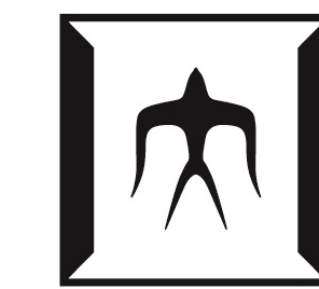
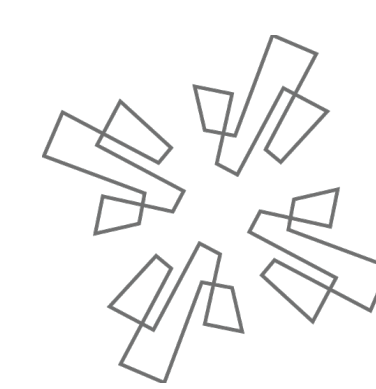


Randomized SVD on TensorCores



Tokyo Tech

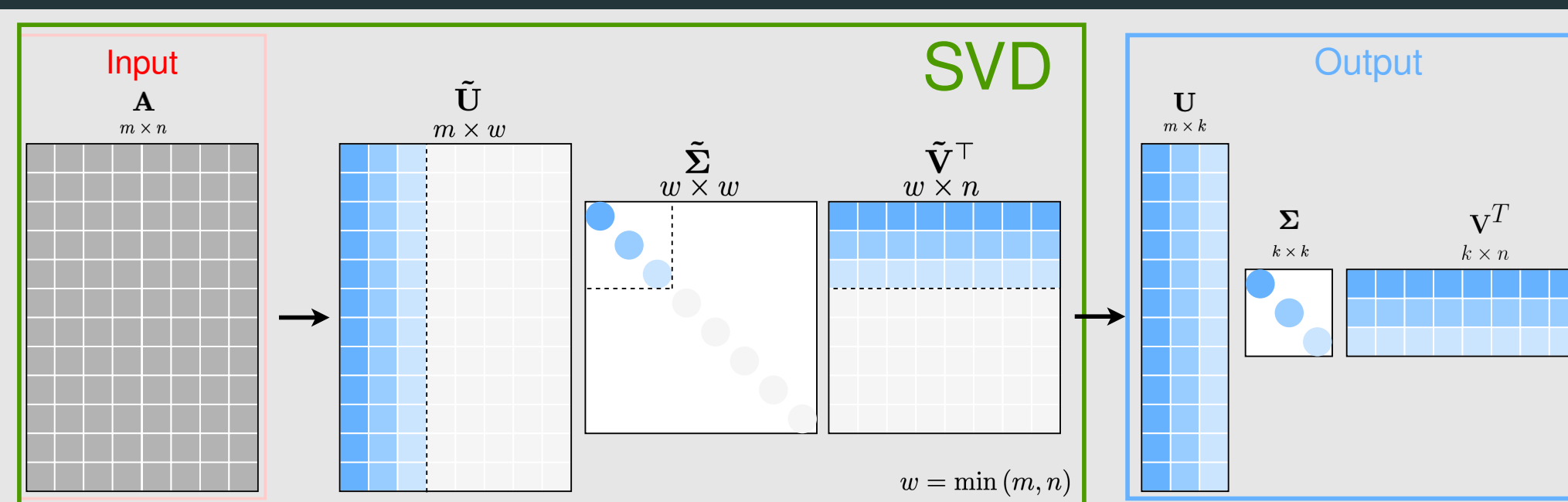
Hiroyuki Ootomo¹, Rio Yokota²

¹School of Computing, TokyoTech ²Global Scientific Information and Computing Center, TokyoTech

