}_1 = \begin{bmatrix} 6 & 2 & 3 \\ 3 & 1 & 3 \\ 0 & 3 & 3 \end{bmatrix} \qquad \widehat{\boldsymbol{P}}\boldsymbol{A}_2 = \begin{bmatrix} 6 & 2 & 3 \\ 3 & 1 & 3 \\ 10^{-20} & 3 & 3 \end{bmatrix}$$

## Theorem

Let $\boldsymbol{PA} = \boldsymbol{LU}$ be the LU factorization with partial pivoting in exact arithmetic. Then the computed versions of these factors $\widetilde{\boldsymbol{P}}$, $\widetilde{\boldsymbol{L}}$, and $\widetilde{\boldsymbol{U}}$ satisfy

$$\widetilde{\boldsymbol{L}}\widetilde{\boldsymbol{U}} = \widetilde{\boldsymbol{P}}\boldsymbol{A} + \delta\boldsymbol{A}$$

where $\frac{\|\delta\boldsymbol{A}\|}{\boldsymbol{A}} = \mathcal{O}(\rho \cdot \varepsilon_{mach})$, and $\rho = \frac{\max_{i,j}|u_{ij}|}{\max_{i,j}|a_{ij}|}$.

$\rho$ is called the **growth factor**. $\rho \approx 1 \implies$ **stable factorization**

## Example (Worst-case instability)

In spite of partial pivoting yielding backwards stability, $\rho$ can be huge. Let $\boldsymbol{A} = \begin{bmatrix} 1 & & & & 1 \\ -1 & 1 & & & 1 \\ -1 & -1 & 1 & & 1 \\ -1 & -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & -1 & 1 \end{bmatrix}$ then

$$\boldsymbol{U} = \begin{bmatrix} 1 & & & & 1 \\ & 1 & & & 2 \\ & & 1 & & 4 \\ & & & 1 & 8 \\ & & & & 16 \end{bmatrix}$$

- Scaling this example up to $2^m \times 2^m$, we see that the largest entry of $\boldsymbol{U}$ will be $2^{m-1}$!

- Previous example is pathological. "Never" happens in practice for serial case
- Can be justified with careful statistical arguments[1]
  - Each step of Gaussian elimination introduces a rank-one correction to remaining submatrix
  - On average, this implies statistical relationships between entries of remaining submatrix tending to retard growth
- **Important:** this average stability relies on operations introducing *rank-one* corrections.
  - When moving to HPC/Parallel setting where we favor BLAS-3 and block operations (i.e., higher-rank corrections), average stability can no longer be assumed.

---

[1]Trefethen and Schreiber. *Average case stability of Gaussian elimination.* 1990

- Matrix $\boldsymbol{A} = \begin{bmatrix} \boldsymbol{A}_{11} & \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} & \boldsymbol{A}_{22} \end{bmatrix} \in \mathbb{R}^{n \times n}$ with $\boldsymbol{A}_{11} \in \mathbb{R}^{b \times b}$

  *panel* (handwritten)

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

$$\boldsymbol{P}_1 \begin{bmatrix} \boldsymbol{A}_{11} \\ \boldsymbol{A}_{21} \end{bmatrix} = \begin{bmatrix} \boldsymbol{L}_{11} \\ \boldsymbol{L}_{21} \end{bmatrix} \boldsymbol{U}_{11}$$

  (handwritten: $n \times n$, $n \times b$, $n \times b$, $b \times b$)

- We obtain the factorization

$$\boldsymbol{P}_1 \boldsymbol{A} = \begin{bmatrix} \boldsymbol{L}_{11} & \\ \boldsymbol{L}_{21} & \boldsymbol{I}_{n-b} \end{bmatrix} \begin{bmatrix} \boldsymbol{U}_{11} & \boldsymbol{U}_{12} \\ & \boldsymbol{A}_{22}^{(1)} \end{bmatrix}$$

  with $\boldsymbol{U}_{12} = \boldsymbol{L}_{11}^{-1} \boldsymbol{A}_{12}$ and $\boldsymbol{A}_{22}^{(1)} = \boldsymbol{A}_{22} - \boldsymbol{L}_{21} \boldsymbol{U}_{12}$

