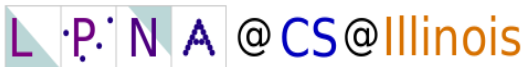


# communication-avoiding Cholesky-QR2 for rectangular matrices (CA-CQR2)

Edward Hutter and Edgar Solomonik

Laboratory for Parallel Numerical Algorithms  
Department of Computer Science  
University of Illinois at Urbana-Champaign

IPDPS 2019



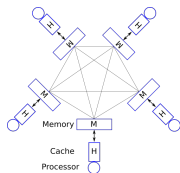
# Motivation for reducing algorithmic communication costs

Communication and synchronization increasingly dominating algorithm performance on modern architectures

$\alpha - \beta - \gamma$  cost model

- $\alpha$  - cost to send zero-byte message
- $\beta$  - cost to inject byte of data into network
- $\gamma$  - cost to perform flop with register-resident data

Architectural trend:  $\alpha \gg \beta \gg \gamma$



$$T_{\text{near-neighbor-exchange}}(n,P) = a + n\beta$$

$$T_{\text{all-reduce}}(n,P) = f(P)\alpha + f(P)n\beta$$

Figure: Horizontal (internode network) communication along critical path

Communication-avoiding algorithms for **most** dense matrix factorizations present in numerical libraries

**Goal: A QR factorization algorithm that prioritizes minimizing synchronization and communication cost**

3D algorithms utilize available extra memory to reduce communication asymptotically.

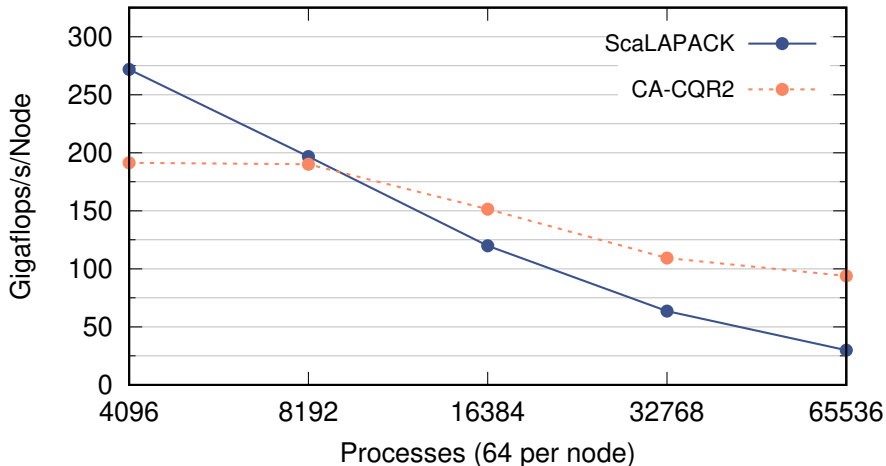
We introduce CA-CQR2, a novel practical 3D QR factorization algorithm

- extends CholeskyQR2 algorithm to arbitrary matrices
- requires  $\mathcal{O}\left(\left(\frac{Pm^2}{n^2}\right)^{\frac{1}{6}}\right)$  less communication than known 2D QR algorithms for  $m \times n$  matrices across  $P$  processes
- obtains 3x speedups over ScaLAPACK on 1024 nodes
- utilizes first distributed-memory implementation of recursive 3D Cholesky factorization

CA-CQR2's asymptotic communication reduction incurs tradeoffs

- increased computation (2 – 4x more flops than Householder QR)
- constrained applicability (matrix must be sufficiently well-conditioned)
- requires  $\mathcal{O}\left(\left(\frac{Pm}{n}\right)^{\frac{1}{3}}\right)$  more memory than known 2D QR algorithms for  $m \times n$  matrices across  $P$  processes

## Strong Scaling on Stampede2, 8388608 x 2048 matrix

Figure: Strong scaling for  $m \times n$  matrices

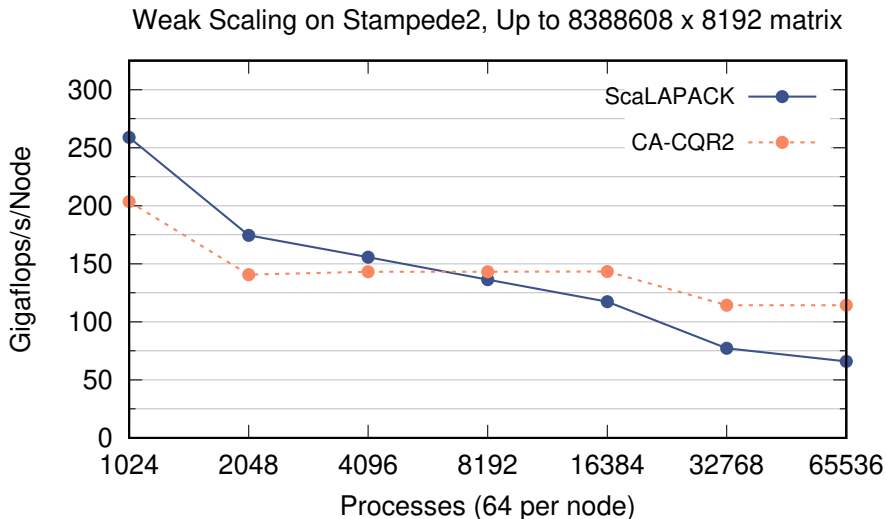


Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

# Competing costs of parallel QR factorization of $A_{m \times n}$

$\alpha - \beta$  model captures communication ( $\beta$ ) and synchronization ( $\alpha$ ) costs over  $P$  processors

ScaLAPACK's PGEQRF is communication-optimal assuming minimal memory (2D)

$$T_{\text{PGEQRF}}^{\alpha, \beta} = \mathcal{O} \left( n \log P \cdot \alpha + \frac{mn}{\sqrt{P}} \cdot \beta \right) \quad M_{\text{PGEQRF}} = \mathcal{O} \left( \frac{mn}{P} \right)$$

CAQR factors panels using TSQR to reduce synchronization<sup>1</sup> (2D)

$$T_{\text{CAQR}}^{\alpha, \beta} = \mathcal{O} \left( \sqrt{P} \log^2 P \cdot \alpha + \frac{mn}{\sqrt{P}} \cdot \beta \right) \quad M_{\text{CAQR}} = \mathcal{O} \left( \frac{mn}{P} \right)$$

CA-CQR2 leverages extra memory to reduce communication (3D)

$$T_{\text{CA-CQR2}}^{\alpha, \beta} = \mathcal{O} \left( \left( \frac{Pn}{m} \right)^{\frac{2}{3}} \log P \cdot \alpha + \left( \frac{n^2 m}{P} \right)^{\frac{2}{3}} \cdot \beta \right) \quad M_{\text{CA-CQR2}} = \mathcal{O} \left( \left( \frac{n^2 m}{P} \right)^{\frac{2}{3}} \right)$$

3D algorithms exist in theory<sup>2 3 4</sup>, but **CA-CQR2 is the first practical approach**

<sup>1</sup>J. Demmel et al., "Communication-optimal Parallel and Sequential QR and LU Factorizations", SISC 2012

<sup>2</sup>A. Tiskin, "Communication-efficient generic pairwise elimination", Future Generation Computer Systems 2007

<sup>3</sup>E. Solomonik et al., "A communication-avoiding parallel algorithm for the symmetric eigenvalue problem", SPAA 2017

<sup>4</sup>G. Ballard et al., "A 3D Parallel Algorithm for QR Decomposition", SPAA 2018

QR factorization algorithms used in practice stem from processes of orthogonal triangularization for their superior numerical stability

$$Q_n Q_{n-1} \dots Q_1 A = R$$

The Cholesky-QR algorithm is a simple algorithm that follows a numerically unstable process of triangular orthogonalization

$$AR_1^{-1} R_2^{-1} \dots R_n^{-1} = Q$$

---

$[Q, R] \leftarrow$  **Cholesky-QR** ( $A$ )

---

$$B \leftarrow A^T A$$

▷  $B$  may be indefinite!

$$R^T R \leftarrow B$$

▷ Possible failure in Cholesky factorization!

$$Q \leftarrow AR^{-1}$$

▷  $R$  may have lost all accuracy!  $Q$  may lost orthogonality!

