# Performance and Accuracy of Mixed-Precision Matrix Factorizations with GPU Tensor Cores

Theo Mary, CNRS
Joint work with P. Blanchard, N. J. Higham, F. Lopez, and S. Pranesh
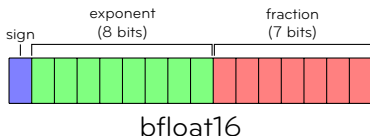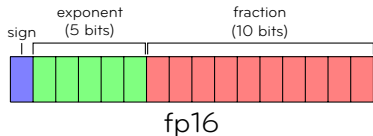
# Today's floating-point precision arithmetics

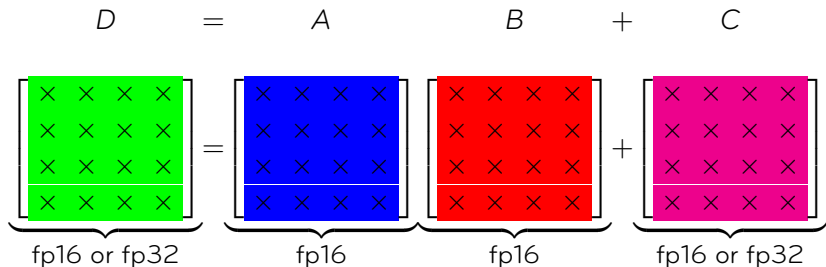| Type | | Bits | Range | $u = 2^{-t}$ |
|------|------|------|-------|------|
| fp128 | quad | 128 | $10^{\pm 4932}$ | $2^{-113} \approx 1 \times 10^{-34}$ |
| fp64 | double | 64 | $10^{\pm 308}$ | $2^{-53} \approx 1 \times 10^{-16}$ |
| fp32 | single | 32 | $10^{\pm 38}$ | $2^{-24} \approx 6 \times 10^{-8}$ |
| fp16 | half | 16 | $10^{\pm 5}$ | $2^{-11} \approx 5 \times 10^{-4}$ |
| bfloat16 | half | 16 | $10^{\pm 38}$ | $2^{-8} \approx 4 \times 10^{-3}$ |

Half precision increasingly **supported by hardware**:

- Present: **NVIDIA** Pascal & Volta GPUs, **AMD** Radeon Instinct MI25 GPU, **Google** TPU, **ARM** NEON
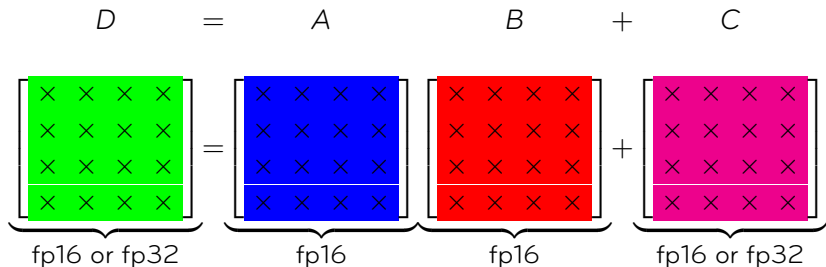- Near future: Fujitsu A64FX ARM, **IBM** AI chips, **Intel** Xeon Cooper Lake and Intel Nervana Neural Network

Mixed-Precision Matrix Factorizations                     Theo Mary

$4 \times 4$ matrix multiplication **in 1 clock cycle**:



- This is a **block fused multiply-add** (FMA) in terms of speed (in terms of accuracy, depends on the implementation)
- $\Rightarrow$ Performance peak $125$ TFlops/s ($8\times$ speedup vs fp32!)
- Algorithms now become intrinsically **mixed precision**—and more complicated to analyze

$4 \times 4$ matrix multiplication **in 1 clock cycle**:



$$D \qquad = \qquad A \qquad\qquad B \qquad + \qquad C$$

fp16 or fp32       fp16       fp16       fp16 or fp32

- This is a **block fused multiply-add** (FMA) in terms of speed (in terms of accuracy, depends on the implementation)
⇒ Performance peak $125$ TFlops/s ($8\times$ speedup vs fp32!)
- Algorithms now become intrinsically **mixed precision**—and more complicated to analyze
⇒ Need for **new analysis** to understand how to **best use these new units**

