

A communication-avoiding parallel algorithm for the symmetric eigenvalue problem

Edgar Solomonik¹, Grey Ballard², James Demmel³, Torsten Hoefer⁴

(1) Department of Computer Science, University of Illinois at Urbana-Champaign

(2) Department of Computer Science, Wake Forest University

(3) Department of Computer Science and Department of Mathematics, University of California, Berkeley

(4) Department of Computer Science, ETH Zurich

Symposium for Parallel Architectures and Algorithms (SPAA), 2017

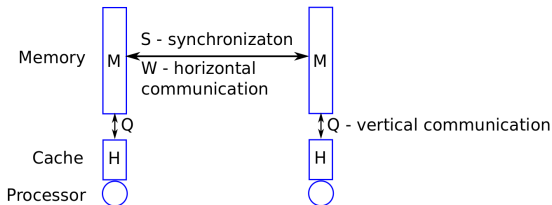
July 24, 2017

 @CS@Illinois

Beyond computational complexity

Algorithms should minimize communication, not just computation

- communication and synchronization cost more **energy** than flops
- two types of communication (data movement):



- **vertical** (intranode memory–cache)
- **horizontal** (internode network transfers)
- parallel algorithm design involves tradeoffs: computation vs communication vs synchronization
- parameterized algorithms provide optimality and flexibility

Cost model for parallel algorithms

We use the **Bulk Synchronous Parallel (BSP) model** (L.G. Valiant 1990)

- execution is subdivided into S supersteps, each associated with a global **synchronization** (cost α)
- at the start of each superstep, processors interchange messages, then they perform local computation
- if the **maximum amount of data** sent or received by any process is m_i (and work done is f_i) at superstep i then the BSP time is

$$T = \sum_{i=1}^S \alpha + m_i \cdot \beta + f_i \cdot \gamma = O(S \cdot \alpha + W \cdot \beta + F \cdot \gamma)$$

We additionally consider **vertical communication cost**

- F – computation cost (local computation)
- Q – vertical communication cost (memory–cache traffic)
- W – horizontal communication cost (interprocessor communication)
- S – synchronization cost (number of supersteps)

Symmetric eigenvalue problem

Given a dense symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ find diagonal matrix \mathbf{D} so

$$\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{D}$$

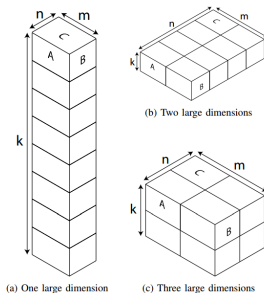
where \mathbf{X} is an orthogonal matrix composed of eigenvectors of \mathbf{A}

- **diagonalization** – reduction of \mathbf{A} to diagonal matrix \mathbf{D}
- computing the SVD has very similar computational structure
- we focus on tridiagonalization (bidiagonalization for SVD), from which standard approaches (e.g. MRRR, see Dhillon, Parlett, Vömel, 2006) can be used
- core building blocks:
 - **matrix multiplication**
 - **QR factorization**
- QR, SVD, diagonalization of large matrices are needed for applications in scientific computing, data analysis and beyond

Communication complexity of matrix multiplication

Multiplication of $\mathbf{A} \in \mathbb{R}^{m \times k}$ and $\mathbf{B} \in \mathbb{R}^{k \times n}$ can be done in $O(1)$ supersteps with **communication cost** $W = O\left(\left(\frac{mnk}{p}\right)^{2/3}\right)$ provided sufficiently memory and sufficiently large p

- when $m = n = k$, 3D blocking gets $O(p^{1/6})$ improvement over $2D^1$
- when m, n, k are unequal, need appropriate processor grid²



¹J. Berntsen, Par. Comp., 1989; A. Aggarwal, A. Chandra, M. Snir, TCS, 1990; R.C. Agarwal, S.M. Balle, F.G. Gustavson, M. Joshi, P. Palkar, IBM, 1995; F.W. McColl, A. Tiskin, Algorithmica, 1999; ...