- Algorithm is applied recursively on $\boldsymbol{A}_{22}^{(1)}$

Compute LU of first **panel**: $P_1 \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} = \begin{bmatrix} L_{11} \\ L_{21} \end{bmatrix} U_{11}$

Apply pivoting matrix $P_1$ to full matrix $\overline{A} = P_1 A$

Solve triangular system $U_{12} = L_{11}^{-1} \overline{A}_{12}$

Update the **trailing matrix**: $\overline{A}_{22}^{(1)} = \overline{A}_{22} - L_{21} U_{12}$

Apply recursively on trailing matrix $\overline{A}_{22}^{(1)}$

**for** $i = 0, b, 2b, \ldots n$ **do**

$\quad \boldsymbol{A}^{(ib)} = \boldsymbol{A}(i : n, i : n)$

$\quad$ Factorize $i$th column panel – find a pivot for $\quad O(n \log_2 P_r)$
$\quad\quad$ each column, do row swap

$\quad$ Broadcast pivot information along all rows; $\quad O(n/b(\log_2 P_c + \log_2 P_r))$
$\quad\quad$ do swaps

$\quad$ Broadcast $\boldsymbol{L}_{11}$ along row to compute $\boldsymbol{U}_{12}$ $\quad O(n/b \log_2 P_c)$

$\quad$ Broadcast $\boldsymbol{L}_{21}$ *along rows* and $\boldsymbol{U}_{12}$ *down*
$\quad\quad$ *columns* to update trailing matrix $\quad O(n/b(\log_2 P_c + \log_2 P_r))$

**end**

[2] **Bottleneck**

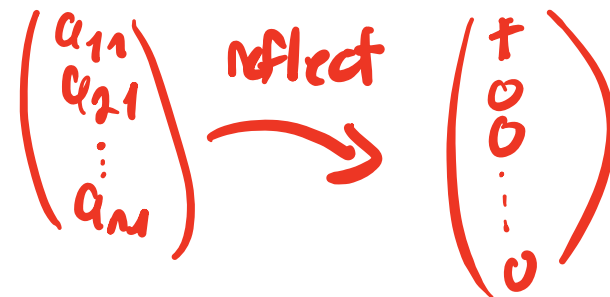- Getting pivot information from entire column across all processors

---

[2]Source: slides from Laura Grigori:
`https://people.eecs.berkeley.edu/~demmel/cs267_Spr15/Lectures/`
`lecture13_densela2_CommAvoid_UCB_Grigori_v3.pdf`

For $\boldsymbol{A} \in \mathbb{R}^{m \times n}$, $m \geq n$, w compute the "skinny" QR factorization

$$A = Q \begin{bmatrix} R \\ \underline{0} \end{bmatrix}$$

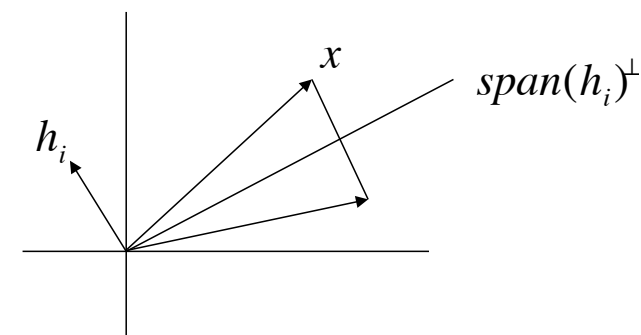- Orthogonal triangularization (Householder) vs. triangular orthogonalization (Gram-Schmidt)
- Householder matrix: $\boldsymbol{H}_i = \boldsymbol{I} - \tau_i \boldsymbol{h}_i \boldsymbol{h}_i^T$
- Reflect a vector across a hyperplane
- Choose $\boldsymbol{h}_i$ (the hyperplane) so that $\boldsymbol{x} \xrightarrow[\boldsymbol{H}_i]{} \|\boldsymbol{x}\| \boldsymbol{e}_1$
- $\boldsymbol{h}_i = \pm (\|x\| \boldsymbol{e}_1 - \boldsymbol{x})$; *sign* chosen for stability