We consider the following framework

- $A \in \mathbb{R}^{b_1 \times b}$, $B \in \mathbb{R}^{b \times b_2}$, and $C \in \mathbb{R}^{b_1 \times b_2}$,

$$\underbrace{D}_{u_{\text{low}} \text{ or } u_{\text{high}}} = \underbrace{C}_{u_{\text{low}} \text{ or } u_{\text{high}}} + \underbrace{A}_{u_{\text{low}}} \underbrace{B}_{u_{\text{low}}}$$

- $AB$ is computed with multiplications in precision $u_{\text{mul}}$ and additions in precision $u_{\text{add}}$, and then rounded to precision $u_{\text{FMA}} = u_{\text{high}}$ or $u_{\text{low}}$

$$|\widehat{D} - D| \lesssim u_{\text{FMA}}(|C| + |A||B|) + ((b-1)u_{\text{add}} + u_{\text{mul}})|A||B|$$

- What choice of $u_{\text{add}}$ and $u_{\text{mul}}$?
  - $u_{\text{add}} = u_{\text{mul}} = 0$: true FMA (only 1 rounding error per element of $D$)
  - $u_{\text{add}} = u_{\text{mul}} = u_{\text{low}}$: not an FMA in terms of accuracy, just speed
  - $u_{\text{add}} = u_{\text{low}}$, $u_{\text{mul}} = u_{\text{high}}$: not really an FMA either
  - $u_{\text{add}} = u_{\text{mul}} = u_{\text{high}}$: almost an FMA (FMA to first order)

|  | $b_1$ | $b$ | $b_2$ | $u_{\text{low}}$ | $u_{\text{high}}$ |
|---|---|---|---|---|---|
| Google TPU v1 | 256 | 256 | 256 | bfloat16 | fp32 |
| Google TPU v2 | 128 | 128 | 128 | bfloat16 | fp32 |
| NVIDIA Volta | 4 | 4 | 4 | fp16 | fp32 |
| Intel NNP-T | 32 | 32 | 32 | bfloat16 | fp32 |
| Armv8-A | 2 | 4 | 2 | bfloat16 | fp32 |

- What are $u_{\text{add}}$ and $u_{\text{mul}}$? $\Rightarrow$ not entirely clear. For tensor cores:

  *Element-wise multiplication of matrix A and B is performed with at least single precision. When .ctype or .dtype is .f32, accumulation of the intermediate values is performed with at least single precision. When both .ctype and .dtype are specified as .f16, the accumulation is performed with at least half precision. The accumulation order, rounding and handling of subnormal inputs is unspecified.*

$\Rightarrow$ In the following we distinguish two variants:
  - TC16 ($u_{\text{FMA}} = u_{\text{add}} = u_{\text{low}} = u_{16}$, $u_{\text{mul}} = u_{\text{high}} = u_{32}$)
  - TC32 ($u_{\text{FMA}} = u_{\text{add}} = u_{\text{mul}} = u_{\text{high}} = u_{32}$)
  - Intermediate variant $u_{\text{FMA}} = u_{\text{mul}} = u_{32}$ and $u_{\text{add}} = u_{16}$ not discussed here

Mixed-Precision Matrix Factorizations
Theo Mary

This algorithm computes $C = AB$ using a block FMA, where $A, B, C \in \mathbb{R}^{n \times n}$, and returns $C$ in precision $u_{\text{FMA}}$

$\widetilde{A} \leftarrow \text{fl}_{\text{low}}(A)$ and $\widetilde{B} \leftarrow \text{fl}_{\text{low}}(B)$ (if necessary)
**for** $i = 1: n/b_1$ **do**
    **for** $j = 1: n/b_2$ **do**
        $C_{ij} = 0$
        **for** $k = 1: n/b$ **do**
            Compute $C_{ij} = C_{ij} + \widetilde{A}_{ik}\widetilde{B}_{kj}$ using a block FMA
        **end for**
    **end for**
**end for**