---

The Cholesky-QR2 algorithm *can* achieve stability through iterative refinement<sup>1</sup>

---

$$[Q, R] \leftarrow \text{Cholesky-QR2}(A)$$

---

$$Z, R_1 \leftarrow \text{CQR}(A)$$

$$Q, R_2 \leftarrow \text{CQR}(Z)$$

$$R \leftarrow R_2 R_1$$

---

- leverages near-perfect conditioning of  $Z$  in a second iteration<sup>1</sup>
- $A = ZR_1 = QR_2R_1$ , from  $A^T A = R_1^T Z^T Z R_1 = R_1^T R_2^T Q^T Q R_2 R_1$ , where  $R_2$  corrects initial  $R_1$
- numerical breakdown still possible if first iteration loses positive definiteness in  $A^T A$  via  $\kappa(A) \leq 1/\sqrt{\epsilon}$

Shifted Cholesky-QR<sup>2</sup> can attain a stable factorization for any matrix  $\kappa(A) \leq 1/\epsilon$

- the eigenvalues of  $A^T A$  are shifted to prevent loss of positive definiteness
- three Cholesky-QR iterations required, essentially 3 – 6x more flops than Householder approaches

---

<sup>1</sup>Y. Yamamoto et al., "Roundoff Error Analysis of the CholeskyQR2 algorithm", Electron. Trans. Numer. Anal. 2015

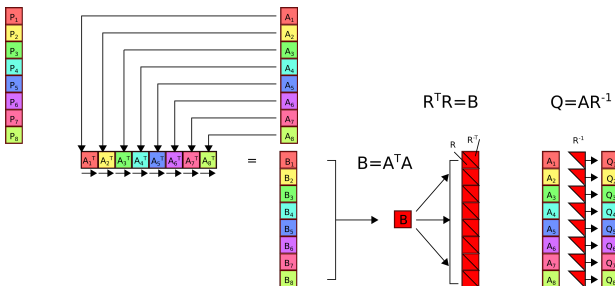
<sup>2</sup>T. Fukaya et al., "Shifted CholeskyQR for computing the QR factorization of ill-conditioned matrices", Arxiv 2018



# Scalability of Cholesky-QR2

Cholesky-QR2 (CQR2) can achieve superior performance on tall-and-skinny matrices<sup>1</sup>

- Householder QR -  $2mn^2 - \frac{2n^3}{3}$  flops, Cholesky-QR2 -  $4mn^2 + \frac{5n^3}{3}$  flops



CQR2 attains minimal communication cost (by  $\mathcal{O}(\log P)$ ), yet simple implementation

$$T_{\text{Cholesky-QR2}}(m, n, P) = \mathcal{O}\left(\log P \cdot \alpha + n^2 \cdot \beta + \left(\frac{n^2 m}{P} + n^3\right) \cdot \gamma\right)$$

CA-CQR2 parallelizes Cholesky-QR2 over a 3D processor grid, **efficiently factoring any rectangular matrix**

<sup>1</sup>T. Fukaya et al., "CholeskyQR2: A communication-avoiding algorithm", ScalA 2014

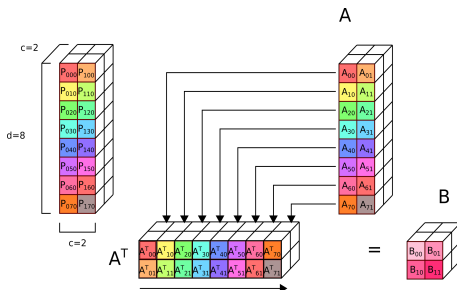
# CA-CQR2's communication-optimal parallelization

CA-CQR2 leverages known 3D algorithms for matrix multiplication<sup>1</sup> and Cholesky factorization<sup>2</sup>

A recursion tree for recursive Cholesky factorization and triangular inverse yields a tradeoff in communication and synchronization<sup>2</sup>

A tunable 3D processor grid of dimensions  $c \times d \times c$  determines the replication factor ( $c$ ), the communication reduction ( $\sqrt{c}$ ), and the number of simultaneous instances of 3D algorithms ( $d/c$ )

Figure: Computation of Gram matrix  $A^T A$



<sup>1</sup>Bersten 1989, "Communication-efficient matrix multiplication on hypercubes", Aggarwal, Chandra, Snir 1990, "Communication complexity of PRAMs", Agarwal et al. 1995, "A three-dimensional approach to parallel matrix multiplication"

<sup>2</sup>A. Tiskin 2007, "Communication-efficient generic pairwise elimination", Future Generation Computer Systems 2007

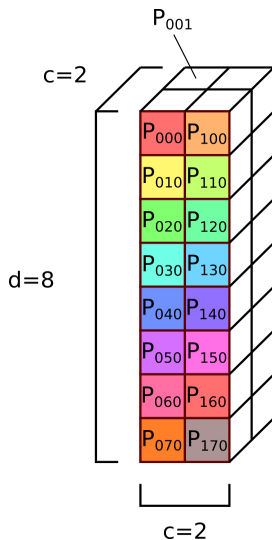
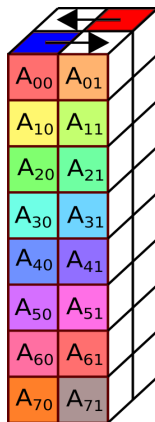
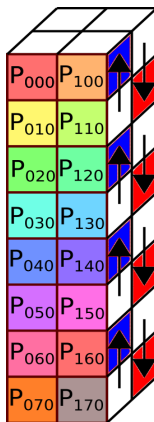
Figure: Start with a tunable  $c \times d \times c$  processor grid

Figure: Broadcast columns of A

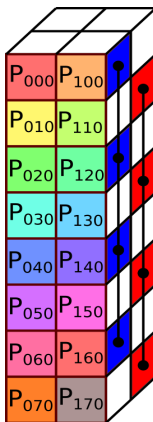


$$\text{Cost: } 2 \log_2 c \cdot \alpha + \frac{2mn}{dc} \cdot \beta$$

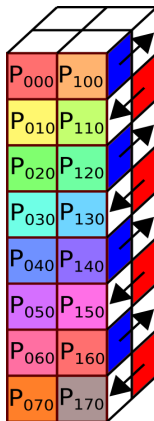
Figure: Reduce contiguous groups of size  $c$ 

$$\text{Cost: } 2 \log_2 c \cdot \alpha + \frac{2n^2}{c^2} \cdot \beta + \frac{n^2}{c^2} \cdot \gamma$$

Figure: Allreduce alternating groups of size  $\frac{d}{c}$

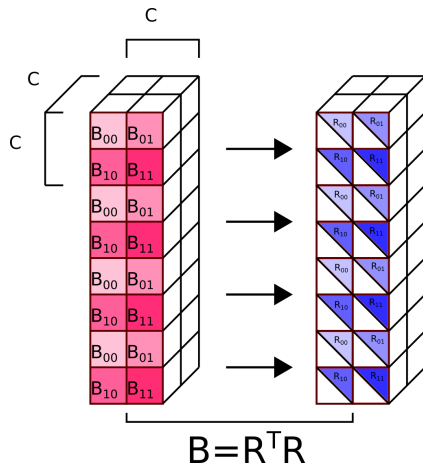


$$\text{Cost: } 2 \log_2 \frac{d}{c} \cdot \alpha + \frac{2n^2}{c^2} \cdot \beta + \frac{n^2}{c^2} \cdot \gamma$$

Figure: Broadcast missing pieces of  $B$  along depth

$$\text{Cost: } 2 \log_2 c \cdot \alpha + \frac{2n^2}{c^2} \cdot \beta$$

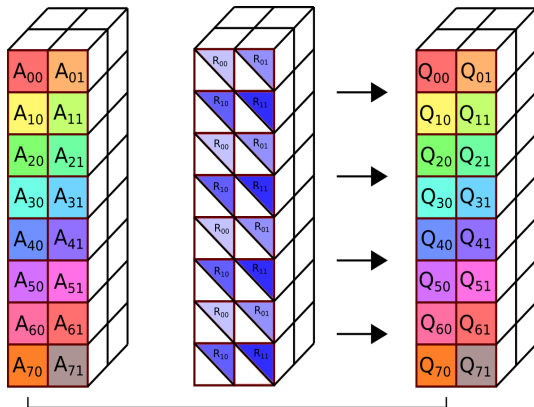
Figure:  $\frac{d}{c}$  simultaneous 3D CholeskyInverse on cubes of dimension  $c$



$$\text{Cost: } \mathcal{O} \left( c^2 \log c^3 \cdot \alpha + \frac{n^2}{2} \cdot \beta + \frac{n^3}{c^3} \cdot \gamma \right)$$