²J. Demmel, D. Eliahu, A. Fox, S. Kamil, B. Lipshitz, O. Schwartz, O. Spillinger 2013

Bandwidth-efficient QR and diagonalization

Goal: achieve the same communication complexity for QR and diagonalization as for matrix multiplication

- synchronization complexity expected to be higher

$$W \cdot S = \Omega(n^2)$$

product of communication and synchronization cost must be greater than the square of the number of columns³

- general strategy
 - 1 use communication-efficient matrix-multiplication for QR
 - 2 use communication-efficient QR for diagonalization

³E.S., E. Carson, N. Knight, J. Demmel, SPAA 2014 (TOPC 2016)

Communication complexity of dense matrix kernels

For $n \times n$ Cholesky with p processors, optimal parallel schedules attain

$$F = O(n^3/p), \quad W = O(n^2/p^\delta), \quad S = O(p^\delta)$$

for any $\delta = [1/2, 2/3]$.

Achieving similar costs for LU, QR, and the symmetric eigenvalue problem requires some algorithmic tweaks.

triangular solve	square TRSM \checkmark^4	rectangular TRSM \checkmark^5
LU with pivoting	pairwise pivoting \checkmark^6	tournament pivoting \checkmark^7
QR factorization	Givens on square \checkmark^3	Householder on rect. \checkmark^8
SVD (sym. eig.)	singular values only \checkmark^8	singular vectors \times

\checkmark means costs attained (synchronization within polylogarithmic factors).

⁴B. Lipshitz, MS thesis 2013

⁵T. Wicky, E.S., T. Hoefler, IPDPS 2017

⁶A. Tiskin, FGCS 2007

⁷E.S., J. Demmel, EuroPar 2011

⁸E.S., G. Ballard, T. Hoefler, J. Demmel, SPAA 2017

QR factorization of tall-and-skinny matrices

Consider the reduced factorization $\mathbf{A} = \mathbf{QR}$ with $\mathbf{A}, \mathbf{Q} \in \mathbb{R}^{m \times n}$ and $\mathbf{R} \in \mathbb{R}^{n \times n}$ when $m \gg n$ (in particular $m \geq np$)

- \mathbf{A} is tall-and-skinny, each processor owns a block of rows
- Householder-QR requires $S = \Theta(n)$ supersteps, $W = O(n^2)$
- Cholesky-QR2, TSQR, and HR-TSQR only $S = \Theta(\log(p))$ supersteps
- TSQR⁹: row-recursive divide-and-conquer, $W = O(n^2 \log(p))$

$$\begin{bmatrix} \mathbf{Q}_1 \mathbf{R}_1 \\ \mathbf{Q}_2 \mathbf{R}_2 \end{bmatrix} = \begin{bmatrix} \text{TSQR}(\mathbf{A}_1) \\ \text{TSQR}(\mathbf{A}_2) \end{bmatrix}, [\mathbf{Q}_{12}, \mathbf{R}] = \text{QR} \left(\begin{bmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \end{bmatrix} \right), \mathbf{Q} = \begin{bmatrix} \mathbf{Q}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{bmatrix} \mathbf{Q}_{12}$$

- TSQR-HR¹⁰: TSQR + Householder-reconstruction, $W = O(n^2 \log(p))$
- Cholesky-QR2¹¹: stable so long as $\kappa(\mathbf{A}) \leq 1/\sqrt{\epsilon}$, achieves $W = O(n^2)$

⁹J. Demmel, L. Grigori, M. Hoemmen, J. Langou 2012

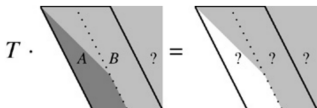
¹⁰G. Ballard, J. Demmel, L. Grigori, M. Jacquelin, H.-D. Nguyen, E.S. 2014

¹¹Y. Yamamoto, Y. Nakatsukasa, Y. Yanagisawa, T. Fukaya 2015

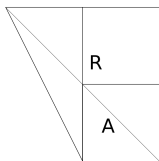
QR factorization of square matrices

Square matrix QR algorithms generally use 1D QR for panel factorization

- algorithms in ScaLAPACK, Elemental, DPLASMA use **2D layout**, generally achieve $W = O(n^2/\sqrt{p})$ cost
- Tiskin's 3D QR algorithm¹² achieves $W = O(n^2/p^{2/3})$ communication