Mixed-Precision Matrix Factorizations

## Matrix multiplication: error analysis

Let $A$ and $B$ already be given in precision $u_{\text{low}}$. For any row $x$ of $A$ and any column $y$ of $B$, computing $s = c + x^T y$ classically produces

$$\widehat{s} = c(1 + \theta_n) + x_1 y_1 (1 + \theta_{n+1}) + x_2 y_2 (1 + \theta'_n)$$
$$+ x_3 y_3 (1 + \theta_{n-1}) + \cdots + x_n y_n (1 + \theta_2),$$

where $z_k = x_k y_k$ and $|\theta_k| \lesssim ku$.

Let $A$ and $B$ already be given in precision $u_{\mathrm{low}}$. For any row $x$ of $A$ and any column $y$ of $B$, computing $s = c + x^T y$ classically produces

$$\widehat{s} = c(1 + \theta_n) + x_1 y_1 (1 + \theta_{n+1}) + x_2 y_2 (1 + \theta'_n)$$
$$+ x_3 y_3 (1 + \theta_{n-1}) + \cdots + x_n y_n (1 + \theta_2),$$

where $z_k = x_k y_k$ and $|\theta_k| \lesssim ku$. With a block FMA, we have instead

$$\widehat{s} = \left( z_1 (1 + \epsilon_1)\left(1 + \theta^{(1)}_{b-1}\right) + \cdots + z_b (1 + \epsilon_b)\left(1 + \theta^{(1)}_1\right) \right) \prod_{i=1}^{n/b} (1 + \delta_i)$$

$$+ \cdots +$$

$$\left( z_{n-b+1} (1 + \epsilon_{n-b+1})\left(1 + \theta^{(n/b)}_{b-1}\right) + \cdots + z_n (1 + \epsilon_n)\left(1 + \theta^{(n/b)}_1\right) \right) (1 + \delta_{n/b})$$

where $|\epsilon_k| \leq u_{\mathrm{mul}}$, $|\theta_k| \lesssim ku_{\mathrm{add}}$, and $|\delta_k| \leq u_{\mathrm{FMA}}$

$$\boxed{\text{Overall: } |s - \widehat{s}| \lesssim \left( \tfrac{n}{b} u_{\mathrm{FMA}} + (b-1) u_{\mathrm{add}} + u_{\mathrm{mul}} \right) |x|^T |y|}$$

If $A$ and $B$ are already given in precision $u_{\text{low}}$:

$$\widehat{C} = AB + \Delta C, \quad |\Delta C| \lesssim \left(\frac{n}{b}u_{\text{FMA}} + (b-1)u_{\text{add}} + u_{\text{mul}}\right)|A||B|$$

If not, we must account for the initial conversion:

$$\widetilde{A} = \text{fl}_{\text{low}}(A) = A + \Delta A, \quad |\Delta A| \leq u_{\text{low}}|A|,$$
$$\widetilde{B} = \text{fl}_{\text{low}}(B) = B + \Delta B, \quad |\Delta B| \leq u_{\text{low}}|B|.$$

$$\widehat{C} = \widetilde{A}\widetilde{B} + \Delta C, \qquad |\Delta C| \lesssim \left(\frac{n}{b}u_{\text{FMA}} + (b-1)u_{\text{add}} + u_{\text{mul}}\right)|\widetilde{A}||\widetilde{B}|,$$
$$= AB + \Delta AB + A\Delta B + \Delta A\Delta B + \Delta C$$
$$= AB + E, \qquad |E| \lesssim \left(2u_{\text{low}} + \frac{n}{b}u_{\text{FMA}} + (b-1)u_{\text{add}} + u_{\text{mul}}\right)|A||B|$$

If $A$ and $B$ are already given in precision $u_{\text{low}}$:

$$\widehat{C} = AB + \Delta C, \quad |\Delta C| \lesssim \big(\frac{n}{b}u_{\text{FMA}} + (b-1)u_{\text{add}} + u_{\text{mul}}\big)|A||B|$$