Figure:  $\frac{d}{c}$  simultaneous 3D matrix multiplication or TRSM on cubes of dimension  $c$



$$\text{Cost: } \mathcal{O}(\log_2 c^3 \cdot \alpha + \left(\frac{mn}{dc} + \frac{n^2+nc}{c^2}\right) \cdot \beta + \frac{n^2 m}{c^2 d} \cdot \gamma)$$

CA-CQR2's cost expression expresses tunable tradeoffs

$$T_{\text{CA-CQR2}}^{\alpha-\beta}(m, n, c, d) = \mathcal{O}\left(c^2 \log(d/c) \cdot \alpha + \left(\frac{mn}{dc} + \frac{n^2}{c^2}\right) \cdot \beta + \left(\frac{mn^2}{c^2 d} + \frac{n^3}{c^3}\right) \cdot \gamma\right)$$

Requiring each processor to own a square submatrix ( $\frac{m}{d} = \frac{n}{c}$ ) and enforcing  $P = c^2 d$ , **CA-CQR2 finds an optimal processor grid that supports minimal communication**

	1D Cholesky-QR2	2D ScaLAPACK	2D CAQR	3D CA-CQR2
messages	$\mathcal{O}(\log P)$	$\mathcal{O}(n \log P)$	$\mathcal{O}(\sqrt{P} \log^2 P)$	$\mathcal{O}\left(\left(\frac{Pn}{m}\right)^{\frac{2}{3}} \log P\right)$
words	$\mathcal{O}(n^2)$	$\mathcal{O}\left(\frac{mn}{\sqrt{P}}\right)$	$\mathcal{O}\left(\frac{mn}{\sqrt{P}}\right)$	$\mathcal{O}\left(\left(\frac{n^2 m}{P}\right)^{\frac{2}{3}}\right)$
flops	$\mathcal{O}\left(\frac{n^2 m}{P} + n^3\right)$	$\mathcal{O}\left(\frac{mn^2}{P}\right)$	$\mathcal{O}\left(\frac{mn^2}{P}\right)$	$\mathcal{O}\left(\frac{n^2 m}{P}\right)$
memory	$\mathcal{O}\left(\frac{mn}{P} + n^2\right)$	$\mathcal{O}\left(\frac{mn}{P}\right)$	$\mathcal{O}\left(\frac{mn}{P}\right)$	$\mathcal{O}\left(\left(\frac{n^2 m}{P}\right)^{\frac{2}{3}}\right)$

Minimal communication cost in a QR factorization is reflected by the surface area of the cubic volume of  $\mathcal{O}(mn^2/P)$  computation

We factor  $m \times n$  matrices with  $m \gg n$  to highlight the effect CA-CQR2's communication reduction and algorithmic tradeoffs have on performance



Scaling studies highlight **interplay between CA-CQR2's increased arithmetic intensity and an architecture's machine balance**

- ratio of peak-flops to network bandwidth is 8x higher on Stampede2<sup>1</sup> than BlueWaters<sup>2</sup>

We show only the **most-performant variants at each node count** of CA-CQR2 and ScaLAPACK's PGEQRF

- ScaLAPACK tuned over 2D processor grid dimensions and block sizes
- CA-CQR2 tuned over processor grid dimensions  $d$  and  $c$
- each tested/tuned over a number of resource configurations
- both algorithms use Householder's flop cost in determining performance

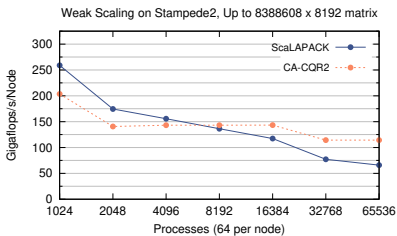
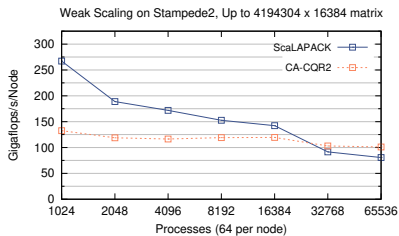
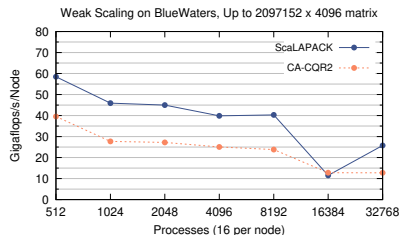
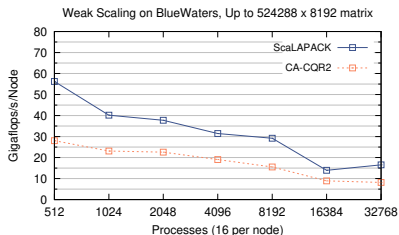
---

<sup>1</sup>Intel Knights Landing (KNL) cluster at TACC

<sup>2</sup>Cray XE/XK hybrid machine at NCSA

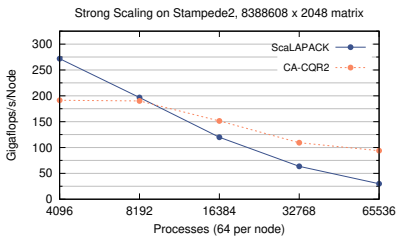
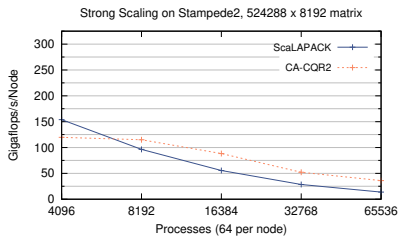
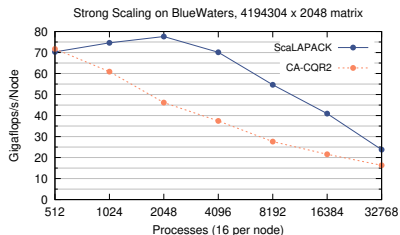
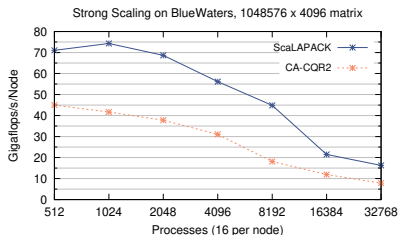
# QR Weak scaling on Stampede2 and Blue Waters

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count



# QR Strong scaling on Stampede2 and Blue Waters

Figure: Strong scaling for  $m \times n$  matrices



CA-CQR2's performance improvements over ScaLAPACK on Stampede2 range from **1.1 - 3.3x at 1024 nodes**

## CA-CQR2 leverages current and future architectural trends

- machines with highest ratio of peak node performance to peak injection bandwidth will benefit most
- asymptotic communication reduction increasingly evident as we scale, despite overheads in synchronization and computation

These results motivate increasingly wide overdetermined systems, a **critical use case for solving linear least squares and eigenvalue problems**

Our study shows that **communication-optimal parallel QR factorizations can achieve superior performance and scaling up to thousands of nodes**<sup>1 2</sup>

---

<sup>1</sup>Our preprint detailing CA-CQR2 can be found at <https://arxiv.org/abs/1710.08471>

<sup>2</sup>Our C++ implementation can be found at <https://github.com/huttered40/CA-CQR2>

I'd like to acknowledge the **Department of Energy** and **Krell Institute** for supporting this research via awarding me a **DOE Computational Science Graduate Fellowship**<sup>1</sup>

We'd also like to acknowledge a number of computing centers for providing benchmarking resources

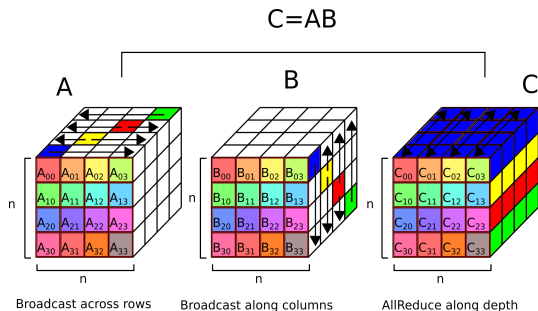
- Texas Advanced Computing Center (TACC) via Stampede<sup>2</sup>
- National Center for Supercomputing Applications (NCSA) via Blue Waters<sup>3</sup>
- Argonne Leadership Computing Facility (Cetus, Mira, Theta) for preliminary benchmarking