**Apply Householder transformations to annihilate subdiagonal entries**

$$A = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} = H_1 \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ & * & * & * \\ & * & * & * \\ & * & * & * \end{bmatrix}$$

$$= H_1 \begin{bmatrix} 1 & \\ & \widetilde{H}_2 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ & r_{22} & r_{23} & r_{24} \\ & & l* & l* \\ & & l* & l* \end{bmatrix} = H_1 \begin{bmatrix} 1 & \\ & \widetilde{H}_2 \end{bmatrix} \begin{bmatrix} 1 & & \\ & 1 & \\ & & \widetilde{H}_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & r_{14} \\ & r_{22} & r_{23} & r_{24} \\ & & r_{33} & r_{34} \\ & & & r_{44} \end{bmatrix}$$

$$\underbrace{H_1 H_2 H_3}_{Q} R$$

For general $A \in \mathbb{R}^{m \times n}$, the factorization is

$$R = H_n H_{n-1} H_{n-2} \cdots H_2 H_1 A$$

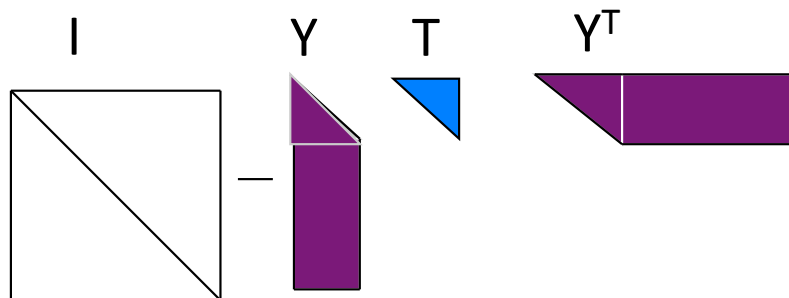$$\Longleftrightarrow (H_n H_{n-1} H_{n-2} \cdots H_2 H_1)^T R = A,$$

and $Q = H_1^T H_2^T \cdots H_{n-1}^T H_n^T$

## Represent $Q$ implicitly

$$Q = H_1 \cdots H_b = \left(I - \tau_1 h_1 h_1^T\right) \cdots \left(I - \tau_b h_b h_b^T\right) = I - Y_b T_b Y_b^T$$

where $Y_b = \begin{bmatrix} h_1 & h_2 & \cdots & h_b \end{bmatrix}$ is <span style="color:red">lower triangular</span> and $T_b$ is upper triangular:

### Example (for $b = 2$)

$$Y = \begin{bmatrix} h_1 & h_2 \end{bmatrix}, T = \begin{bmatrix} \tau_1 & -\tau_2 h_1^T h_2 \tau_2 \\ & \tau_2 \end{bmatrix}$$

$$A = \begin{bmatrix} \boldsymbol{A}_{11} & \boldsymbol{A}_{12} \\ \boldsymbol{A}_{21} & \boldsymbol{A}_{22} \end{bmatrix} \text{ where } \boldsymbol{A}_{11} \in \mathbb{R}^{b \times b}$$

## Block QR Algebra

- Compute the first panel QR-factorization
$$\boldsymbol{Q}_1^T \begin{bmatrix} \boldsymbol{A}_{11} \\ \boldsymbol{A}_{21} \end{bmatrix} = \begin{bmatrix} \boldsymbol{R}_{11} \\ \underline{0} \end{bmatrix}$$

