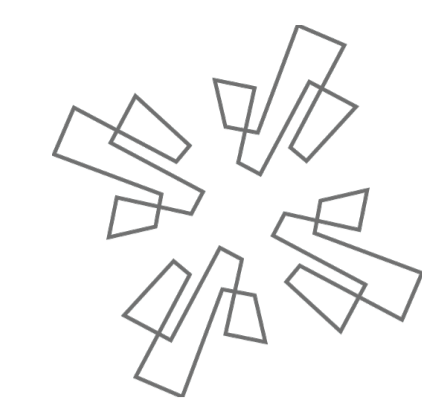


TSQR on TensorCores

Hiroyuki Ootomo¹, Rio Yokota²

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



Tokyo Tech

<https://github.com/enp1s0/tsqr-gpu>

Abstract

Tall-Skinny QR (TSQR) is an efficient algorithm for calculating the QR decomposition of $m \times n$ matrices where $m \gg n$, which is done by recursively performing QR decomposition on subdivided blocks of the tall and skinny matrix. Such operations are useful for low-rank approximation methods, which are replacing more and more dense linear algebra in both scientific computing and machine learning fields. The present work focuses on the implementation of this important algorithm on Tensor Cores, which are available on the latest NVIDIA GPUs. We evaluate the speed, accuracy, and stability of TSQR on TensorCores.

Test Environment

Machine

- ▶ Intel Xeon CPU E5-2630 v3 @ 2.40GHz x2
- ▶ NVIDIA Tesla V100-PCIE-16GB
- ▶ 64GB RAM
- ▶ Ubuntu 18.04
- ▶ CUDA 10.1