Abstract

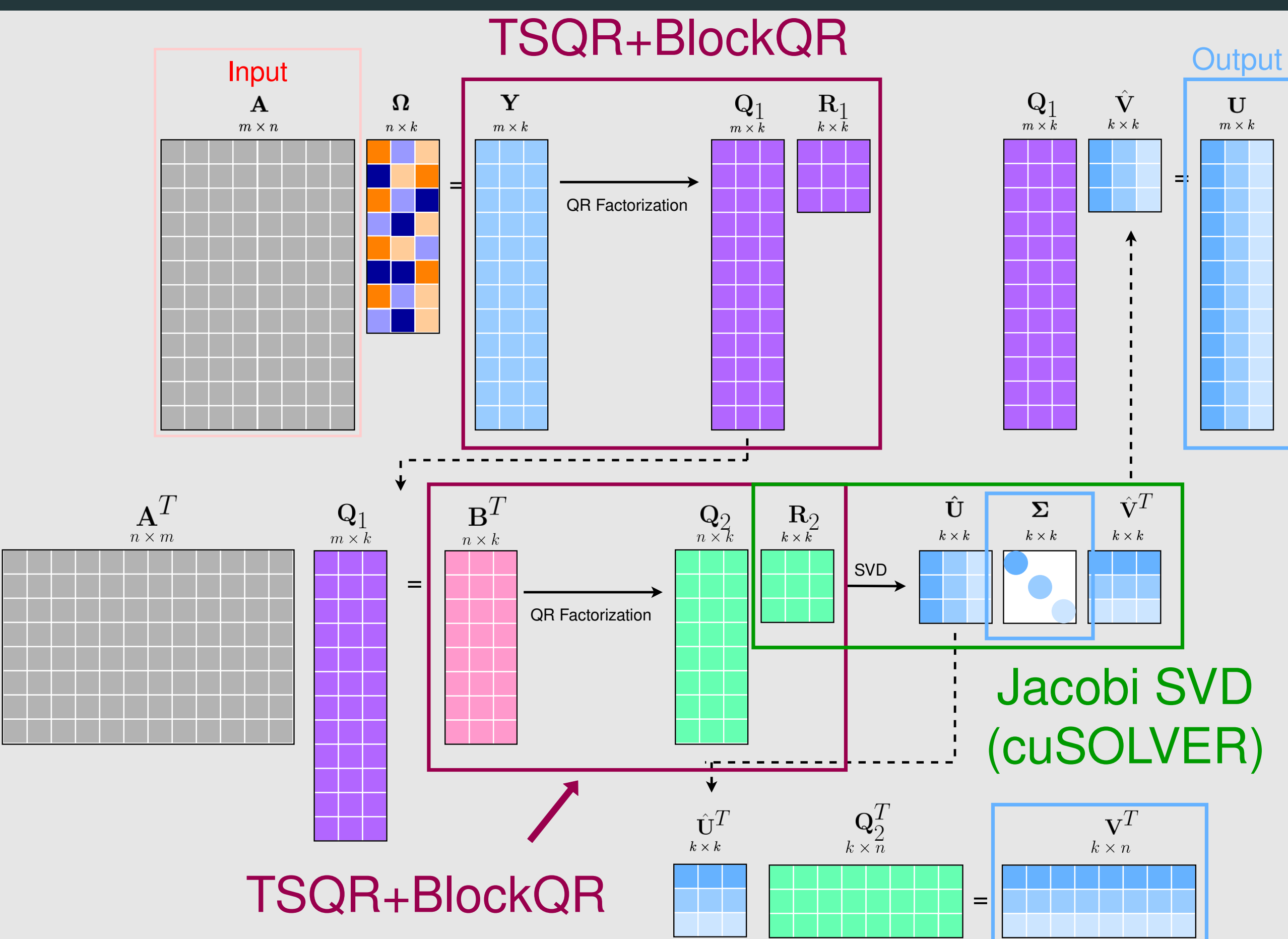
Low-rank approximation is vital and widely used for data compression, dimensionality reduction and noise reduction. Randomized SVD is an efficient algorithm for computing low-rank approximation. This algorithm uses a QR factorization of a tall skinny matrix. We evaluate the speed, accuracy, and stability of Randomized SVD on TensorCores.

Low-rank approximation using SVD



1. $\tilde{U}, \tilde{\Sigma}, \tilde{V}^T \stackrel{SVD}{\leftarrow} A$.
2. Reduce unnecessary singular values and singular vectors.

Randomized SVD



Randomized SVD computes low-rank approximation with specified rank k .