If not, we must account for the initial conversion:

$$\widetilde{A} = \text{fl}_{\text{low}}(A) = A + \Delta A, \quad |\Delta A| \leq u_{\text{low}}|A|,$$
$$\widetilde{B} = \text{fl}_{\text{low}}(B) = B + \Delta B, \quad |\Delta B| \leq u_{\text{low}}|B|.$$

$$\widehat{C} = \widetilde{A}\widetilde{B} + \Delta C, \qquad |\Delta C| \lesssim \big(\frac{n}{b}u_{\text{FMA}} + (b-1)u_{\text{add}} + u_{\text{mul}}\big)|\widetilde{A}||\widetilde{B}|,$$
$$= AB + \Delta AB + A\Delta B + \Delta A\Delta B + \Delta C$$
$$= AB + E, \qquad |E| \lesssim \big( \underbrace{2u_{\text{low}}}_{\text{Conversion}} + \underbrace{\frac{n}{b}u_{\text{FMA}} + (b-1)u_{\text{add}} + u_{\text{mul}}}_{\text{Accumulation}}\big)|A||B|$$

# Matrix multiplication: error bounds interpretation

| Evaluation method ($u_{mul} = u_{high}$) | | | Bound |
|---|---|---|---|
| Standard in precision $u_{low}$ | | | $nu_{low}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = u_{low}$ | $(n/b + b)u_{low}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = u_{high}$ | $(n/b)u_{low} + bu_{high}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = 0$ | $(n/b)u_{low}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = u_{low}$ | $(b + 2)u_{low} + (n/b)u_{high}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = u_{high}$ | $2u_{low} + (n/b + b)u_{high}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = 0$ | $2u_{low} + (n/b)u_{high}$ |
| Standard in precision $u_{high}$ | | | $nu_{high}$ |

| Evaluation method ($u_{mul} = u_{high}$) | | | Bound |
|---|---|---|---|
| Standard in precision $u_{low}$ | | | $nu_{low}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = u_{low}$ | $(n/b + b)u_{low}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = u_{high}$ | $(n/b)u_{low} + bu_{high}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = 0$ | $(n/b)u_{low}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = u_{low}$ | $(b + 2)u_{low} + (n/b)u_{high}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = u_{high}$ | $2u_{low} + (n/b + b)u_{high}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = 0$ | $2u_{low} + (n/b)u_{high}$ |
| Standard in precision $u_{high}$ | | | $nu_{high}$ |

- $u_{FMA} = u_{add} = u_{low} \Rightarrow$ reduction by factor $b$ from blocked sum

| Evaluation method ($u_{\mathrm{mul}} = u_{\mathrm{high}}$) | | | Bound |
|:---:|:---:|:---:|:---:|
| Standard in precision $u_{\mathrm{low}}$ | | | $nu_{\mathrm{low}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{low}}$ | $u_{\mathrm{add}} = u_{\mathrm{low}}$ | $(n/b + b)u_{\mathrm{low}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{low}}$ | $u_{\mathrm{add}} = u_{\mathrm{high}}$ | $(n/b)u_{\mathrm{low}} + bu_{\mathrm{high}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{low}}$ | $u_{\mathrm{add}} = 0$ | $(n/b)u_{\mathrm{low}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{high}}$ | $u_{\mathrm{add}} = u_{\mathrm{low}}$ | $(b + 2)u_{\mathrm{low}} + (n/b)u_{\mathrm{high}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{high}}$ | $u_{\mathrm{add}} = u_{\mathrm{high}}$ | $2u_{\mathrm{low}} + (n/b + b)u_{\mathrm{high}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{high}}$ | $u_{\mathrm{add}} = 0$ | $2u_{\mathrm{low}} + (n/b)u_{\mathrm{high}}$ |
| Standard in precision $u_{\mathrm{high}}$ | | | $nu_{\mathrm{high}}$ |

- $u_{\mathrm{FMA}} = u_{\mathrm{add}} = u_{\mathrm{low}} \Rightarrow$ reduction by factor $b$ from blocked sum
- $u_{\mathrm{FMA}} = u_{\mathrm{low}}$, $u_{\mathrm{add}} \ll u_{\mathrm{low}} \Rightarrow$ smaller $u_{\mathrm{add}}$ not very useful