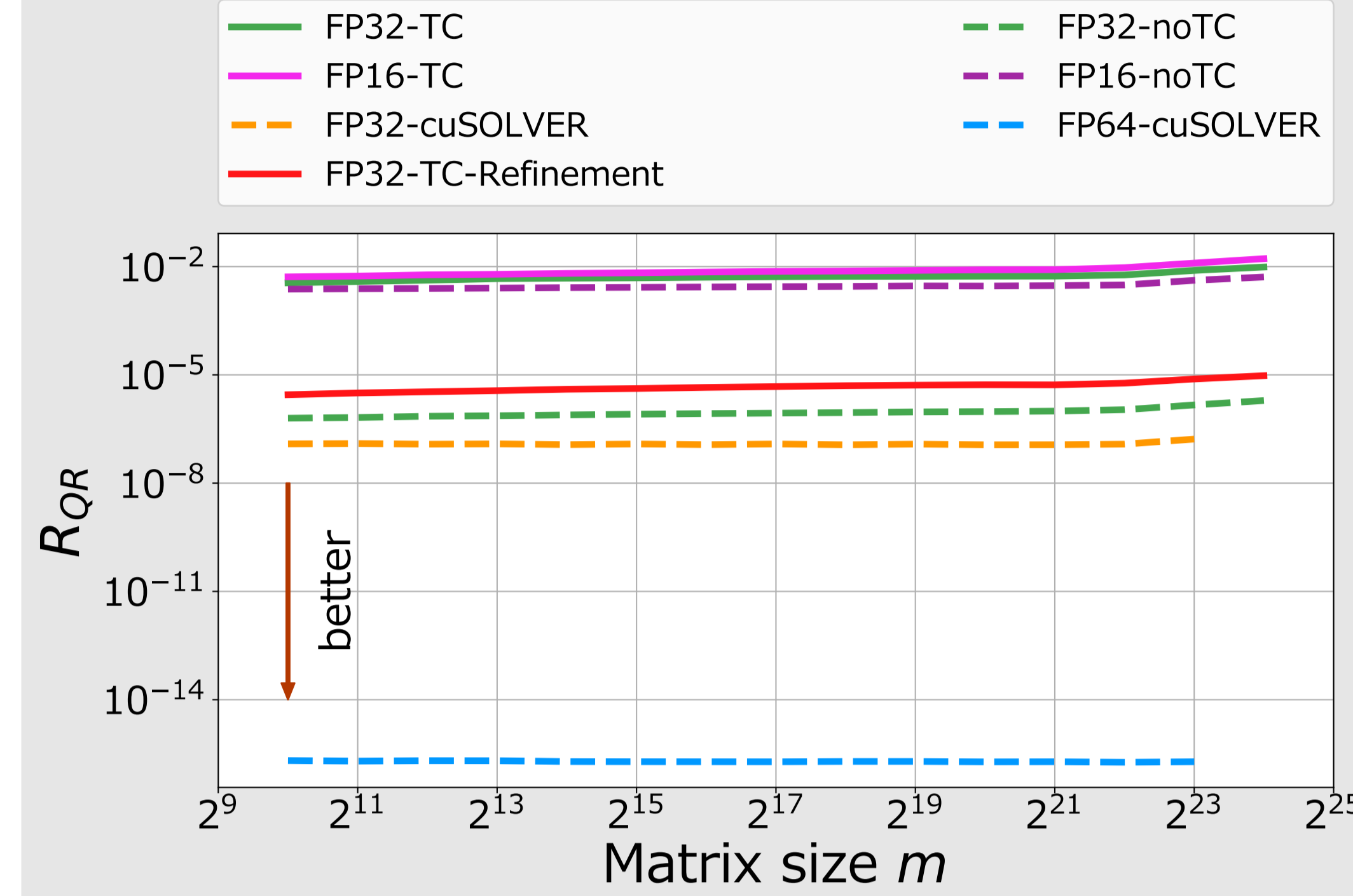
Input Matrix

- ▶ $m \times n$ matrix ($n = 16$ fixed)
- ▶ Randomized with $[-1, 1]$ uniform distribution