- Update $\boldsymbol{Q}_1^T \boldsymbol{A} = \begin{bmatrix} \boldsymbol{R}_{11} & \boldsymbol{R}_{12} \\ & \widetilde{\boldsymbol{A}}_{22} \end{bmatrix}$ and apply algorithm

  recursively on trailing matrix $\widetilde{\boldsymbol{A}}_{22}$

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = Q_1 \begin{bmatrix} R_{11} & R_{12} \\ & \widetilde{A}_{12} \end{bmatrix}$$

Compute panel factorization: $Q_1^T \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} = \begin{bmatrix} R_{11} \\ \underline{0} \end{bmatrix}$

Compute compact representation: $Q_1 = I - Y_1 T_1 Y_1^T$

Update trailing matrix:

$$Q_1^T \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} = \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} - Y_1 \left( T_1^T \left( Y_1^T \begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} \right) \right) = \begin{bmatrix} R_{12} \\ \widetilde{A}_{22} \end{bmatrix}$$

Continue recursively on trailing matrix

**Bottleneck**

- We must use the entire column to build Householder vector

- We compute the QR-factorization of $\boldsymbol{W} \in \mathbb{R}^{m \times b}$ with $b \ll m$

- $p$ processors: each processor owns a block of rows of $\boldsymbol{W}$
- Classic Parallel TSQR (ScaLAPACK)
  - Compute Householder vector for each column (Bottleneck!!)
  - Compute norm of each column: $\mathcal{O}(b \log p)$ messages
- Communication Avoiding Algorithm
  - Perform local QR-factorizations, then reduce; repeat
  - $\mathcal{O}(\log p)$ messages

- We compute the QR-factorization of $\boldsymbol{W} \in \mathbb{R}^{m \times b}$ with $b \ll m$

- $p$ processors: each processor owns a block of rows of $\boldsymbol{W}$

- Classic Parallel TSQR (ScaLAPACK)
  - Compute Householder vector for each column (Bottleneck!!)
  - Compute norm of each column: $\mathcal{O}(b \log p)$ messages

- Communication Avoiding Algorithm
  - Perform local QR-factorizations, then reduce; repeat
  - $\mathcal{O}(\log p)$ messages

- We compute the QR-factorization of $\boldsymbol{W} \in \mathbb{R}^{m \times b}$ with $b \ll m$

- $p$ processors: each processor owns a block of rows of $\boldsymbol{W}$
- Classic Parallel TSQR (ScaLAPACK)
  - Compute Householder vector for each column (Bottleneck!!)
  - Compute norm of each column: $\mathcal{O}(b \log p)$ messages
- Communication Avoiding Algorithm
  - Perform local QR-factorizations, then reduce; repeat
  - $\mathcal{O}(\log p)$ messages

- We compute the QR-factorization of $\boldsymbol{W} \in \mathbb{R}^{m \times b}$ with $b \ll m$

- $p$ processors: each processor owns a block of rows of $\boldsymbol{W}$

- Classic Parallel TSQR (ScaLAPACK)
  - Compute Householder vector for each column (Bottleneck!!)
  - Compute norm of each column: $\mathcal{O}(b \log p)$ messages