- however, requires **slanted-panel matrix embedding**



which is highly inefficient for rectangular (tall-and-skinny) matrices

¹²A. Tiskin 2007, "Communication-efficient generic pairwise elimination"

Communication-avoiding rectangular QR

For $\mathbf{A} \in \mathbb{R}^{m \times n}$ existing algorithms are optimal when $m = n$ and $m \gg n$

- cases with $n < m < np$ underdetermined equations are important
- new algorithm
 - subdivide p processors into m/n groups of pn/m processors
 - perform row-recursive QR (TSQR) with tree of height $\log_2(m/n)$
 - compute each tree-node elimination $\text{QR}\left(\begin{bmatrix} \mathbf{R}_1 \\ \mathbf{R}_2 \end{bmatrix}\right)$ using Tiskin's QR with pn/m or more processors
- note: interleaving rows of \mathbf{R}_1 and \mathbf{R}_2 gives a slanted panel!
- obtains ideal communication cost for any m, n , generally

$$W = O\left(\left(\frac{mn^2}{p}\right)^{2/3}\right)$$

Reducing the symmetric matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ to a tridiagonal matrix

$$\mathbf{T} = \mathbf{Q}^T \mathbf{A} \mathbf{Q}$$

via a **two-sided orthogonal transformation** is most costly in diagonalization

- can be done by **successive column QR factorizations**

$$\mathbf{T} = \underbrace{\mathbf{Q}_1^T \cdots \mathbf{Q}_n^T}_{\mathbf{Q}^T} \mathbf{A} \underbrace{\mathbf{Q}_1 \cdots \mathbf{Q}_n}_{\mathbf{Q}}$$

- two-sided updates harder to manage than one-sided
- can use n/b QRs on panels of b columns to go to band-width $b + 1$
- $b = 1$ gives direct tridiagonalization

Multi-stage tridiagonalization

Writing the orthogonal transformation in Householder form, we get

$$\underbrace{(I - \mathbf{U}\mathbf{T}\mathbf{U}^T)^T}_{\mathbf{Q}^T} \mathbf{A} \underbrace{(I - \mathbf{U}\mathbf{T}\mathbf{U}^T)}_{\mathbf{Q}} = \mathbf{A} - \mathbf{U}\mathbf{V}^T - \mathbf{V}\mathbf{U}^T$$

where \mathbf{U} are Householder vectors and \mathbf{V} is

$$\mathbf{V}^T = \mathbf{T}\mathbf{U}^T + \frac{1}{2} \mathbf{T}^T \mathbf{U}^T \underbrace{\mathbf{A}\mathbf{U}}_{\text{challenge}} \mathbf{T}\mathbf{U}^T$$

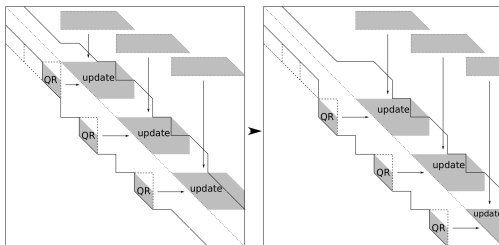
- when performing two-sided updates, computing $\mathbf{A}\mathbf{U}$ dominates cost
- if $b = 1$, \mathbf{U} is a column-vector, and $\mathbf{A}\mathbf{U}$ is dominated by [vertical communication cost](#) (moving \mathbf{A} between memory and cache)
- [idea](#): reduce to banded matrix ($b \gg 1$) first¹³

¹³T. Auckenthaler, H.-J. Bungartz, T. Huckle, L. Krämer, B. Lang, P. Willems 2011

Successive band reduction (SBR)

After reducing to a banded matrix, we need to transform the banded matrix to a tridiagonal one

- fewer nonzeros lead to lower computational cost, $F = O(n^2 b/p)$
- however, transformations introduce **fill/bulges**
- bulges must be chased down the band¹⁴



