DGEMM on Integer Tensor Cores

Hiroyuki Ootomo Tokyo Institute of Technology

NHR PerfLab Seminar Sep. 05, 2023

\$ whoami

- Hiroyuki Ootomo
- Ph.D. candidate at Tokyo Institute of Technology
- Research:
 - Mixed-precision computing
 - Randomized numerical linear algebra
 - Quantum circuit simulation
 - HPC processors



Motivation

- Question

Can deep learning processors be used for HPC applications?

- Deep learning hardware
 - Low- and mixed-precision
 - Matrix multiplication
- Usable in HPC apps?
 - Not easily

Myth 11: HPC Will Pivot to Low or Mixed Precision!

A high-performance language is nothing without proper data types, but high-precision operations such as fp64 come at a significant cost in terms of silicon area, energy and speed, according to Myth 6. Lowering this precision can save costs

S. Matsuoka, et al. (2022)

Myths and Legends in High-Performance Computing



FIGURE 1. Overview of FP representations.

H. Ltaief, et al., (2022)

Responsibly Reckless Matrix Algorithms for HPC Scientific Applications

Outline

1 SGEMM Emulation on Tensor Cores

2 DGEMM Emulation on Integer Tensor Cores

NVIDIA Tensor Core

- Mixed-precision matrix multiplication and addition unit $(\mathbf{A} \cdot \mathbf{B} + \mathbf{C})$.
- On the latest NVIDIA GPUs.



Can we compute SGEMM on Tensor Cores?

• We need to convert the input FP32 matrices (A, B) to low-precision.



- The type conversion results in loss of accuracy in the resulting matrix
 - Almost the same as half-precision
- ⇒ Our previous study proposes a method to refine the accuracy by preserving and utilizing the mantissa that is lost during type conversion.

SGEMM emulation on Tensor Cores

- Matrix multiplication almost equivalent to single-precision on Tensor Cores.
- Faster than the theoretical peak throughput of FP32 computing units.



while surpassing the FP32 theoretical peak performance, IJHPCA, 2022

SGEMM emulation on Tensor Cores ($C_{FP32} \leftarrow A_{FP32} \cdot B_{FP32}$)



 $v_{\text{FP32}} \approx v_{\text{FP16}} + \Delta v_{\text{FP16}}$





7/26

Can we emulate DGEMM in the same manner?

No

- The accumulation precision limits the emulation accuracy.
 - In the case of Tensor Cores, it is FP32.
- If there were FP64 accumulation Tensor Core, we could emulate DGEMM
- But we can compute it using another method called the **Ozaki scheme**



Hiroyuki Ootomo, Katsuhisa Ozaki, Rio Yokota



1 SGEMM Emulation on Tensor Cores

2 DGEMM Emulation on Integer Tensor Cores

DGEMM emulation on Integer Tensor Cores

DGEMM emulation on Integer Tensor Cores using the Ozaki scheme



Ref: Ootomo, Ozaki, and Yokota, *DGEMM on Integer Matrix Multiplication Unit*, arXiv, 2023

DGEMM on Integer Tensor Cores Integer Tensor Cores





- \blacksquare The throughput is $2\times$ higher than FP16 Tensor Cores
- The inference computation in deep learning

Ozaki scheme (Ozaki et al., 2012)

 The Ozaki scheme is a high-precision matrix multiplication algorithm using lower-precision arithmetic



Ozaki scheme (Ozaki et al., 2012)

 The Ozaki scheme is a high-precision matrix multiplication algorithm using lower-precision arithmetic



Examle (in decimal number)

- **a** = $\begin{bmatrix} a_1 & a_2 \end{bmatrix}$ where $a_1 = 1.11111111 \times 10^5, a_2 = 2.2222222 \times 10^3$
- Split mantissa by 3 digits

× Similar method to the SGEMM emulation (Not the Ozaki scheme)

$$\begin{aligned} a_1 &= 1.11 \times \quad 10^5 + 1.11 \times \quad 10^2 + 1.11 \times \quad 10^{-1} \\ a_2 &= 2.22 \times \quad 10^3 + 2.22 \times \quad 10^0 + 2.22 \times \quad 10^{-3} \\ \hline \mathbf{a}^{(1)} + \quad \mathbf{a}^{(2)} + \quad \mathbf{a}^{(3)} \end{aligned}$$