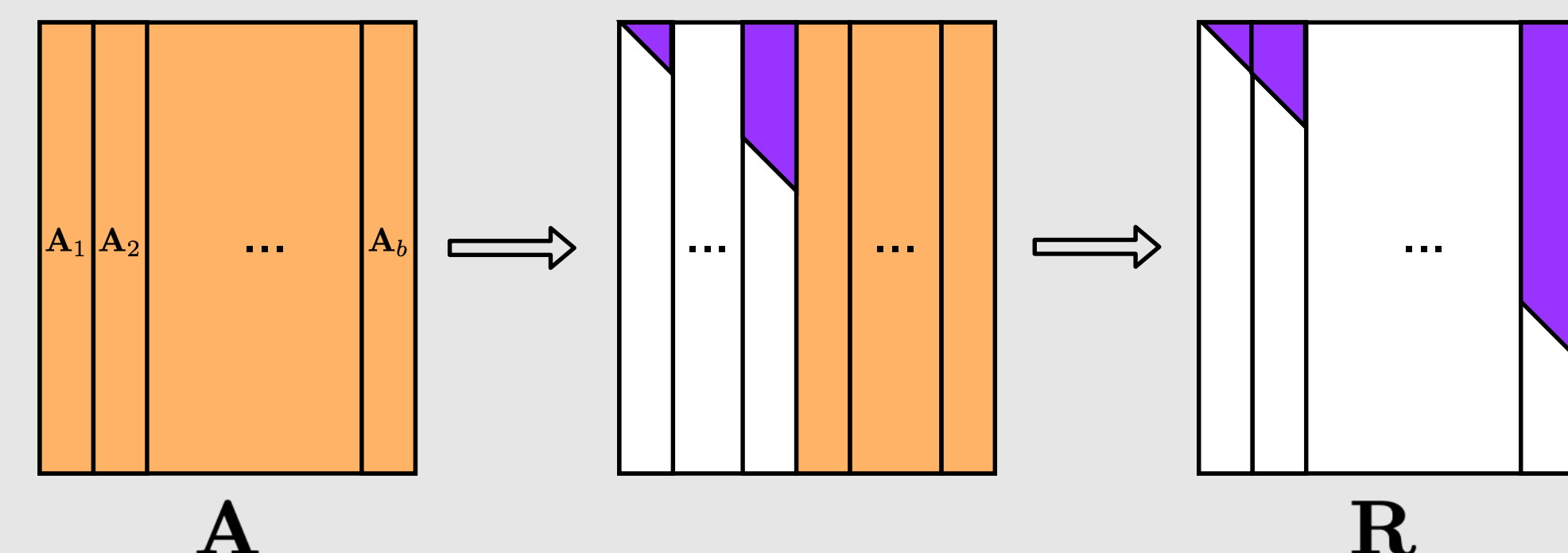
Algorithm

1. $Y \leftarrow A\Omega$
2. $Q_1 R_1 \leftarrow QR(Y)$
3. $B^T \leftarrow A^T Q_1$
4. $Q_2 R_2 \leftarrow QR(B^T)$
5. $\hat{U} \hat{\Sigma} \hat{V}^T \leftarrow SVD(R_2)$
6. $U \leftarrow Q_1 \hat{U}$
7. $V \leftarrow Q_2 \hat{U}$