- Communication Avoiding Algorithm
  - Perform local QR-factorizations, then reduce; repeat
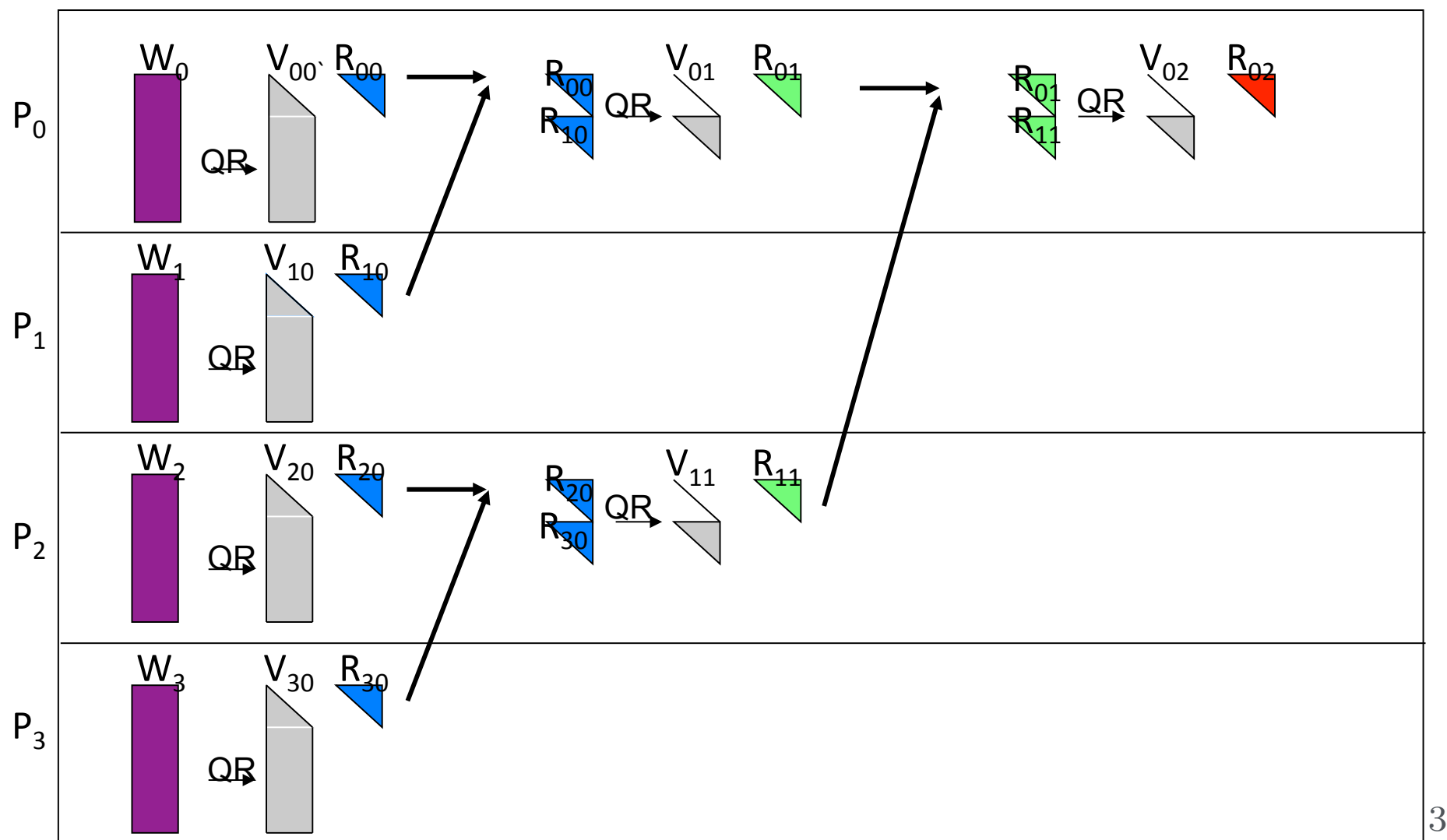  - $\mathcal{O}(\log p)$ messages

- Stage 1:

$$W = \begin{bmatrix} W_0 \\ \hline W_1 \\ \hline W_2 \\ \hline W_3 \end{bmatrix} = \begin{bmatrix} Q_{00}R_{00} \\ \hline Q_{10}R_{10} \\ \hline Q_{20}R_{20} \\ \hline Q_{30}R_{30} \end{bmatrix} = \begin{bmatrix} Q_{00} & & & \\ \hline & Q_{10} & & \\ \hline & & Q_{20} & \\ \hline & & & Q_{30} \end{bmatrix} \cdot \begin{bmatrix} R_{00} \\ \hline R_{10} \\ \hline R_{20} \\ \hline R_{30} \end{bmatrix}$$

$$b = \begin{pmatrix} b_0 \\ \hline b_1 \\ \hline b_2 \\ \hline b_3 \end{pmatrix}$$