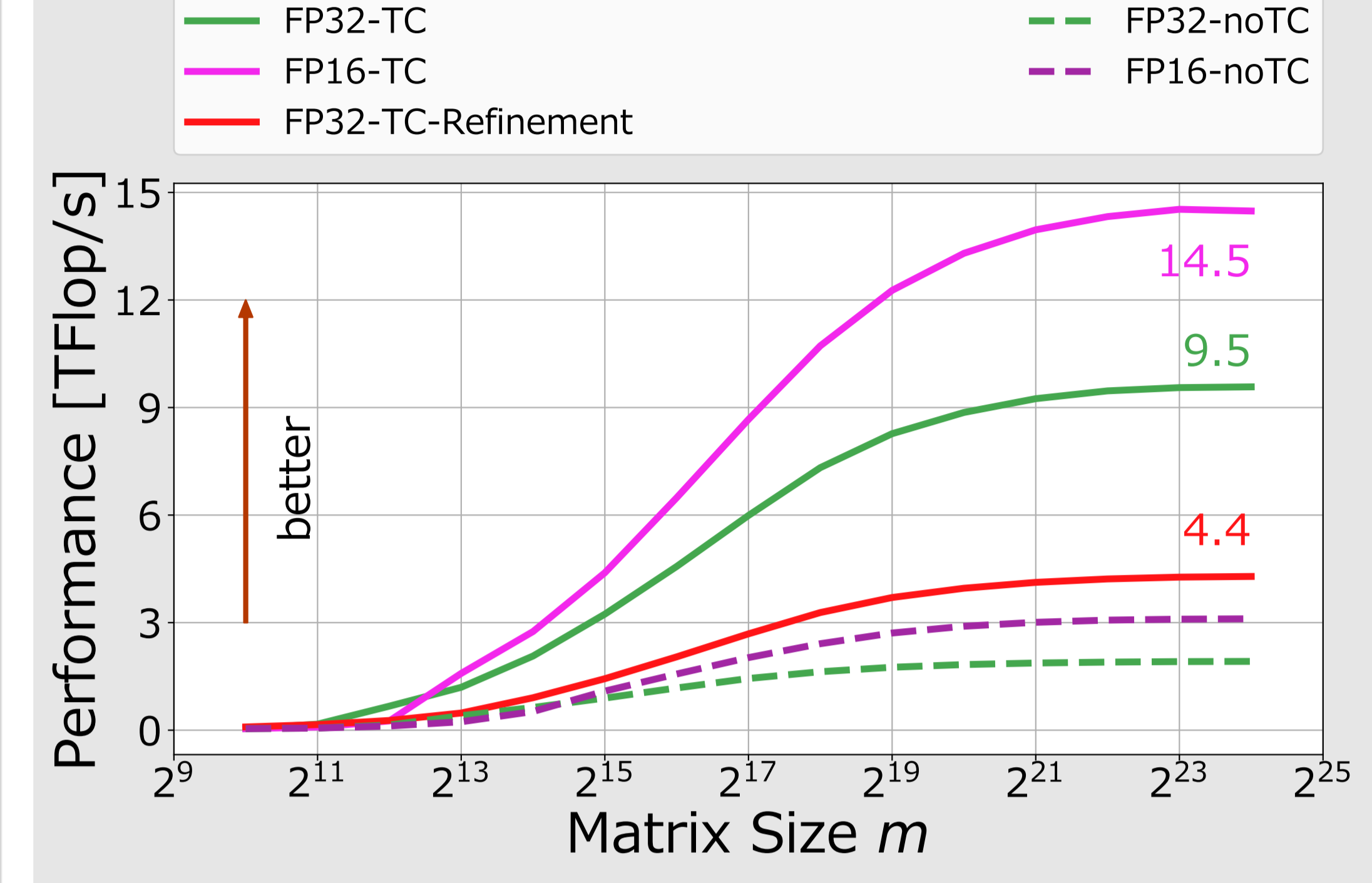
Comparison

- ▶ cuSOLVER (FP32, FP64)