QR factorization

- ▶ TSQR + BlockQR
- ▶ cuSOLVER

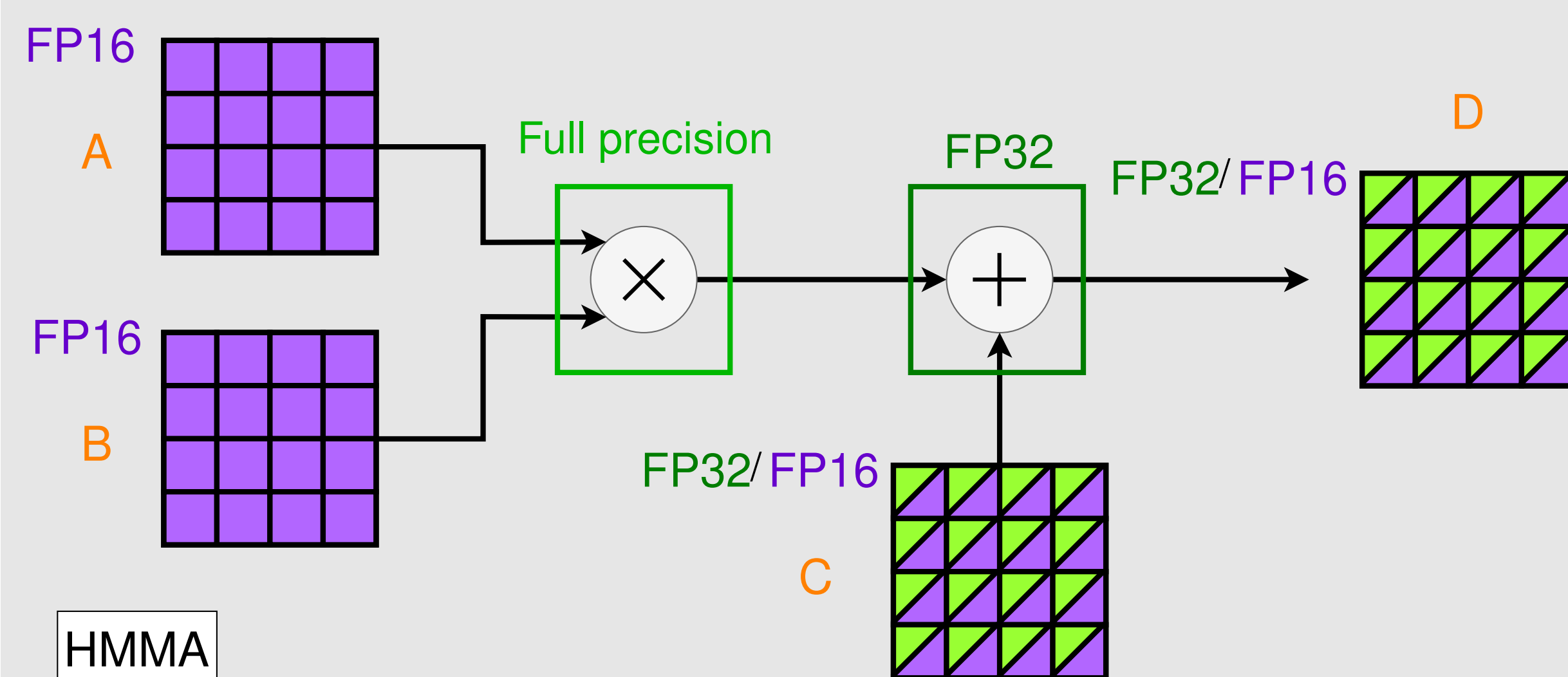
BlockQR



1. Divide input matrix A into A_1, A_2 by column.
2. Compute $Q_1, R_1 \leftarrow QR(A_1)$
3. Merge matrices. $Q = \begin{bmatrix} Q_1 & Q_2 \\ 0 & R_2 \end{bmatrix}$
4. Compute this two divided QR recursively. $R = \begin{bmatrix} R_1 & R_{11} \\ 0 & R_2 \end{bmatrix}$

FP32 matmul with accuracy correction on TensorCores

TensorCore



▶ Mixed precision matrix product and addition circuit
Accuracy correction of FP32 matrix product

$$A_{FP16} \leftarrow A_{FP16}, B_{FP16} \leftarrow B_{FP16}$$

- ▶ Without correction

$$C_{FP32} \leftarrow A_{FP16} B_{FP16}$$

- ▶ With correction

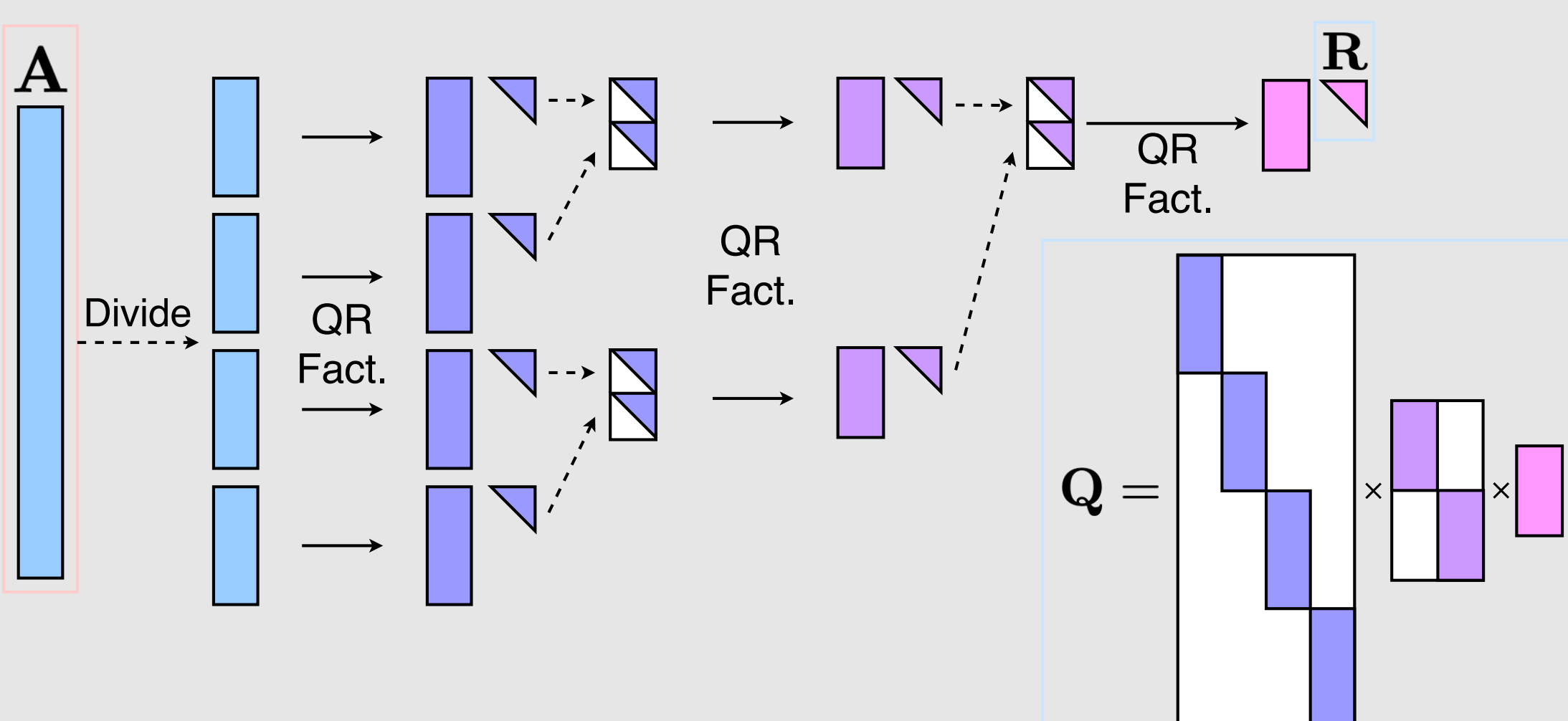
$$C_{FP32} \leftarrow A_{FP16} B_{FP16} + \underbrace{\Delta A_{FP16} B_{FP16} + A_{FP16} \Delta B_{FP16}}_{\text{Correction term}}$$