- communication- and synchronization-efficient **1D SBR algorithm** known for small band-width¹⁵

¹⁴ B. Lang 1993; C. Bischof, B. Lang, X. Sun 2000

¹⁵ G. Ballard, J. Demmel, N. Knight 2012

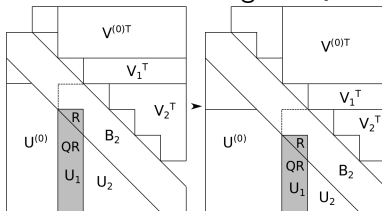
Communication-efficient eigenvalue computation

Previous work (start-of-the-art): **two-stage tridiagonalization**

- implemented in ELPA, can outperform ScaLAPACK¹⁶
- with $n = n/\sqrt{p}$, 1D SBR gives $W = O(n^2/\sqrt{p})$, $S = O(\sqrt{p} \log^2(p))$ ¹⁷

We show the benefits of **many-stage tridiagonalization**

- use $\Theta(\log(p))$ intermediate band-widths to achieve $W = O(n^2/p^{2/3})$
- leverage communication-efficient rectangular QR with processor groups



- 3D SBR (each QR and matrix multiplication update parallelized)

¹⁶ T. Auckenthaler, H.-J. Bungartz, T. Huckle, L. Krämer, B. Lang, P. Willems 2011

¹⁷ G. Ballard, J. Demmel, N. Knight 2012

Symmetric eigensolver results summary

Algorithm	W	Q	S
ScaLAPACK	n^2/\sqrt{p}	n^3/p	$n \log(p)$
ELPA	n^2/\sqrt{p}	-	$n \log(p)$
two-stage + 1D-SBR	n^2/\sqrt{p}	$n^2 \log(n)/\sqrt{p}$	$\sqrt{p}(\log^2(p) + \log(n))$
many-stage	$n^2/p^{2/3}$	$n^2 \log p/p^{2/3}$	$p^{2/3} \log^2 p$

- costs are asymptotic (same computational cost F for eigenvalues)
- W – horizontal (interprocessor) communication
- Q – vertical (memory-cache) communication, excluding $W + F/\sqrt{H}$
- S – synchronization cost (number of supersteps)

Summary of contributions

- communication-efficient **QR factorization** algorithm
 - optimal communication cost for any matrix dimensions
 - variants that trade-off some accuracy guarantees for performance
- communication-efficient **symmetric eigensolver** algorithm
 - reduce matrix to successively smaller band-width
 - uses concurrent executions of 3D matrix multiplication and 3D QR

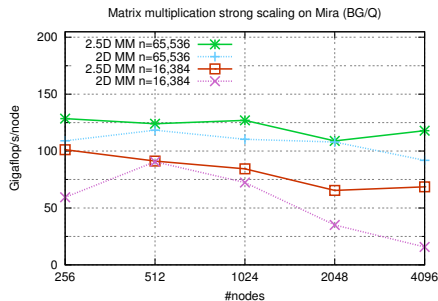
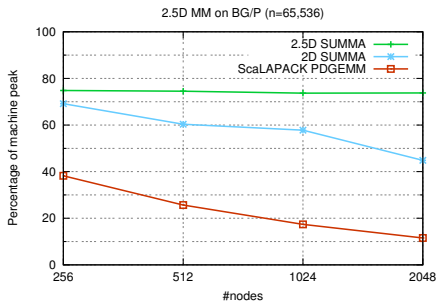
Practical implications

- ELPA demonstrated efficacy of two-stage approach, **our work motivates 3+ stages**
- partial parallel implementation is competitive but no speed-up

Future work

- back-transformations to compute **eigenvectors** in less computational complexity than $F = O(n^3 \log(p)/p)$
- QR with **column pivoting** / low-rank SVD

Communication-efficient matrix multiplication



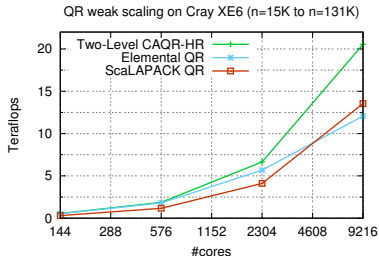
12X speed-up, 95% reduction in comm. for $n = 8K$ on 16K nodes of BG/P