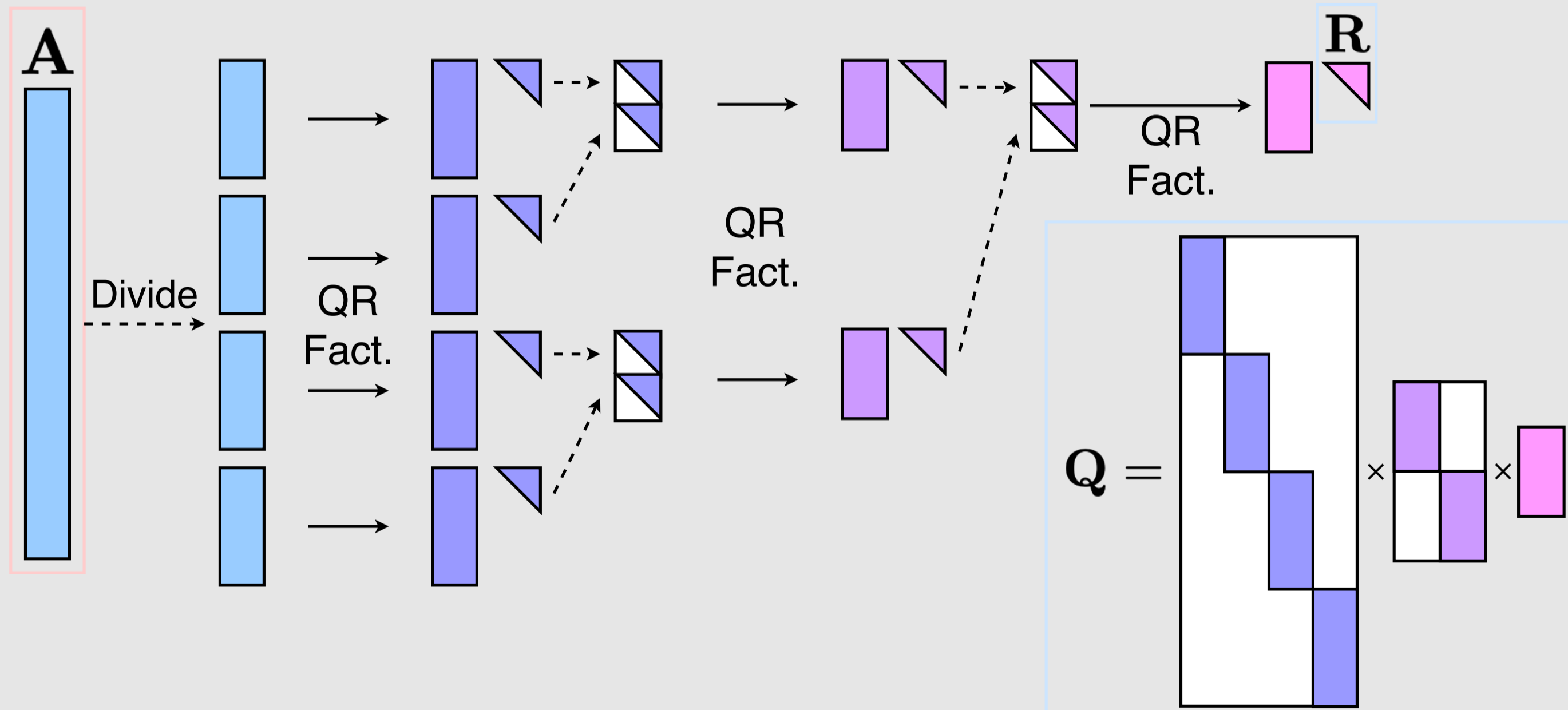
Residual Evaluation



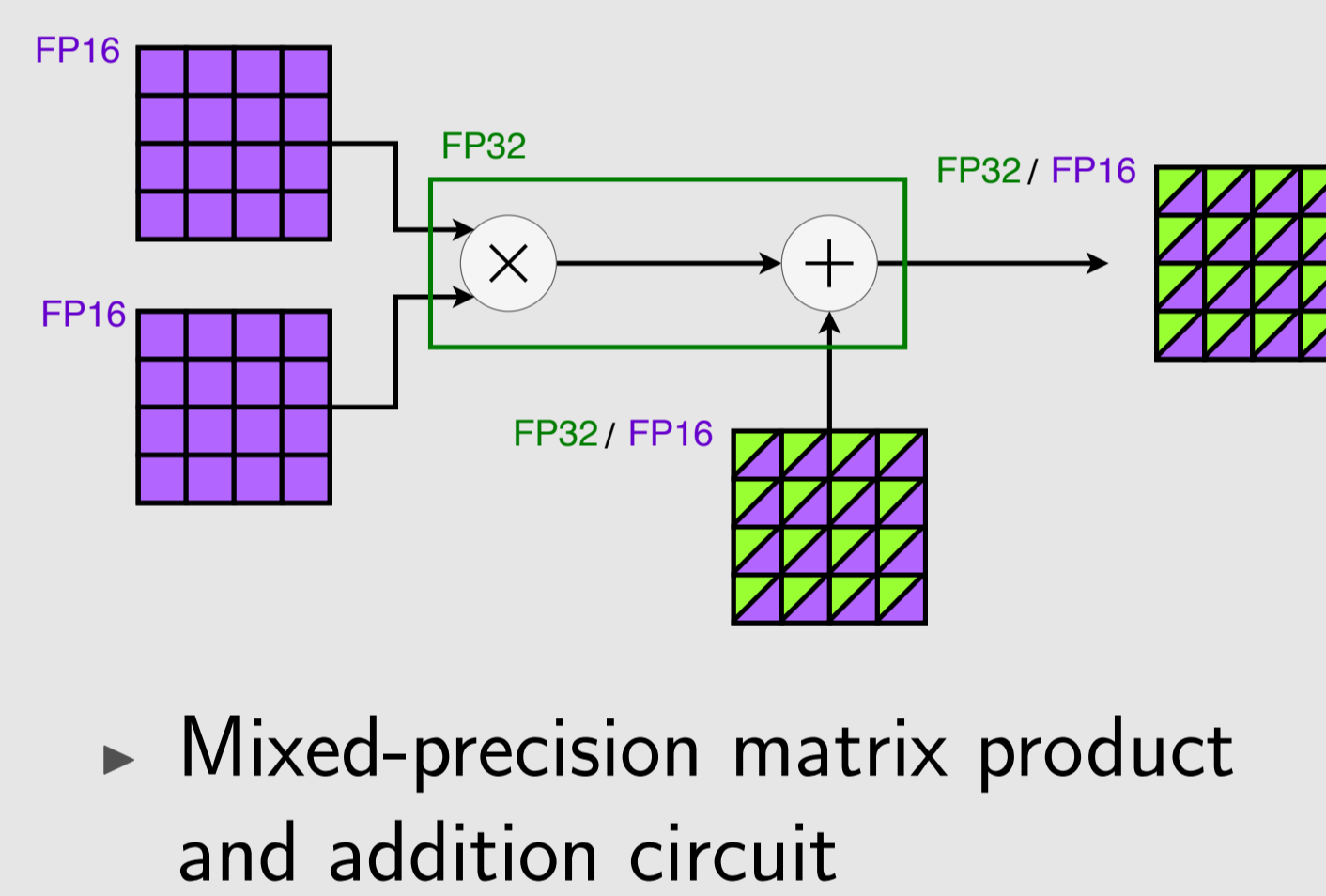
Performance Evaluation



TSQR on TensorCores



TensorCore



Refinement using TensorCores

$$\begin{aligned} C_{FP32} &\leftarrow A_{FP16} B_{FP16} \\ &+ \Delta A_{FP16} B_{FP16} \\ &+ A_{FP16} \Delta B_{FP16} \end{aligned}$$

where

$$\begin{aligned} M_{FP16} &\leftarrow \text{narrowing } M_{FP32} \\ \Delta M_{FP16} &\leftarrow \text{narrowing } M_{FP32} - M_{FP16} \end{aligned}$$

TSQR Implementation

Name	Type	TensorCore
FP32-TC	float	Used
FP16-TC	half	Used
FP32-noTC	float	Not used
FP16-noTC	half	Not used

FP32-TC with refinement

- ▶ $m \times n$ ($n \leq 16$)

Where to use TensorCore

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

Algorithm 1. TSQR

1. Divide the input matrix A .
2. Calculate QR decomposition for each subdivided matrices 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.

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 implementation

- ▶ Parallel QR Factorization for some $m \times n$ matrices ($16 \leq m \leq 32, n \leq 16$)
- ▶ Householder QR