---

<sup>1</sup>Grant number DE-SC0019323

<sup>2</sup>Allocation TG-CCR180006

<sup>3</sup>Awards OCI-0725070 and ACI-1238993

Figure: 3D algorithm for square matrix multiplication <sup>1 2 3</sup>

$$T_{3D\_MM}(n, P) = \mathcal{O}\left(\log P \cdot \alpha + \frac{n^2}{P^{\frac{2}{3}}} \cdot \beta + \frac{n^3}{P} \cdot \gamma\right)$$

<sup>1</sup>Bersten 1989, "Communication-efficient matrix multiplication on hypercubes"

<sup>2</sup>Aggarwal, Chandra, Snir 1990, "Communication complexity of PRAMs"

<sup>3</sup>Agarwal et al. 1995, "A three-dimensional approach to parallel matrix multiplication"



We can embed the recursive definitions of Cholesky factorization and triangular inverse to find matrices  $R, R^{-1}$

Tuning the recursion tree yields a tradeoff in horizontal bandwidth and synchronization<sup>1</sup>

---


$$[L, L^{-1}] \leftarrow \text{CholeskyInverse}(A)$$


---

$$\begin{bmatrix} L_{11} & L_{11}^{-1} \end{bmatrix} \leftarrow \text{CholeskyInverse}(A_{11})$$

$$L_{21} \leftarrow A_{21} L_{11}^{-T}$$

$$\begin{bmatrix} L_{22} & L_{22}^{-1} \end{bmatrix} \leftarrow \text{CholeskyInverse}(A_{22} - L_{21} L_{21}^T)$$

$$L_{21}^{-1} \leftarrow -L_{22}^{-1} L_{21} L_{11}^{-1}$$


---

$$T_{\text{CholeskyInverse3D}}(n, P) = \mathcal{O} \left( P^{\frac{2}{3}} \log P \cdot \alpha + \frac{n^2}{P^{\frac{2}{3}}} \cdot \beta + \frac{n^3}{P} \cdot \gamma \right)$$

$$T_{\text{ScaLAPACK}}(n, P) = \mathcal{O} \left( \sqrt{P} \log P \cdot \alpha + \frac{n^2}{\sqrt{P}} \cdot \beta + \frac{n^3}{P} \cdot \gamma \right)$$

---

<sup>1</sup>A. Tiskin 2007, "Communication-efficient generic pairwise elimination"

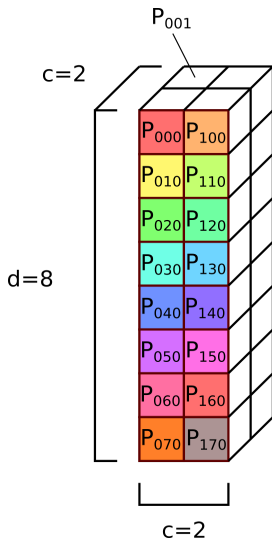
Figure: Start with a tunable  $c \times d \times c$  processor grid

Figure: Compute Gram matrix  $A^T A$

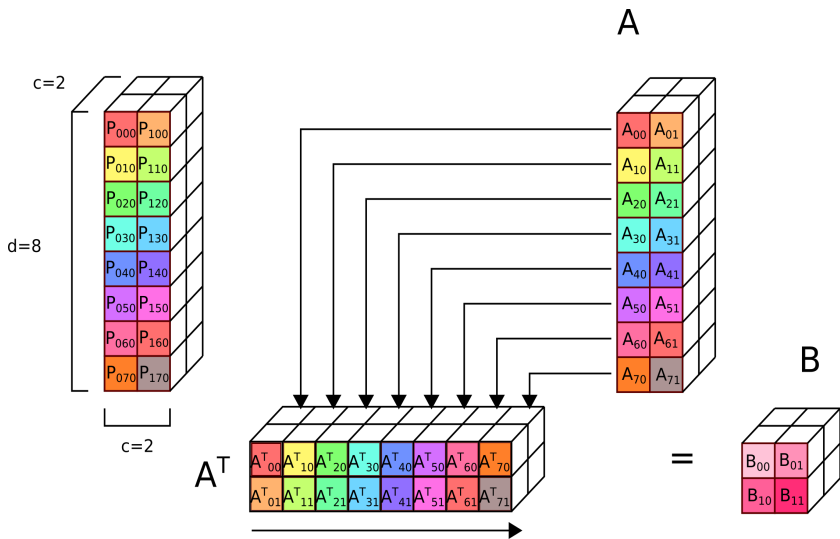


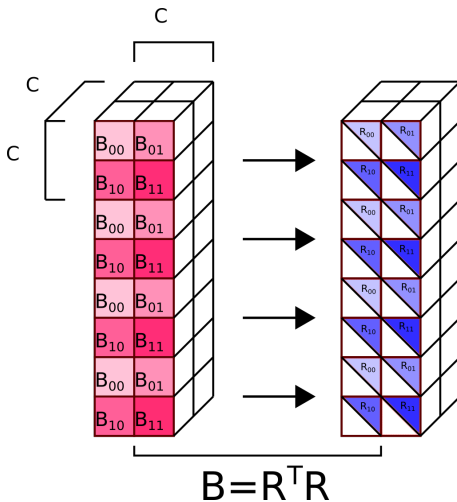
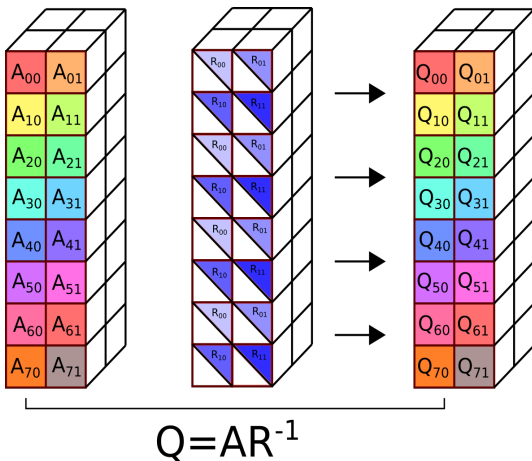
Figure: Compute 3D CholeskyInverse on processor grid cubes of dimension  $c$ 

Figure: Compute 3D matrix multiplication or 3D TRSM on processor grid cubes of dimension  $c$



# Optimum cost of CholesyQR2-Tunable

The advantage of using a tunable grid lies in the ability to frame the shape of the grid around the shape of rectangular  $m \times n$  matrix  $A$ . Optimal communication can be attained by ensuring that the grid perfectly fits the dimensions of  $A$ , or that the dimensions of the grid are proportional to the dimensions of the matrix. We derive the cost for the optimal ratio  $\frac{m}{d} = \frac{n}{c}$  below. Using equation  $P = c^2 d$  and

$\frac{m}{d} = \frac{n}{c}$ , solve for  $d, c$  in terms of  $m, n, P$ . Solving the system of equations yields  $c = \left(\frac{Pn}{m}\right)^{\frac{1}{3}}$ ,  $d = \left(\frac{Pm^2}{n^2}\right)^{\frac{1}{3}}$ . We can plug these values into the cost of Cholesky-QR2-Tunable to find the optimal cost.

$$\begin{aligned}
 T_{\text{Cholesky-QR2-Tunable}}^{\alpha-\beta} \left( m, n, \left(\frac{Pn}{m}\right)^{\frac{1}{3}}, \left(\frac{Pm^2}{n^2}\right)^{\frac{1}{3}} \right) &= \mathcal{O} \left( \left(\frac{Pn}{m}\right)^{\frac{2}{3}} \log P \cdot \alpha \right. \\
 &+ \frac{\left(\frac{Pn}{m}\right)^{\frac{1}{3}} mn + n^2 \left(\frac{Pm^2}{n^2}\right)^{\frac{1}{3}}}{\left(\frac{Pm^2}{n^2}\right)^{\frac{1}{3}} \left(\frac{Pn}{m}\right)^{\frac{2}{3}}} \cdot \beta + \frac{n^3 \left(\frac{Pm^2}{n^2}\right)^{\frac{1}{3}} + n^2 m \left(\frac{Pn}{m}\right)^{\frac{1}{3}}}{\left(\frac{Pn}{m}\right) \left(\frac{Pm^2}{n^2}\right)^{\frac{1}{3}}} \cdot \gamma \Big) \\
 &= \mathcal{O} \left( \left(\frac{Pn}{m}\right)^{\frac{2}{3}} \log P \cdot \alpha + \left(\frac{n^2 m}{P}\right)^{\frac{2}{3}} \cdot \beta + \frac{n^2 m}{P} \cdot \gamma \right)
 \end{aligned} \tag{1}$$