- Stages 2 & 3:

$$\underbrace{\begin{bmatrix} R_{00} \\ R_{10} \\ R_{20} \\ R_{30} \end{bmatrix} = \begin{bmatrix} Q_{01}R_{01} \\ Q_{11}R_{11} \end{bmatrix} = \begin{bmatrix} Q_{01} & \\ & Q_{11} \end{bmatrix} \cdot \begin{bmatrix} R_{01} \\ R_{11} \end{bmatrix}}_{\text{Stage2}}; \quad \underbrace{\begin{bmatrix} R_{01} \\ R_{11} \end{bmatrix} = Q_{02}R_{02}}_{\text{Stage3}}$$

- $Q = \begin{bmatrix} Q_{00} & & & \\ & Q_{10} & & \\ & & Q_{20} & \\ & & & Q_{30} \end{bmatrix} \cdot \begin{bmatrix} Q_{01} & \\ & Q_{11} \end{bmatrix} \cdot Q_{02}$ and $R = R_{02}$

- **Important:** The QR-factorization is <u>unique</u> up to the sign of the columns of $Q$.

## Flexibility of TSQR and CAQR algorithms

**Parallel:** $W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix}$

$R_{00}$   $R_{01}$

$R_{10}$     $R_{02}$

$R_{20}$

$R_{30}$   $R_{11}$

**Sequential:** $W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix}$

$R_{00}$   $R_{01}$   $R_{02}$   $R_{03}$

**Dual Core:** $W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix}$

$R_{00}$   $R_{01}$   $R_{02}$   $R_{03}$

$R_{01}$   $R_{11}$   $R_{11}$

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

Solving LS: $Rx = (Q^T b)_{1:n}$

what if I want to keep $Q$ & $R$?

- $A \in \mathbb{R}^{m \times n}$, how do we best lay out on processors? What mapping?

- In some of your readings, they show how any processor layout can accomodation communication-avoiding factorization.

- Optimal block size $b$?

- Use communication-avoiding TSQR for panel factorization

- Broadcast $Q_{ij}$ across appropriate rows after each panel factorization?

$Q = A R^{-1}$

Keep only $R$.
If you need $Q$

**Tall skinny matrix $W$:**

$$W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ W_3 \end{bmatrix} = \underbrace{\begin{bmatrix} \Pi_{00} & & & \\ & \Pi_{10} & & \\ & & \Pi_{20} & \\ & & & \Pi_{30} \end{bmatrix}}_{\Pi_0} \cdot \begin{bmatrix} L_{00} & & & \\ & L_{10} & & \\ & & L_{20} & \\ & & & L_{30} \end{bmatrix} \cdot \begin{bmatrix} U_0 \\ U_1 \\ U_2 \\ U_3 \end{bmatrix}$$
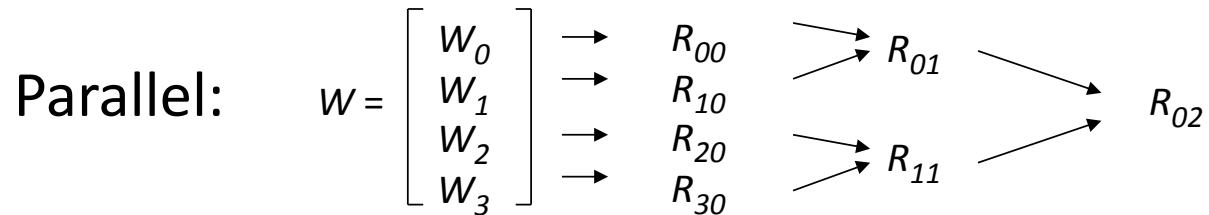
$$\begin{bmatrix} U_0 \\ U_1 \\ U_2 \\ U_3 \end{bmatrix} = \underbrace{\begin{bmatrix} \Pi_{01} & \\ & \Pi_{11} \end{bmatrix}}_{\Pi_1} \cdot \begin{bmatrix} L_{01} & \\ & L_{11} \end{bmatrix} \cdot \begin{bmatrix} U_{01} \\ U_{11} \end{bmatrix}$$