$$\text{where } \Delta M_{FP16} \leftarrow \text{Convert to FP16 } M_{FP32} - M_{FP16}$$

Test environment

- ▶ Intel Xeon CPU E5-2630 v3 @ 2.40GHz x2
- ▶ NVIDIA Tesla V100-PCIE-16GB
- ▶ CUDA 10.1

TSQR on TensorCore^[1]



Algo 1. TSQR

1. Divide the input matrix A .
2. Calculate QR decomposition for each subdivided matrix to get R s and Q s.
3. Merge consecutive R blocks.
4. Repeat 2-3 until there is only one R .
5. Calculate Q from Q s which are calculated in 2-4.

Algo 2. Householder QR

Require: $m, n \in \mathbb{N}, A \in \mathbb{R}^{m \times n}$

Ensure: $Q \in \mathbb{R}^{m \times m}, R \in \mathbb{R}^{m \times n}$

1. $Q' \leftarrow I, R \leftarrow A$
2. for $i \leftarrow 0$ to $n - 1$ do
3. $u \leftarrow [0 \cdots 0 \ R_{i,i} \cdots R_{m-1,i}]^T$
4. $u_i \leftarrow u_i \pm |u_i|$
5. $H \leftarrow I - 2 \frac{uu^T}{|u|^2}$
6. $R \leftarrow HR$
7. $Q' \leftarrow HQ'$
8. end for
9. $Q \leftarrow Q'^T$

Implementation

- ▶ **Step 2-4**
Some QR decompositions can be calculated in parallel.
⇒ Batched QR
- ▶ **Step 5**
Calculate matrix multiplication implicitly using batched matmul.

Batched QR

To compute each QR factorization, we use HouseholderQR.

Where to use TensorCores

- ▶ Calculate H (Algo 2. line 5) ▶ Batched Matmul
- ▶ Update Q, R (Algo 2. line 6, 7)

[1] Hiroyuki Ootomo, Rio Yokota "TSQR on TensorCores", SC'19 Research Poster

Implementation

- ▶ Input/Output type is FP32
 - ▶ SVD operations using gesvdj subroutine of cuSOLVER
 - ▶ Implement TSQR with different configurations to be used within the QR factorizations (BlockQR) of Randomized SVD
- | | TensorCore | Correction |
|----------|------------|------------|
| TC-noCor | YES | NO |
| TC-Cor | YES | YES |
| noTC | NO | NO |
- ▶ Implement Randomized SVD using QR factorizations subroutine of cuSOLVER for comparison