Grid shape	Metric	Cost
optimal	# of messages	$\mathcal{O} \left( \left(\frac{Pn}{m}\right)^{\frac{2}{3}} \log P \right)$
	# of words	$\mathcal{O} \left( \left(\frac{n^2 m}{P}\right)^{\frac{2}{3}} \right)$
	# of flops	$\mathcal{O} \left( \frac{n^2 m}{P} \right)$
	Memory footprint	$\mathcal{O} \left( \left(\frac{n^2 m}{P}\right)^{\frac{2}{3}} \right)$

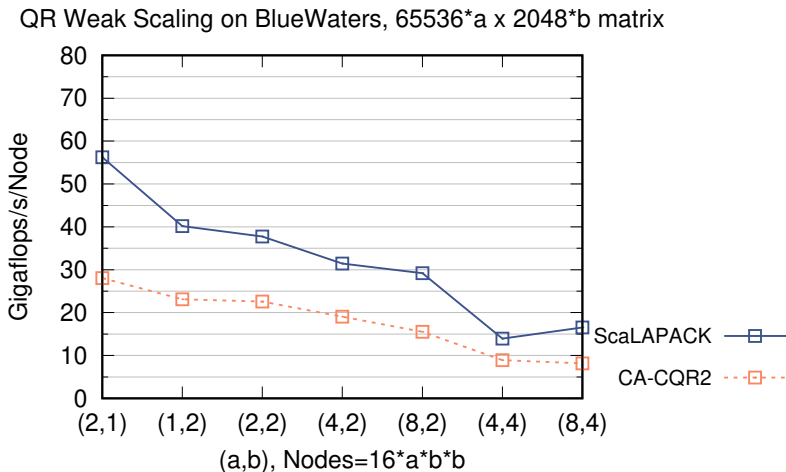


Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

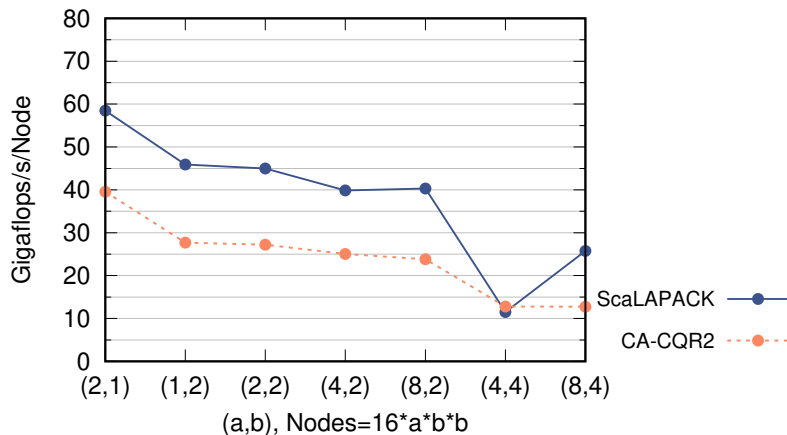
QR Weak Scaling on BlueWaters,  $262144 \times a \times 1024 \times b$  matrix

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count



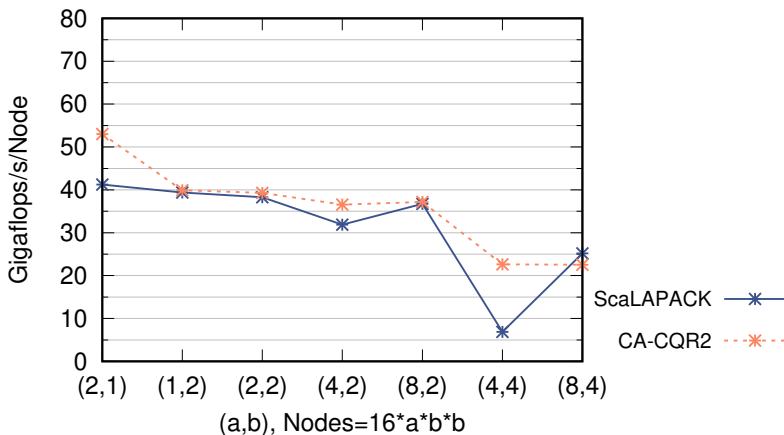
QR Weak Scaling on BlueWaters,  $1048576 \times a \times 512 \times b$  matrix

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

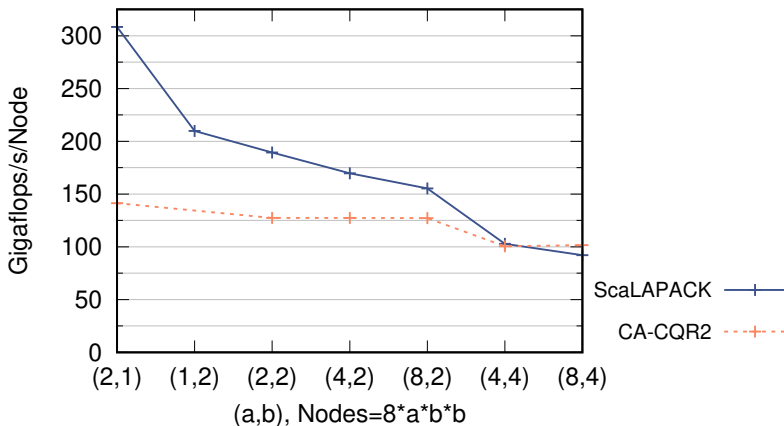
QR Weak Scaling on Stampede2,  $262144^*a \times 8192^*b$  matrix

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

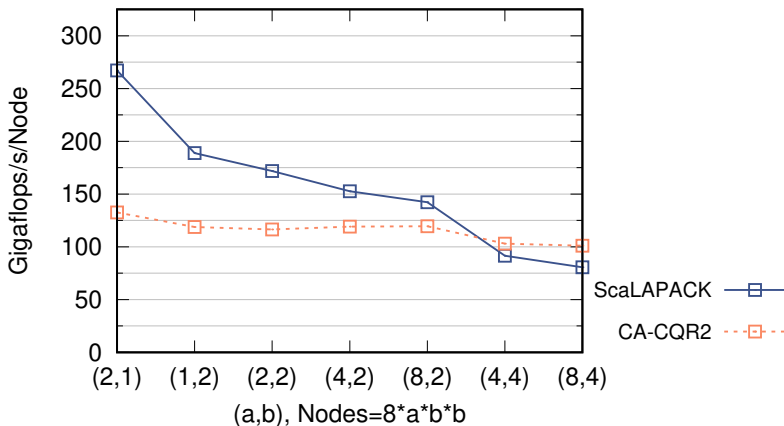
QR Weak Scaling on Stampede2,  $524288 \times a \times 4096 \times b$  matrix

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

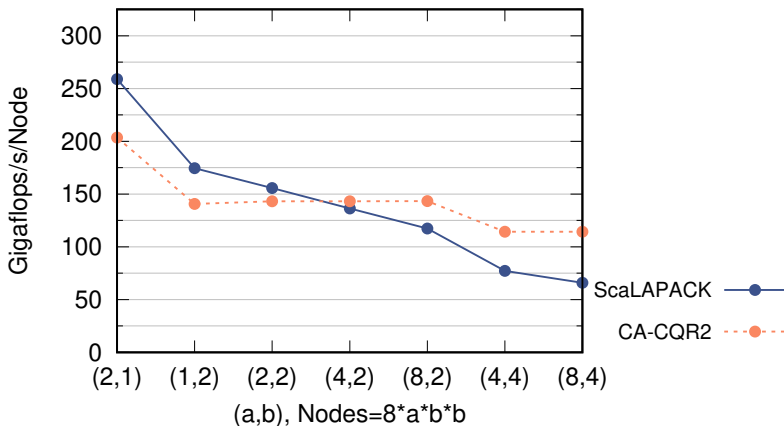
QR Weak Scaling on Stampede2,  $1048576 \times a \times 2048 \times b$  matrix