$$\begin{bmatrix} U_{01} \\ U_{11} \end{bmatrix} = \Pi_2 L_{02} U_{02}$$

# Block parallel pivoting



average growth factor (partial pivoting;b= 1,2,4,8,16,32)

- Recall: LU stability controlled by a growth factor $\rho$ of the entries of $\boldsymbol{U}$

- Unstable for large numbers of processors

- p = # of rows corresponds to "parallel pivoting" which is known to be unstable [Trefethen and Schreiber, '90]

[4]Source: slides from Laura Grigori

At each iteration, we have

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \text{with} W = \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} \in \mathbb{R}^{n \times b}$$

- Preprocess $W$ to find good pivots, getting $P$
- Perform all permutations, i.e., compute $PA$
- Perform LU without pivots for updated $W$, update trailing matrix

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

- We do a reduction with LU at each step, **but only to get pivots**, and **we only pick the top $b$ pivot rows** for each block to broadcast.

- Compute $LU$ of each block $\boldsymbol{W}_i$ to get
$\boldsymbol{\Pi}_0 = \operatorname{diag} \operatorname{curl} \boldsymbol{\Pi}_{00}, \boldsymbol{\Pi}_{10}, \boldsymbol{\Pi}_{20}, \boldsymbol{\Pi}_{30}$

$$
\boldsymbol{W} = \begin{bmatrix} \boldsymbol{W}_0 \\ \hline \boldsymbol{W}_1 \\ \hline \boldsymbol{W}_2 \\ \hline \boldsymbol{W}_3 \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Pi}_{00} \boldsymbol{L}_{00} \boldsymbol{U}_{00} \\ \hline \boldsymbol{\Pi}_{10} \boldsymbol{L}_{10} \boldsymbol{U}_{10} \\ \hline \boldsymbol{\Pi}_{20} \boldsymbol{L}_{20} \boldsymbol{U}_{20} \\ \hline \boldsymbol{\Pi}_{30} \boldsymbol{L}_{30} \boldsymbol{U}_{30} \end{bmatrix}
$$

- Apply pivots to get $\boldsymbol{W}_{i0} = \left( \boldsymbol{\Pi}_{i0}^T \boldsymbol{W}_i \right)_{1:b}$. Get $\boldsymbol{\Pi}_1$

$$
\begin{bmatrix} \boldsymbol{W}_{00} \\ \boldsymbol{W}_{10} \\ \hline \boldsymbol{W}_{20} \\ \boldsymbol{W}_{30} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Pi}_{01} \boldsymbol{L}_{01} \boldsymbol{U}_{01} \\ \hline \boldsymbol{\Pi}_{11} \boldsymbol{L}_{11} \boldsymbol{U}_{11} \end{bmatrix}
$$

- Apply pivots to get $\boldsymbol{W}_{i1} = \left( \boldsymbol{\Pi}_{i1}^T \boldsymbol{W}_{i0} \right)_{1:b}$. Get $\boldsymbol{\Pi}_2$ with
$\begin{bmatrix} \boldsymbol{W}_{i0} \\ \boldsymbol{W}_{i1} \end{bmatrix} = \boldsymbol{\Pi}_2 \boldsymbol{L}_{02} \boldsymbol{U}^{02}$.

- Computed the *unpivoted* LU factorization $\boldsymbol{\Pi}_2^T \boldsymbol{\Pi}_1^T \boldsymbol{\Pi}_0^T \boldsymbol{W} = \boldsymbol{L}\boldsymbol{U}$.