# Matrix multiplication: error bounds interpretation

| Evaluation method ($u_{\mathrm{mul}} = u_{\mathrm{high}}$) | | | Bound |
|---|---|---|---|
| Standard in precision $u_{\mathrm{low}}$ | | | $nu_{\mathrm{low}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{low}}$ | $u_{\mathrm{add}} = u_{\mathrm{low}}$ | $(n/b + b)u_{\mathrm{low}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{low}}$ | $u_{\mathrm{add}} = u_{\mathrm{high}}$ | $(n/b)u_{\mathrm{low}} + bu_{\mathrm{high}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{low}}$ | $u_{\mathrm{add}} = 0$ | $(n/b)u_{\mathrm{low}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{high}}$ | $u_{\mathrm{add}} = u_{\mathrm{low}}$ | $(b + 2)u_{\mathrm{low}} + (n/b)u_{\mathrm{high}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{high}}$ | $u_{\mathrm{add}} = u_{\mathrm{high}}$ | $2u_{\mathrm{low}} + (n/b + b)u_{\mathrm{high}}$ |
| Block FMA | $u_{\mathrm{FMA}} = u_{\mathrm{high}}$ | $u_{\mathrm{add}} = 0$ | $2u_{\mathrm{low}} + (n/b)u_{\mathrm{high}}$ |
| Standard in precision $u_{\mathrm{high}}$ | | | $nu_{\mathrm{high}}$ |

- $u_{\mathrm{FMA}} = u_{\mathrm{add}} = u_{\mathrm{low}} \Rightarrow$ reduction by factor $b$ from blocked sum

- $u_{\mathrm{FMA}} = u_{\mathrm{low}}$, $u_{\mathrm{add}} \ll u_{\mathrm{low}} \Rightarrow$ smaller $u_{\mathrm{add}}$ not very useful

- $u_{\mathrm{FMA}} = u_{\mathrm{high}} \Rightarrow$ reduction by factor $\min(n/2, u_{\mathrm{low}}/u_{\mathrm{high}})$,
  $u_{\mathrm{add}} \ll u_{\mathrm{low}}$ useful, $u_{\mathrm{add}} = 0$ not useful

| Evaluation method ($u_{mul} = u_{high}$) | | | Bound |
|---|---|---|---|
| Standard in precision $u_{low}$ | | | $nu_{low}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = u_{low}$ | $(n/b + b)u_{low}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = u_{high}$ | $(n/b)u_{low} + bu_{high}$ |
| Block FMA | $u_{FMA} = u_{low}$ | $u_{add} = 0$ | $(n/b)u_{low}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = u_{low}$ | $(b + 2)u_{low} + (n/b)u_{high}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = u_{high}$ | $2u_{low} + (n/b + b)u_{high}$ |
| Block FMA | $u_{FMA} = u_{high}$ | $u_{add} = 0$ | $2u_{low} + (n/b)u_{high}$ |
| Standard in precision $u_{high}$ | | | $nu_{high}$ |

- $u_{FMA} = u_{add} = u_{low} \Rightarrow$ reduction by factor $b$ from blocked sum
- $u_{FMA} = u_{low}$, $u_{add} \ll u_{low} \Rightarrow$ smaller $u_{add}$ not very useful
- $u_{FMA} = u_{high} \Rightarrow$ reduction by factor $\min(n/2, u_{low}/u_{high})$,
  $u_{add} \ll u_{low}$ useful, $u_{add} = 0$ not useful
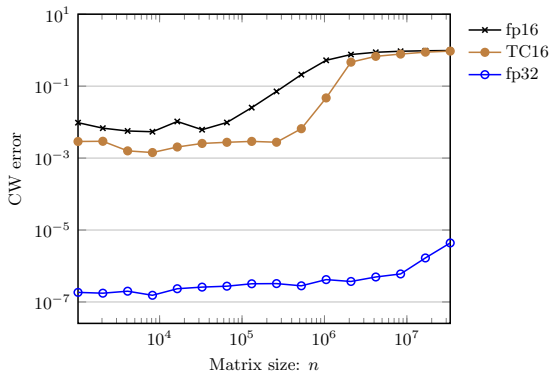
> **Conclusion: choice of $u_{FMA}$ critical, $u_{add}$ less so**