Communication-efficient QR factorization

- Householder form can be reconstructed quickly from TSQR¹⁸

$$Q = I - YTY^T \quad \Rightarrow \quad LU(I - Q) \rightarrow (Y, TY^T)$$

- Householder aggregation yields performance improvements

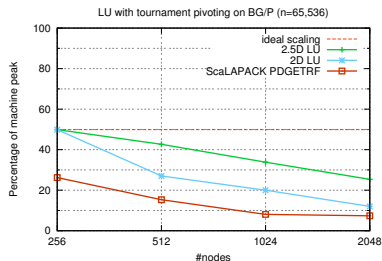


¹⁸Ballard, Demmel, Grigori, Jacquelin, Nguyen, S., IPDPS, 2014

Communication-efficient LU factorization

For any $c \in [1, p^{1/3}]$, use cn^2/p memory per processor and obtain

$$W_{LU} = O(n^2/\sqrt{cp}), \quad S_{LU} = O(\sqrt{cp})$$



- LU with pairwise pivoting¹⁹ extended to tournament pivoting²⁰
- first implementation of a communication-optimal LU algorithm¹¹

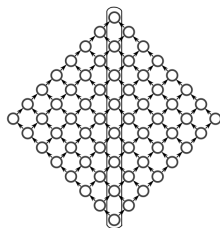
¹⁹Tiskin, FGCS, 2007

²⁰S., Demmel, Euro-Par, 2011

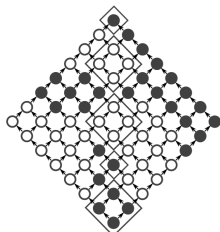
Tradeoffs in the diamond DAG

Computation vs synchronization tradeoff for the $n \times n$ diamond DAG,²¹

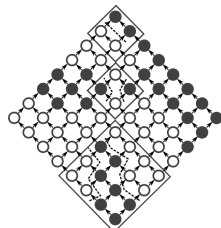
$$F \cdot S = \Omega(n^2)$$



Dependency chain P



Monochrome dependency intervals



Multicolored dependency intervals

We generalize this idea²²

- additionally consider horizontal communication
- allow arbitrary (polynomial or exponential) interval expansion

²¹Papadimitriou, Ullman, SIAM JC, 1987

²²S., Carson, Knight, Demmel, SPAA 2014 (extended version, JPDC 2016)

Tradeoffs involving synchronization

We apply tradeoff lower bounds to dense linear algebra algorithms, represented via dependency hypergraphs:²³

For triangular solve with an $n \times n$ matrix,

$$F_{\text{TRSV}} \cdot S_{\text{TRSV}} = \Omega(n^2)$$

For Cholesky of an $n \times n$ matrix,

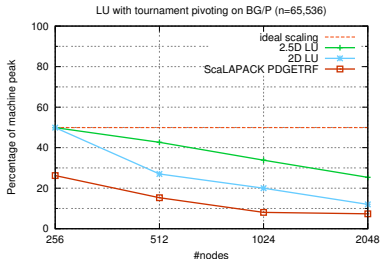
$$F_{\text{CHOL}} \cdot S_{\text{CHOL}}^2 = \Omega(n^3) \quad W_{\text{CHOL}} \cdot S_{\text{CHOL}} = \Omega(n^2)$$

²³S., Carson, Knight, Demmel, SPAA 2014 (extended version, JPDC 2016)

Communication-efficient LU factorization

For any $c \in [1, p^{1/3}]$, use cn^2/p memory per processor and obtain

$$W_{\text{LU}} = O(n^2/\sqrt{cp}), \quad S_{\text{LU}} = O(\sqrt{cp})$$



- LU with pairwise pivoting²⁴ extended to tournament pivoting²⁵
- first implementation of a communication-optimal LU algorithm¹⁰

²⁴Tiskin, FGCS, 2007

²⁵S., Demmel, Euro-Par, 2011