$P_0$
$$W_0 = \begin{pmatrix} 2 & 4 \\ 0 & 1 \\ 2 & 0 \\ 1 & 2 \end{pmatrix} = \Pi_0 L_0 U_0 \qquad \Pi_0^T W_0 = \begin{pmatrix} 2 & 4 \\ 2 & 0 \end{pmatrix}$$

$$\overline{W}_0 = \begin{pmatrix} 2 & 4 \\ 2 & 0 \\ 4 & 1 \\ 2 & 0 \end{pmatrix} = \overline{\Pi}_0 \overline{L}_0 \overline{U}_0 \qquad \overline{\Pi}_0^T \overline{W}_0 = \begin{pmatrix} 4 & 1 \\ 2 & 4 \end{pmatrix}$$

$$\underline{W}_0 = \begin{pmatrix} 4 & 1 \\ 2 & 4 \\ 4 & 2 \\ 1 & 4 \end{pmatrix} = \underline{\Pi}_0 \underline{L}_0 \underline{U}_0 \qquad \underline{\Pi}_0^T \underline{W}_0 = \begin{pmatrix} 4 & 1 \\ 1 & 4 \end{pmatrix}$$

Good pivots for factorizing W

$P_1$
$$W_1 = \begin{pmatrix} 2 & 0 \\ 0 & 0 \\ 4 & 1 \\ 1 & 0 \end{pmatrix} = \Pi_1 L_1 U_1 \qquad \Pi_1^T W_1 = \begin{pmatrix} 4 & 1 \\ 2 & 0 \end{pmatrix}$$

$P_2$
$$W_2 = \begin{pmatrix} 0 & 1 \\ 1 & 4 \\ 0 & 0 \\ 0 & 2 \end{pmatrix} = \Pi_2 L_2 U_2 \qquad \Pi_2^T W_2 = \begin{pmatrix} 1 & 4 \\ 0 & 2 \end{pmatrix}$$

$$\overline{W}_2 = \begin{pmatrix} 1 & 4 \\ 0 & 2 \\ 4 & 2 \\ 0 & 2 \end{pmatrix} = \overline{\Pi}_2 \overline{L}_2 \overline{U}_2 \qquad \overline{\Pi}_2^T \overline{W}_2 = \begin{pmatrix} 4 & 2 \\ 1 & 4 \end{pmatrix}$$

$P_3$
$$W_3 = \begin{pmatrix} 2 & 1 \\ 0 & 2 \\ 1 & 0 \\ 4 & 2 \end{pmatrix} = \Pi_3 L_3 U_3 \qquad \Pi_3^T W_3 = \begin{pmatrix} 4 & 2 \\ 0 & 2 \end{pmatrix}$$

time →

- Experimentally: CALU has been demonstrated stable (with low growth factor $\rho$) for large random matrices with normally distributed entries
- Theoretically: CALU equivalent **in exact arithmetic** to BLAS-3 serial block LU with p.p. applied to an auxiliary matrix.
- Example: Let

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

- One round of tournament pivoting produces pivot matrices $\Pi_0$ and $\Pi_1$. Let $\begin{bmatrix} \overline{A}_{11} & \overline{A}_{12} \\ \overline{A}_{21} & \overline{A}_{22} \\ \overline{A}_{31} & \overline{A}_{32} \end{bmatrix} = \Pi_1^T \Pi_0^T A$

- Equivalent to LU with p.p. applied to
$$G = \begin{bmatrix} \overline{A}_{11} & & \overline{A}_{12} \\ A_{21} & A_{21} & \\ & -A_{31} & A_{32} \end{bmatrix}$$

## Example

You will be asked to implement a parallel communication-avoiding QR factorization code, use this code to solve some least-squares problems, and run some timing and scaling tests.

- The two most widely used matrix factorizations LU and QR, which are used for directly solving, resp., linear systems and least squares problems

- The stability of LU with partial pivoting

- Parallel and then communication avoiding QR

- The naive analog of CALU followed by the stable version using tournament pivoting