# Matrix multiplication with tensor cores

| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|:---:|:---:|:---:|:---:|
| $nu_{16}$ | $(n/4)u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |

Mixed-Precision Matrix Factorizations

Theo Mary

# Matrix multiplication with tensor cores

| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|:---:|:---:|:---:|:---:|
| $nu_{16}$ | $(n/4)u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |



Mixed-Precision Matrix Factorizations Theo Mary

| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|---|---|---|---|
| $nu_{16}$ | $(n/4)u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |

| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|:---:|:---:|:---:|:---:|
| $nu_{16}$ | $(n/4)u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |

| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|:---:|:---:|:---:|:---:|
| $nu_{16}$ | $(n/4)u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |



Matrix size: $n$

- - - perf. peak
- TC16
- TC32
- fp16
- fp32

**Conclusion: TC32 significantly more accurate than TC16, with almost no performance loss**
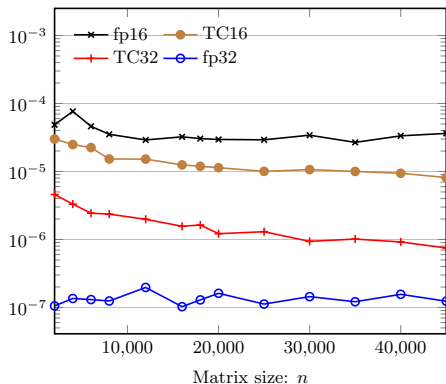
Mixed-Precision Matrix Factorizations
Theo Mary

This algorithm computes $A = LU$ using a block FMA, where $A \in \mathbb{R}^{n \times n}$ is given in precision $u_{\text{high}}$, and $L$ and $U$ are returned in precision $u_{\text{FMA}}$

---

**for** $k = 1: n/b$ **do**
    Factorize $L_{kk}U_{kk} = A_{kk}$
    **for** $i = k + 1: n/b$ **do**
        Solve $L_{ik}U_{kk} = A_{ik}$ and $L_{kk}U_{ki} = A_{ki}$ for $L_{ik}$ and $U_{ki}$
    **end for**
    **for** $i = k + 1: n/b$ **do**
        **for** $j = k + 1: n/b$ **do**
            $\widetilde{L}_{ik} \leftarrow \text{fl}_{\text{low}}(L_{ik})$ and $\widetilde{U}_{ki} \leftarrow \text{fl}_{\text{low}}(U_{ki})$
            $A_{ij} \leftarrow A_{ij} - \widetilde{L}_{ik}\widetilde{U}_{kj}$ using a block FMA
        **end for**
    **end for**
**end for**

---

Error analysis for LU follows from matrix multiplication analysis and gives same bounds to first order (minor changes in the constants)

| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|:---:|:---:|:---:|:---:|
| $nu_{16}$ | $n/4\,u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |

Mixed-Precision Matrix Factorizations · Theo Mary

# LU factorization with tensor cores

Error analysis for LU follows from matrix multiplication analysis and gives same bounds to first order (minor changes in the constants)
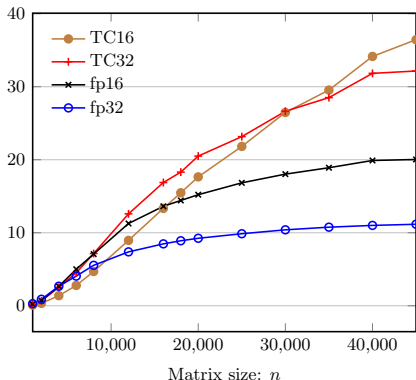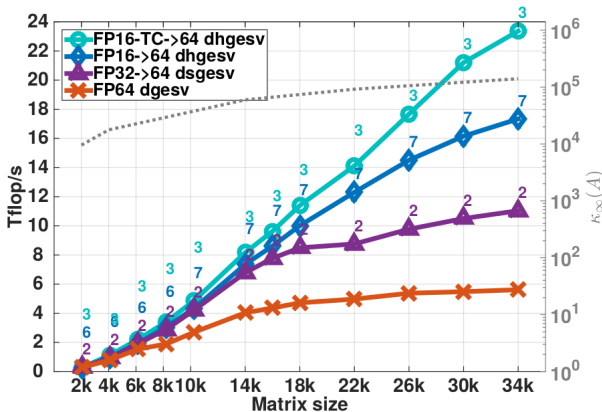
| Standard fp16 | Tensor core TC16 | Tensor core TC32 | Standard fp32 |
|:---:|:---:|:---:|:---:|
| $nu_{16}$ | $n/4\,u_{16}$ | $2u_{16} + (n/4)u_{32}$ | $nu_{32}$ |