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

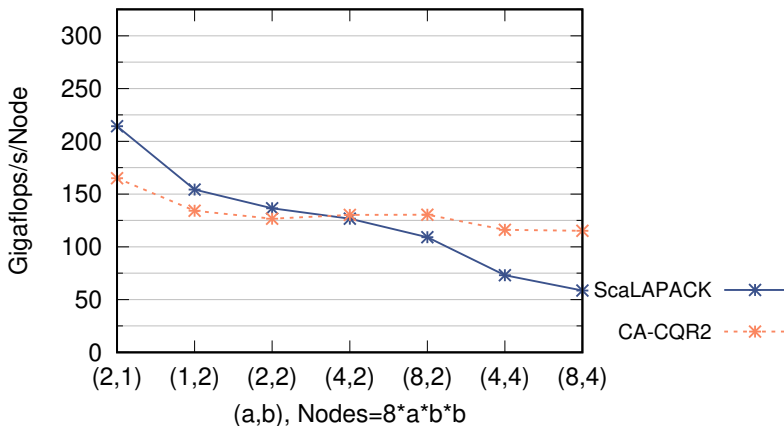
QR Weak Scaling on Stampede2,  $2097152^*a \times 1024^*b$  matrix

Figure: Weak scaling for  $m \times n$  matrices so  $mn^2$  scales linearly with node count