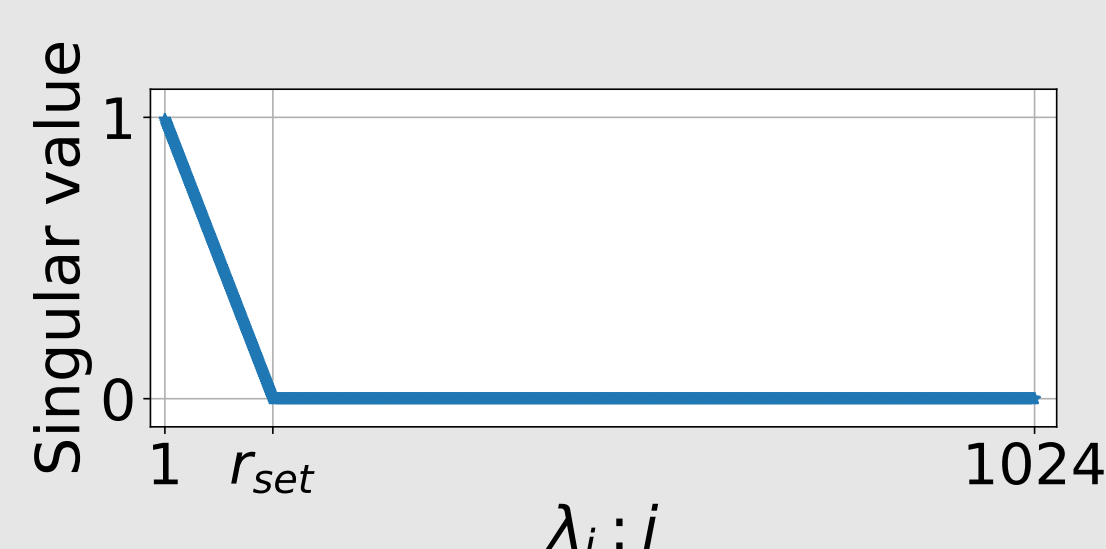
Evaluation of Randomized SVD on TensorCores

Input matrix

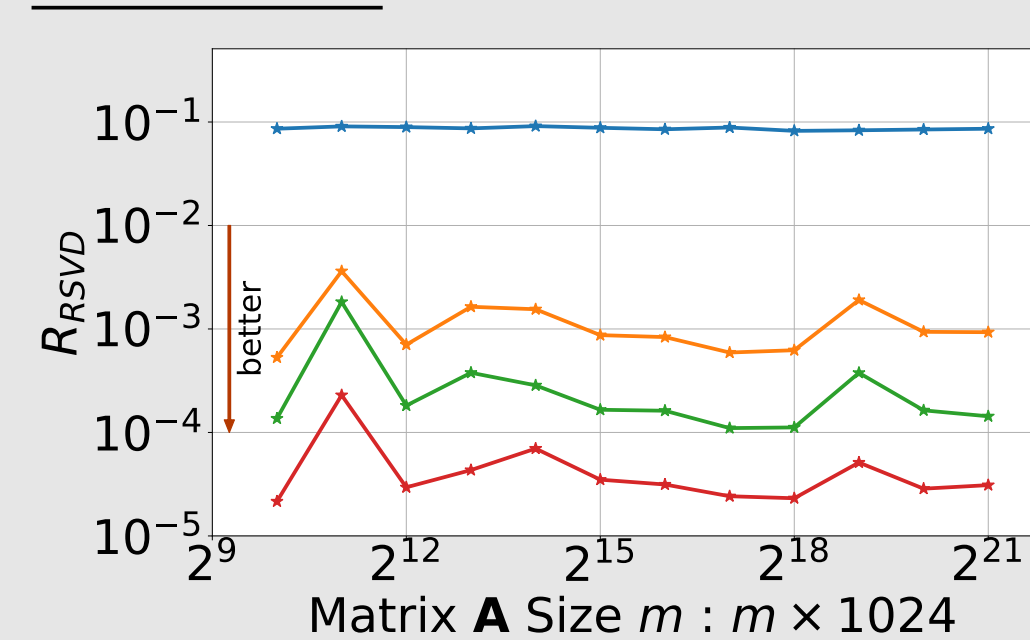
Random $m \times 1024$ FP32 matrix A whose singular values are calculated by eq (1).

$$\lambda_i = (1 - \alpha) \max \left(1 - \frac{i-1}{r_{set}}, 0 \right) + \alpha \quad (1)$$