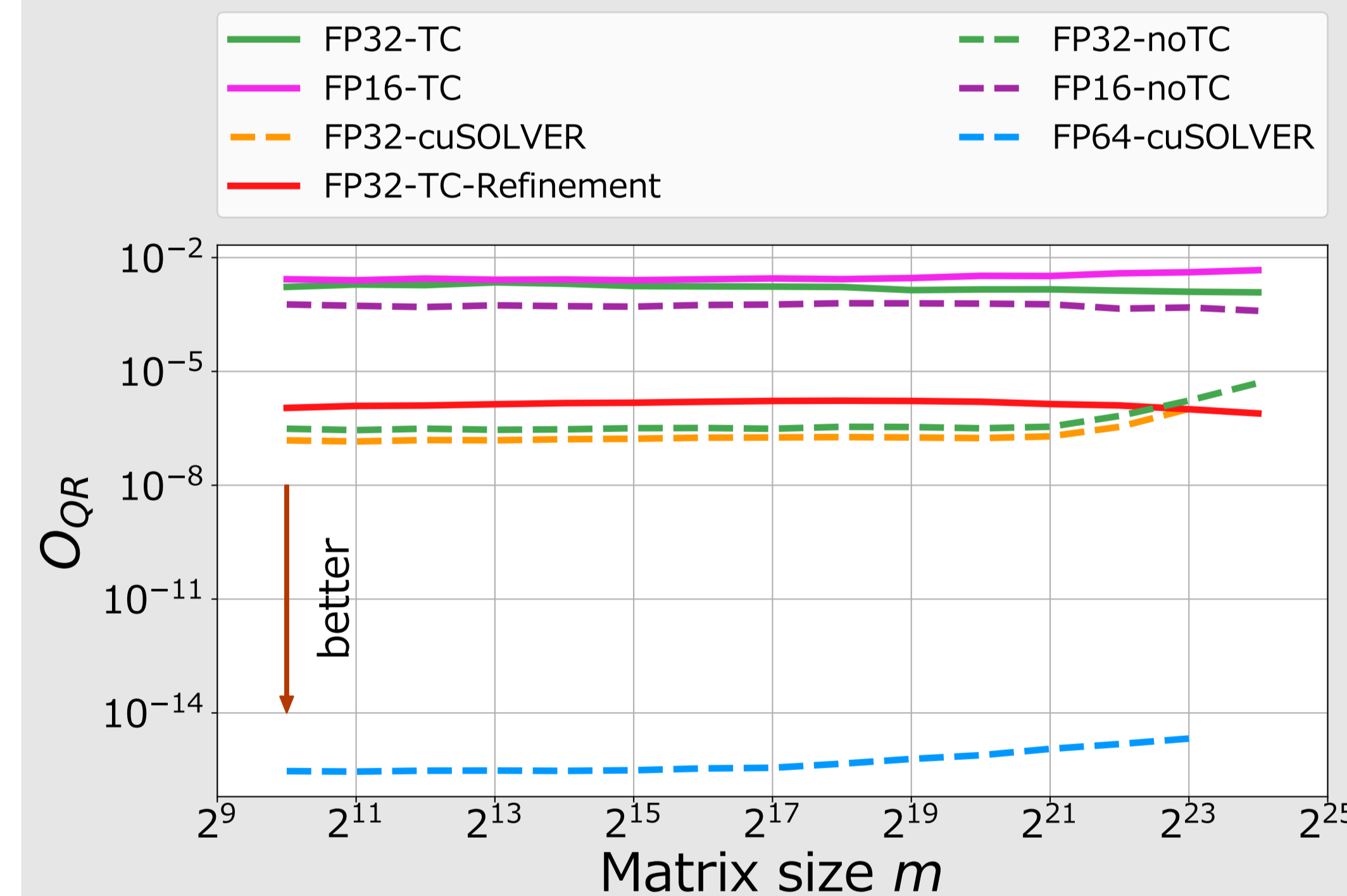
Algorithm 2. Householder QR

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