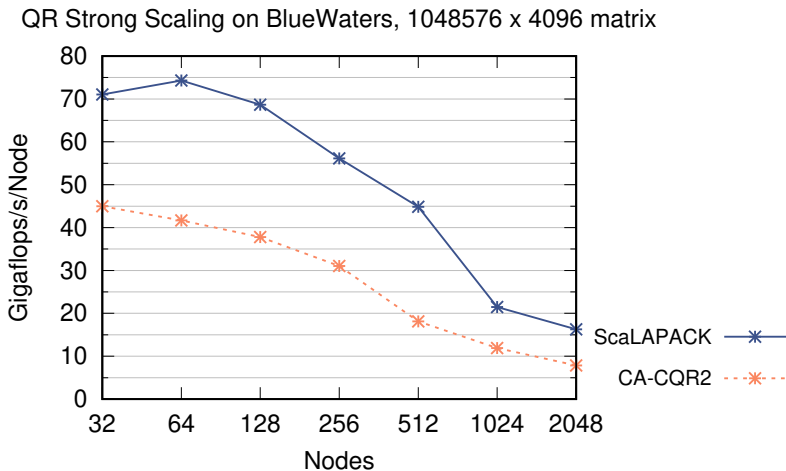


Figure: Strong scaling for QR factorization

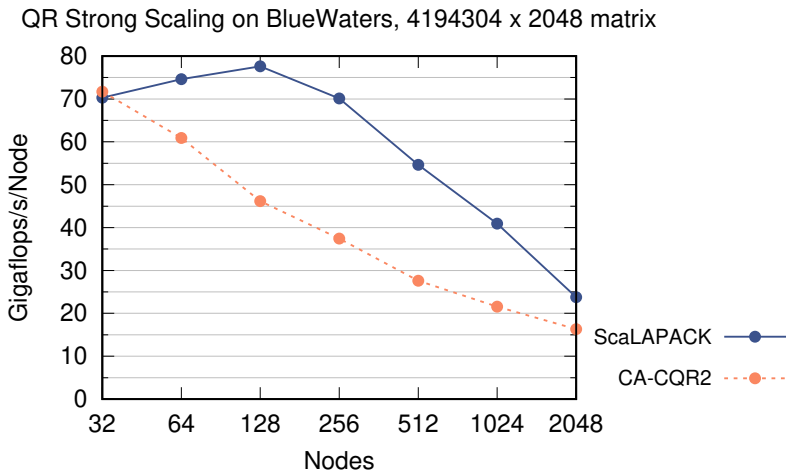


Figure: Strong scaling for QR factorization

QR Strong Scaling on Stampede2, 524288 x 8192 matrix

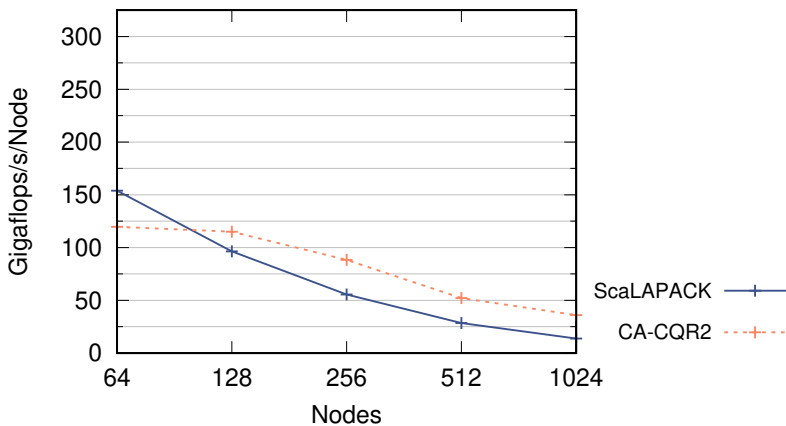


Figure: Strong scaling for QR factorization



QR Strong Scaling on Stampede2, 2048576 x 4096 matrix

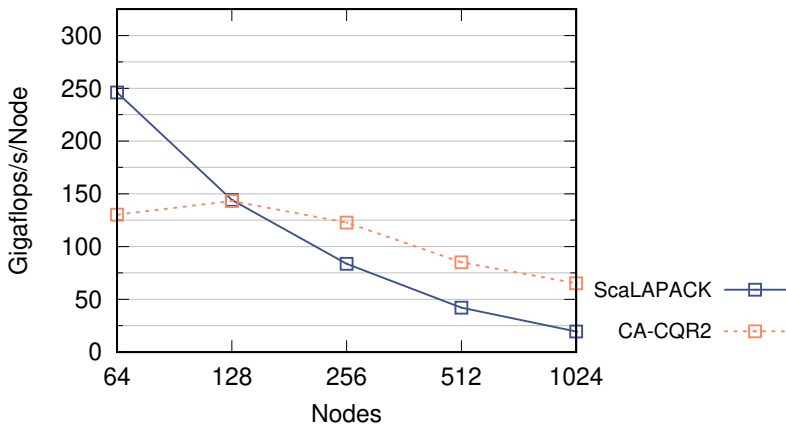


Figure: Strong scaling for QR factorization

QR Strong Scaling on Stampede2, 8388608 x 2048 matrix

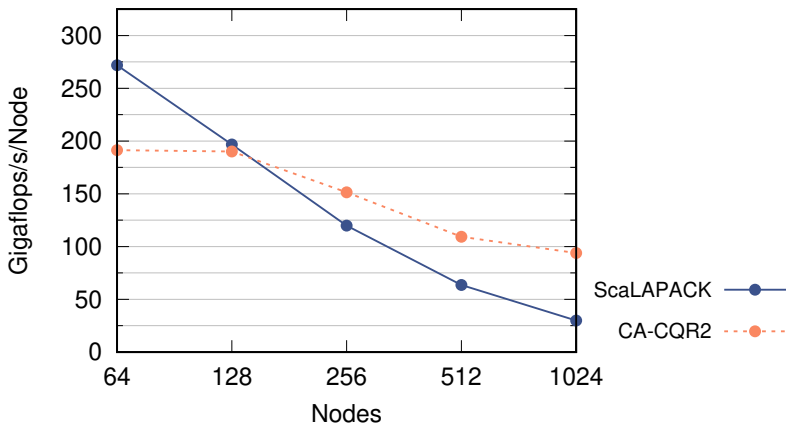


Figure: Strong scaling for QR factorization

QR Strong Scaling on Stampede2, 33554432 x 1024 matrix

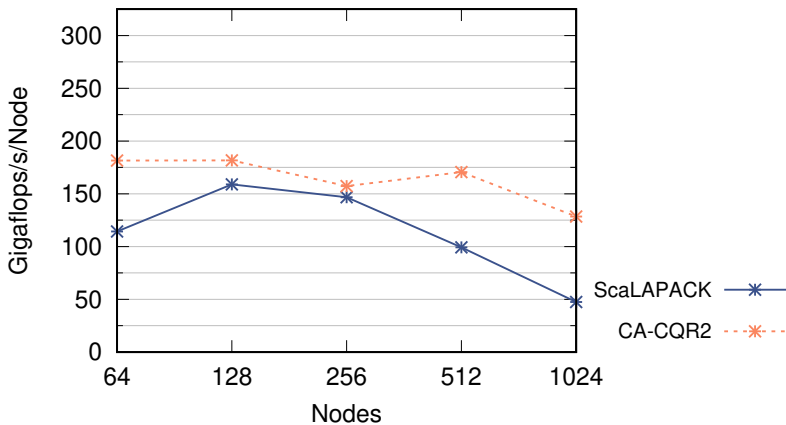


Figure: Strong scaling for QR factorization

Weak Scaling,  $1048576^*a \times 512^*b$

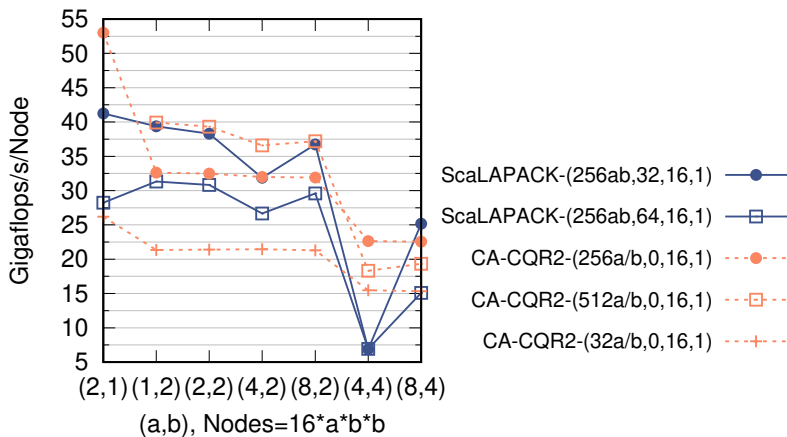


Figure: Weak scaling for matrices with dimensions given in legend

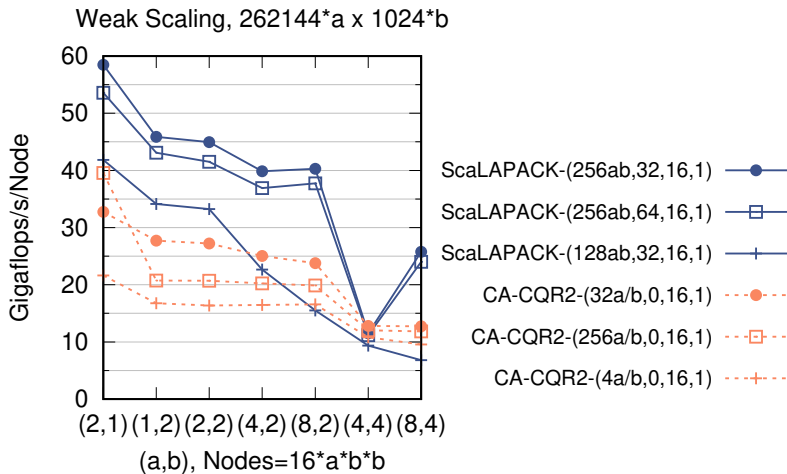


Figure: Weak scaling for matrices with dimensions given in legend

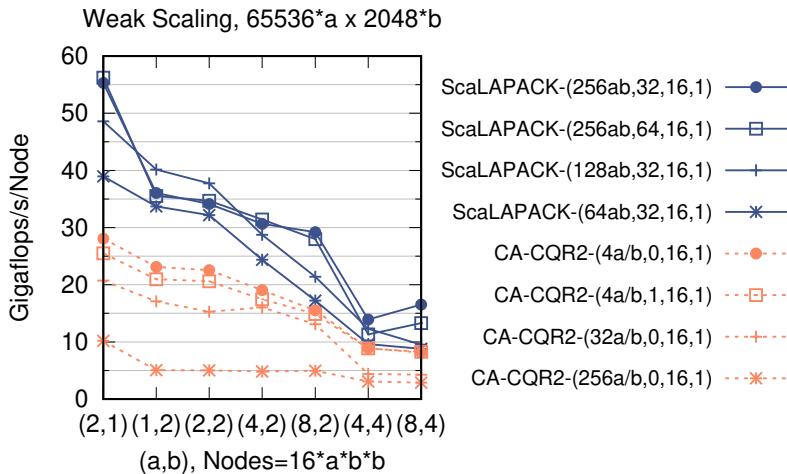


Figure: Weak scaling for matrices with dimensions given in legend