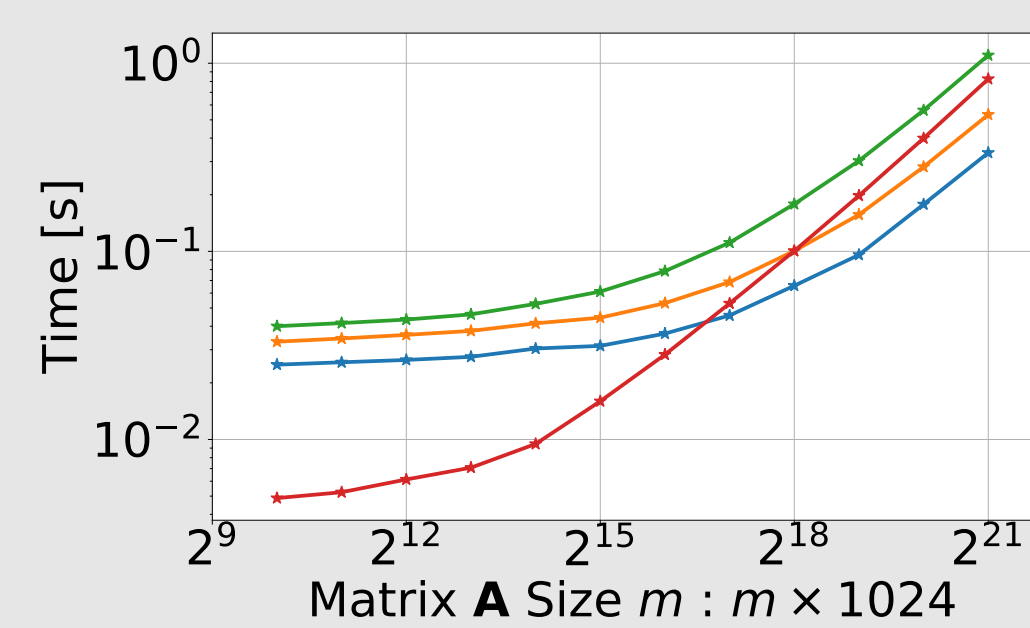
($r_{set} = 128, \alpha = 10^{-2}$)