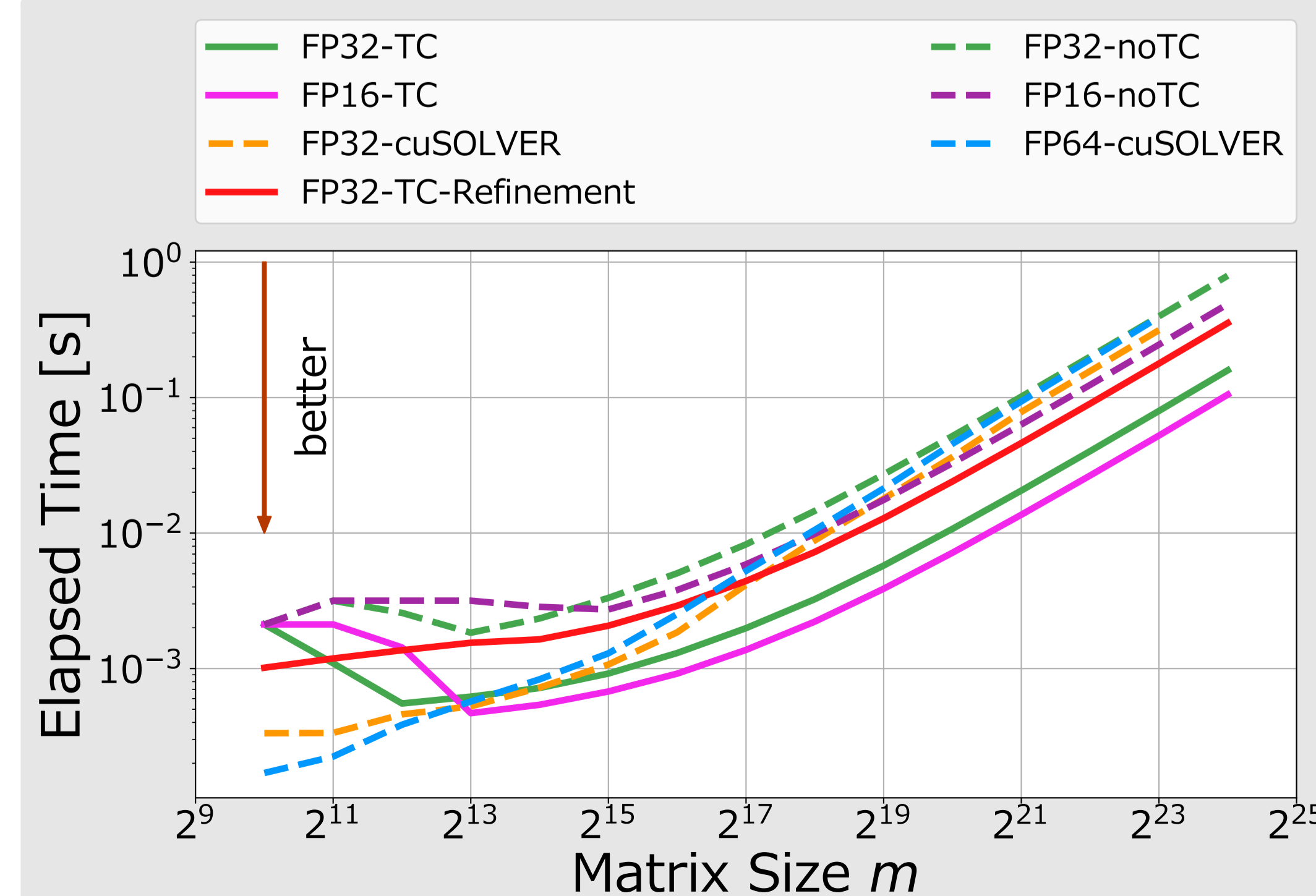
Ensure: $Q \in \mathbb{R}^{m \times m}, R \in \mathbb{R}^{m \times n}$