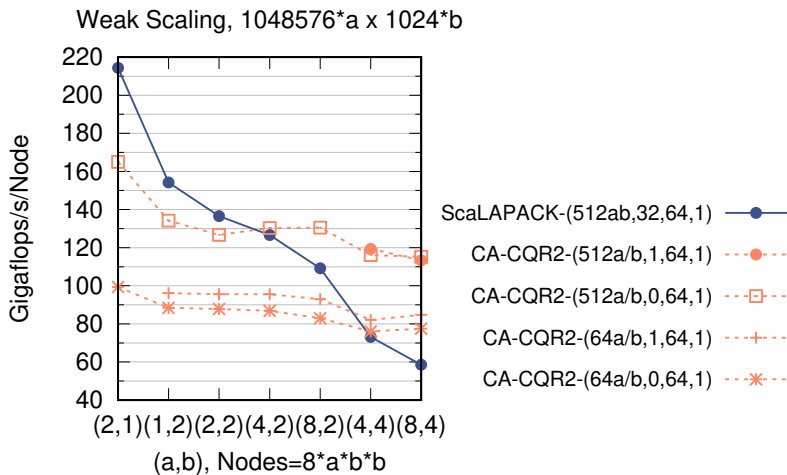


Figure: Weak scaling for matrices with dimensions given in legend

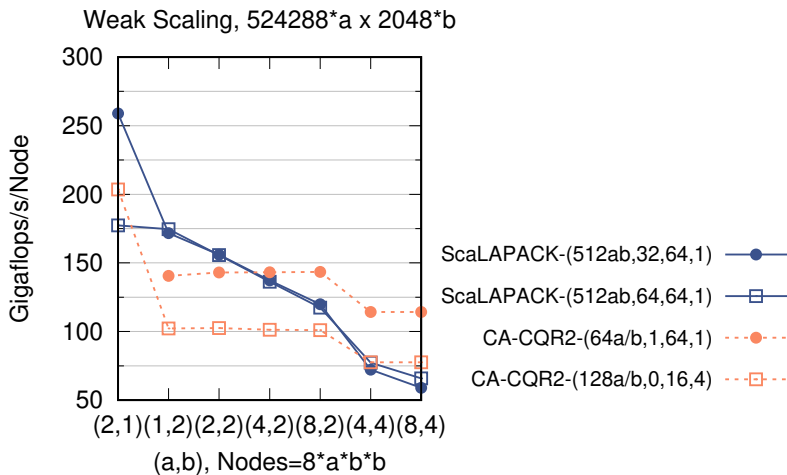


Figure: Weak scaling for matrices with dimensions given in legend



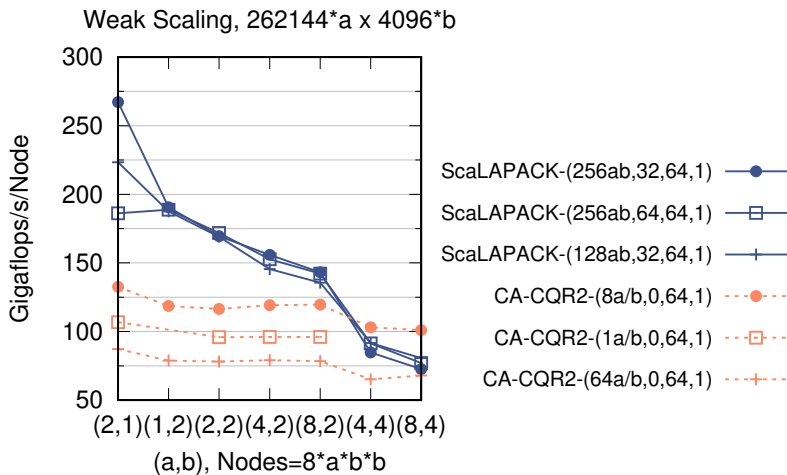


Figure: Weak scaling for matrices with dimensions given in legend

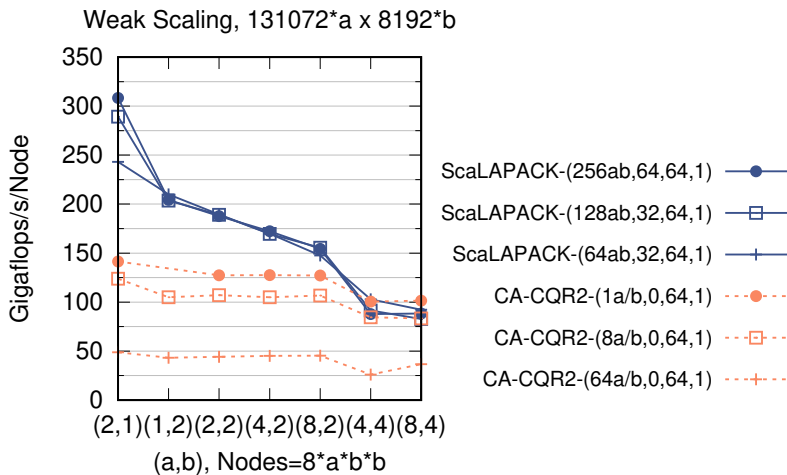


Figure: Weak scaling for matrices with dimensions given in legend

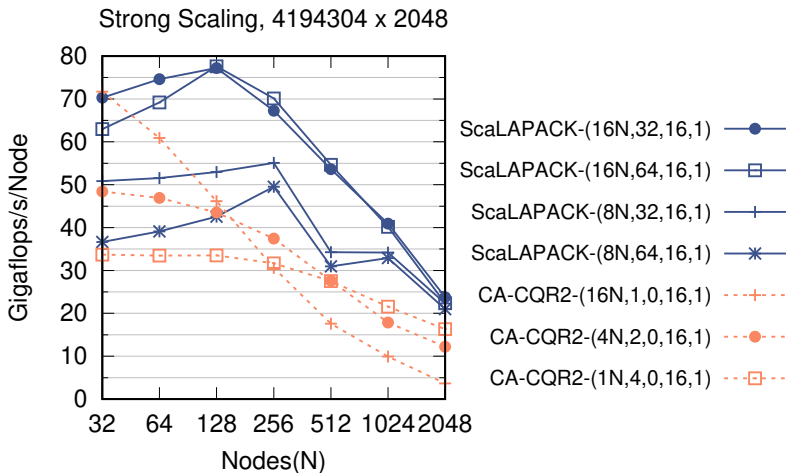


Figure: Strong scaling for matrices with dimensions given in legend

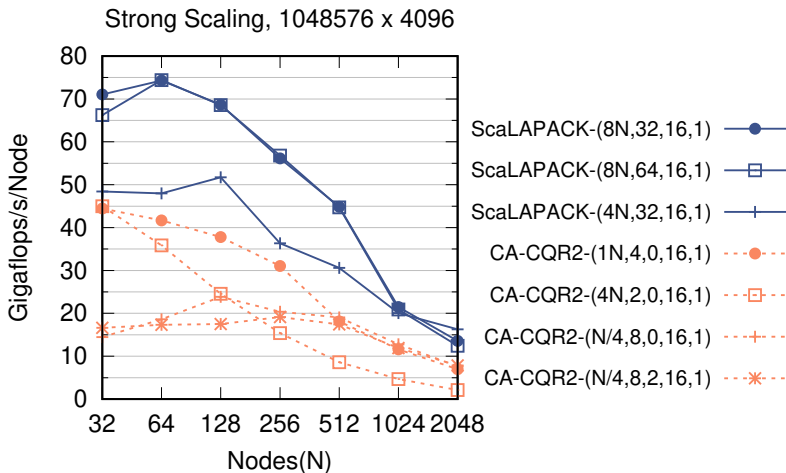


Figure: Strong scaling for matrices with dimensions given in legend

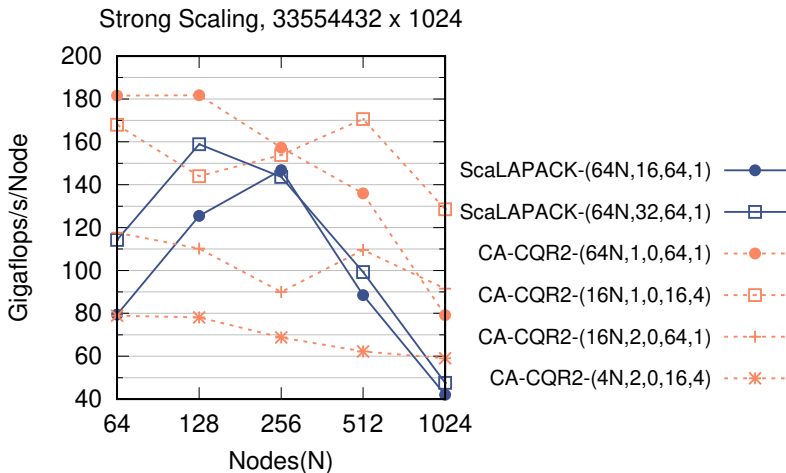


Figure: Strong scaling for matrices with dimensions given in legend

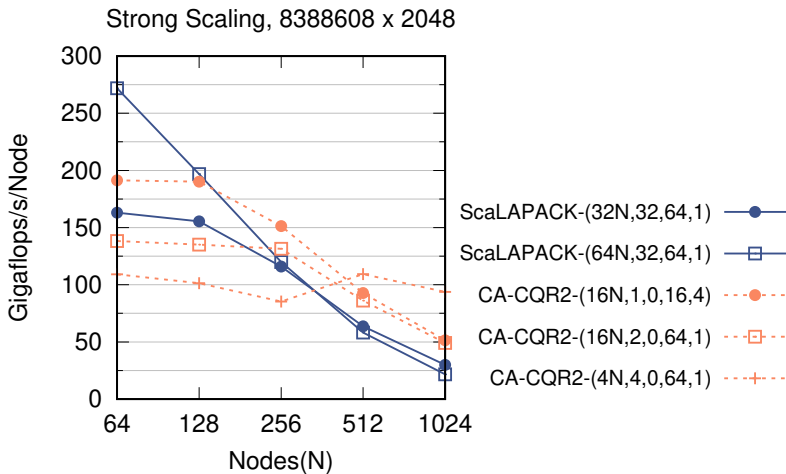


Figure: Strong scaling for matrices with dimensions given in legend

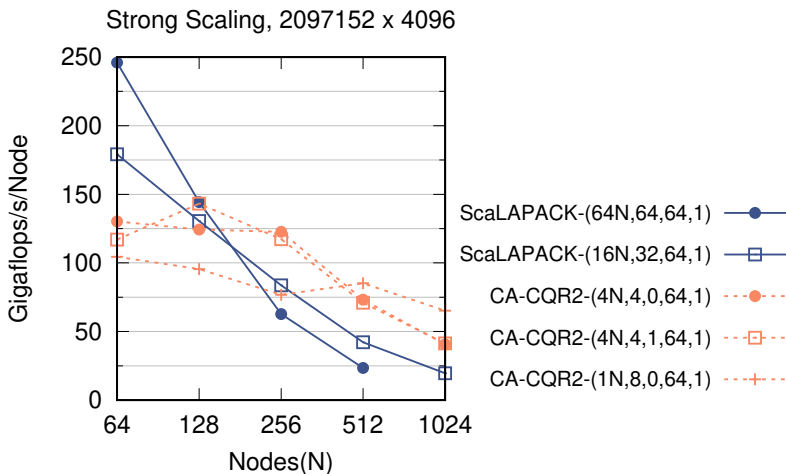


Figure: Strong scaling for matrices with dimensions given in legend

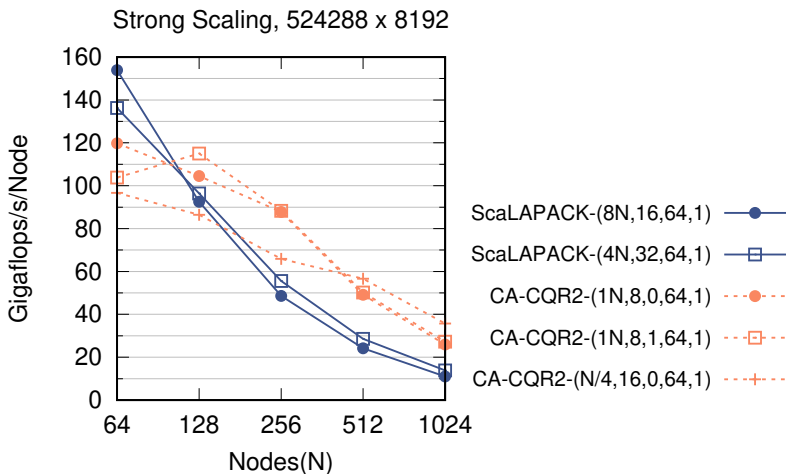


Figure: Strong scaling for matrices with dimensions given in legend