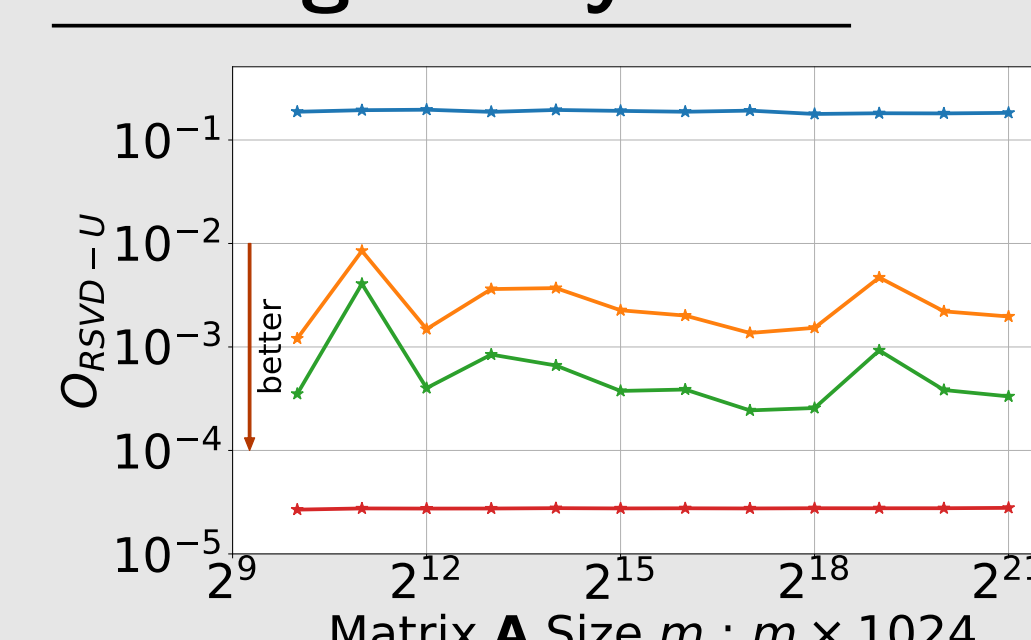
Residual



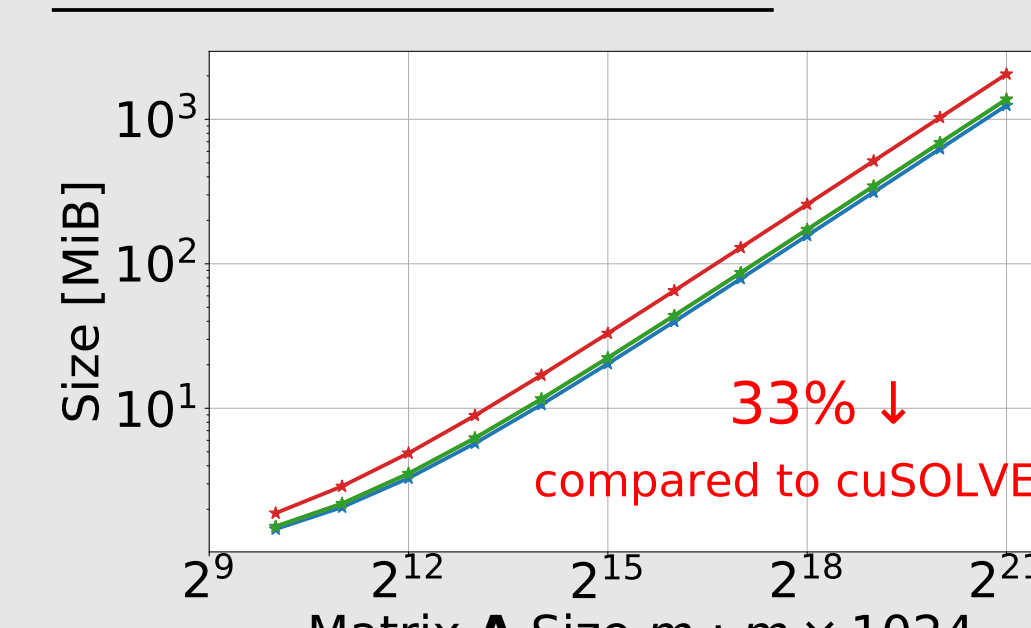
Time



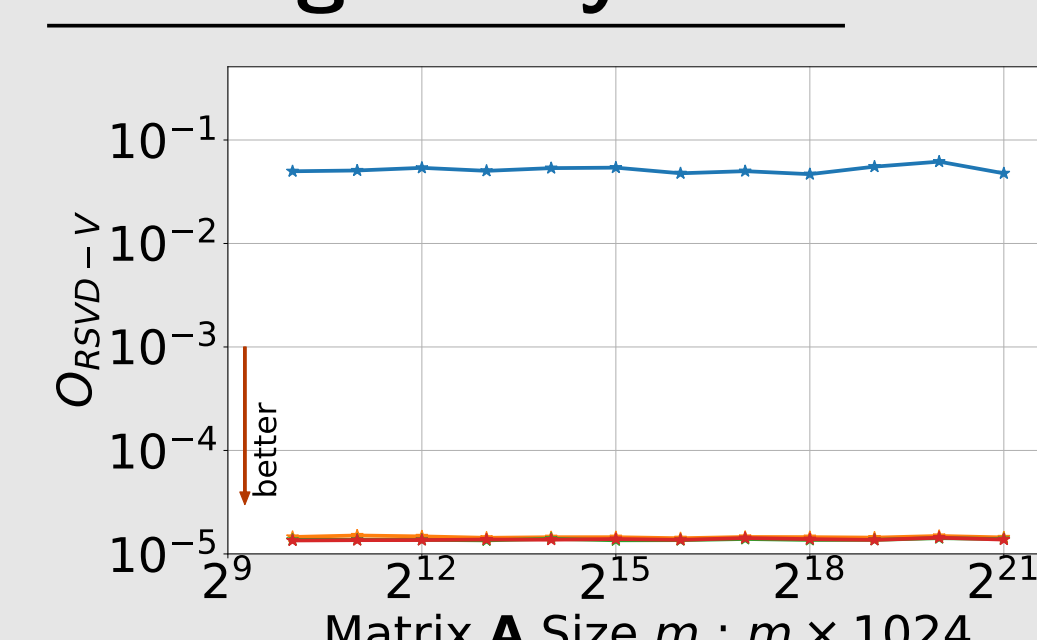
Orthogonality of U



Working memory



Orthogonality of V



- ▶ Residual $R = \frac{\|A - U\Sigma V^T\|_F}{\|A\|_F}$
- ▶ Orthogonality of M $O = \frac{\|I - M^T M\|_F}{\sqrt{n}}$

Conclusion

- ▶ Using TensorCores and correction techniques, our approach can calculate Randomized SVD without much loss of accuracy.
- ▶ Our approach provides 1.5x faster performance and reduces working memory by 33% compared to cuSOLVER.

Future work

- ▶ Investigate the effect of using FP16 to compute the Randomized SVD in depth.
- ▶ Experiment with different singular value decay patterns.
- ▶ Apply the Randomized SVD using TensorCores to some applications and investigate stability and performance.