Matrix size: $n$

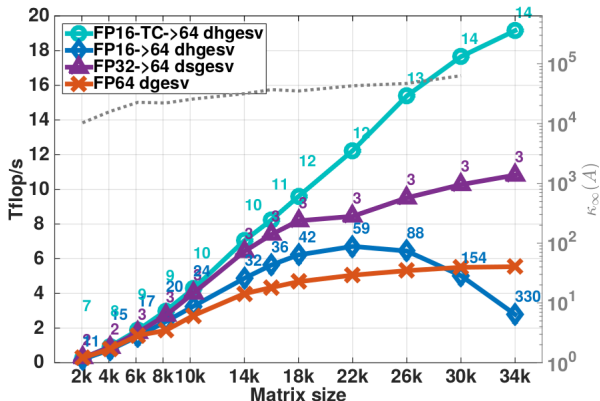Mixed-Precision Matrix Factorizations          Theo Mary

Results from Haidar et al, *Harnessing GPU tensor cores for fast FP16 arithmetic to speed up mixed-precision iterative refinement solvers* (SC'2018):



- Performance TC boost not fully translated in their implementation

Results from Haidar et al, *Harnessing GPU tensor cores for fast FP16 arithmetic to speed up mixed-precision iterative refinement solvers* (SC'2018):



- Performance TC boost not fully translated in their implementation
- But accuracy boost sometimes critical!

- Mixed-precision units increasingly available in hardware
- Proposed a general mixed-precision block FMA framework, should be applicable to existing and future units
- Performed error analysis for matrix mult. and LU factorization
- Application to NVIDIA GPU tensor cores: compared two variants, TC16 and TC32 (different computeType parameter)

$\Rightarrow$ TC32 is significantly more accurate than fp16 and TC16: reduction from $O(nu_{16})$ to $2u_{16} + O(nu_{32})$, while being almost as fast as TC16 (and much faster than fp16)

**Take-home message:** we recommend using TC32 over TC16

### Preprint and slides

See our preprint: *Mixed Precision Block Fused Multiply-Add: Error Analysis and Application to GPU Tensor Cores*
Slides available on my website: `bit.ly/tmaryLIP6`

# Ongoing work on LU factorization with tensor cores

## Partitioned LU factorization with tensor cores

The LU analysis assumes a panel of size $b = 4 \Rightarrow$ not realistic for performance, where $b = O(100)$. Possible solutions:

- Do the panel factorization in fp32 (Haidar et al, 2018) $\Rightarrow$ suboptimal performance
- Do the panel factorization in fp16 $\Rightarrow$ suboptimal accuracy
- **Our proposed solution:** use a double panel hierarchy to use mixed precision TC32 in the panel factorization

## Storing the LU factors in fp16

The LU analysis assumes the LU factors to be stored in fp32 precision $\Rightarrow$ no storage gain! Possible solutions:

- Store them in fp16 $\Rightarrow$ repeated rounding to fp16 (after each update) causes great loss of accuracy (TC16 even with computeType=fp32)
- **Our proposed solution:** use a left-looking factorization with a temporary fp32 buffer to accumulate updates $\Rightarrow$ no accuracy loss!

Mixed-Precision Matrix Factorizations Theo Mary