- 1: $Q' \leftarrow I, R \leftarrow A$
- 2: for $i \leftarrow 0$ to $n - 1$ do
- 3: $u \leftarrow [0 \dots 0 R_{i,i} \dots R_{m-1,i}]^T$
- 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: $Q \leftarrow Q'^T$

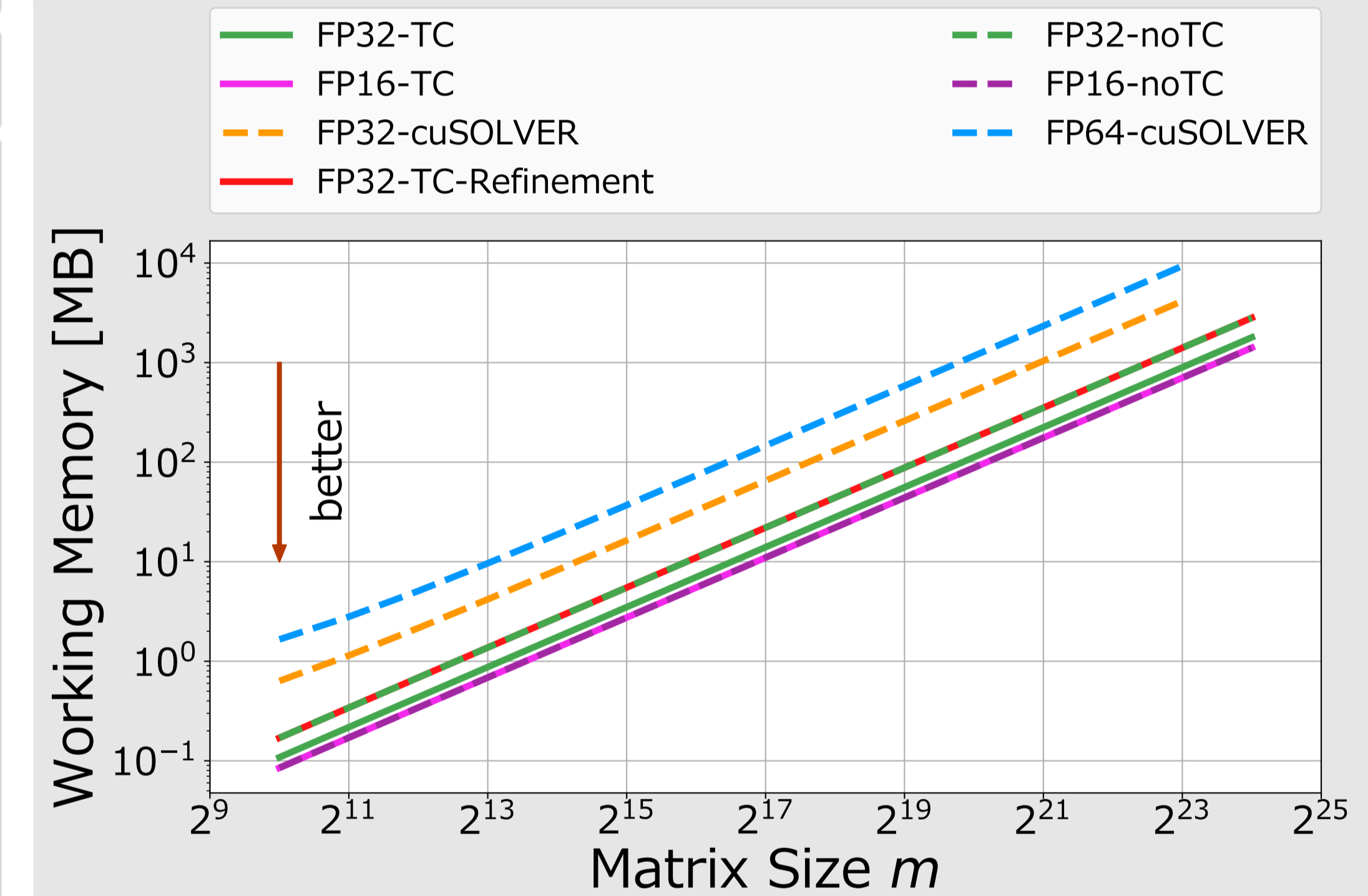
Orthogonality Evaluation



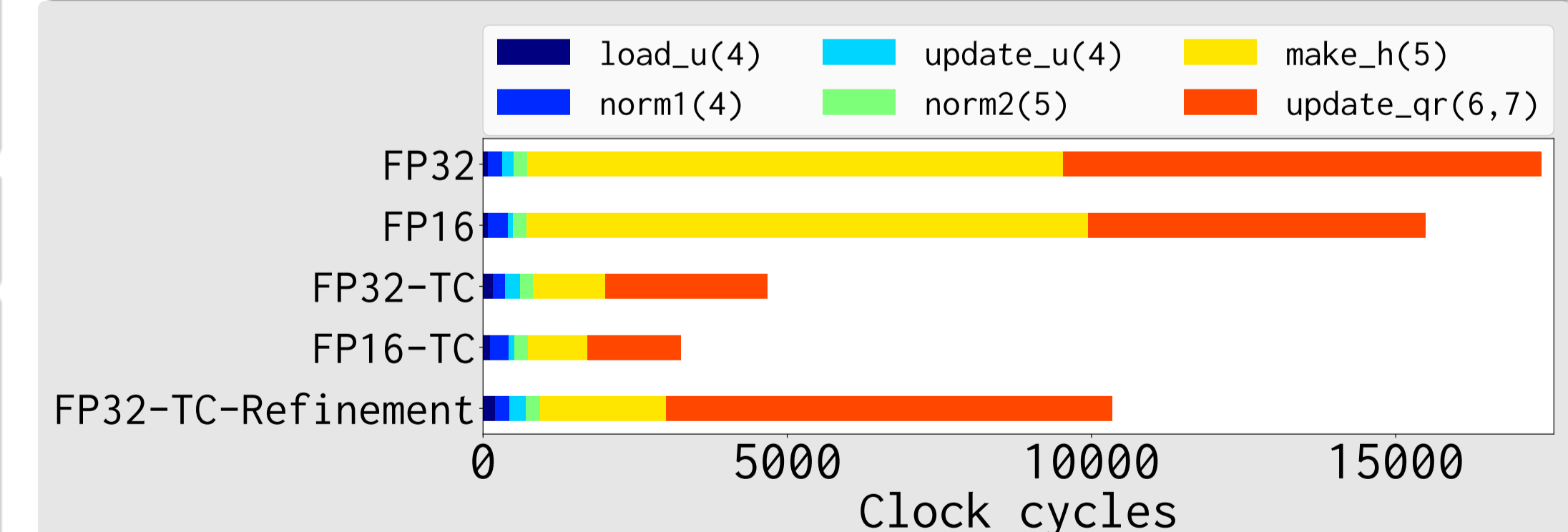
Speed Evaluation



Working Memory Size Evaluation



Computing Time Profile



The number in brackets corresponds to a line number in Algo 2.

Conclusions

- ▶ Using TensorCores and refinement, Our approach can calculate TSQR efficiently without much loss of accuracy
- ▶ Our approach provides 3.4x faster performance compared to cuSOLVER
- ▶ Our approach reduces about 80% of working memory compared to cuSOLVER