The Ozaki scheme

$$\begin{aligned} a_1 &= 1.11 \times \quad 10^5 + 1.11 \times \quad 10^2 + 1.11 \times \quad 10^{-1} + 0.00 \times \quad 10^{-4} \\ a_2 &= 0.02 \times \quad 10^5 + 2.22 \times \quad 10^2 + 2.22 \times \quad 10^{-1} + 2.20 \times \quad 10^{-4} \\ \hline \mathbf{a}^{(1)} + \quad \mathbf{a}^{(2)} + \quad \mathbf{a}^{(3)} + \quad \mathbf{a}^{(4)} \end{aligned}$$

Split the mantissa while aligning the exponent.

- Rounding error in multiplication of two floating-point values



- Rounding error in addition of two floating-point values



- Consider an inner product of two length-k vectors \mathbf{a} and \mathbf{b} .
- We can compute their inner product without rounding error if
 - **1** the mantissa lengths of all elements a_i , b_i are small enough

2 the exponents of $a_i b_i$ for i = 1..k are close enough

Intuitive example: (not rigorous)

(For simplicity, assume that all elements in a vector have the same exponent)



- Consider an inner product of two length-k vectors \mathbf{a} and \mathbf{b} .
- We can compute their inner product without rounding error if
 - 1 the mantissa lengths of all elements a_i , b_i are small enough

2 the exponents of $a_i b_i$ for i = 1..k are close enough

Intuitive example: (not rigorous)

(For simplicity, assume that all elements in a vector have the same exponent)



- Consider an inner product of two length-k vectors \mathbf{a} and \mathbf{b} .
- We can compute their inner product without rounding error if
 - 1 the mantissa lengths of all elements a_i , b_i are small enough

2 the exponents of $a_i b_i$ for i = 1..k are close enough

Intuitive example: (not rigorous)

(For simplicity, assume that all elements in a vector have the same exponent)



■ If $17 + \lceil \log_2 k \rceil \le$ the mantissa length of the output format, no rounding error occurs 16 / 26

Split the mantissa so that no rounding error occurs during inner product computation

 $[a_1 \ a_2 \ \cdots \ a_k] \text{ is a row vector of matrix } \mathbf{A}.$





18/26



- Compute high-precision GEMM by the accumulation of lower-precision GEMMs
- 2 Can be applied to arbitrary precision GEMM

- Weak point

- **1** Requires additional memory space to store split matrices (typically large)
- 2 More splits are required when the exponent distribution is large

Ozaki scheme on FP16 Tensor Cores (by Mukunoki et al.)

■ Use FP16 Tensor Cores for each matrix multiplication



Mukunoki et al. DGEMM Using Tensor Cores, and Its Accurate and Reproducible Versions, ISC 2020

20/26

Ozaki scheme on FP16 Tensor Cores (by Mukunoki et al.)



- $\blacksquare \text{ Larger } \phi$
 - \Rightarrow larger exponent distribution
 - \Rightarrow more splitting is required
 - $\Rightarrow {\rm throughput} \ {\rm degradation}$
- Faster than cublasDgemm

Problem

- FP16 is inefficient for the Ozaki scheme.
- The effective mantissa length is 5 ~ 7 bits, whereas FP16 has 11 bit mantissa.

Ozaki scheme on Integer Tensor Cores



■ a⁽ⁱ⁾ · b^(j) can be calculated by integer operations
■ A⁽ⁱ⁾ · B^(j) can be calculated by integer Tensor Cores!

Ozaki scheme on Integer Tensor Cores (Our work)

- Use Int8 Tensor Cores for each matrix multiplication
 - Int8-input Int32-accumulation



Pros/Cons of using Integer Tensor Cores in the Ozaki scheme

⊂ Pros

- Less working memory usage than FP16 Tensor Cores
- Int8 Tensor Cores have higher throughput than FP16
- Higher efficiency in the actual mantissa usage
- Fewer computations are required than FP16 to obtain the same accuracy

etc

Cons No significant Cons.

Details are in the paper.

Quantum circuit simulation using DGEMM emu.



NVIDIA RTX A6000 AdaThe accuracy is almost the same

25/26

Summary

- Question

Can deep learning processors be used for HPC applications?

\Rightarrow Yes

- We can compute double-precision equivalent GEMM using integer matrix multiplication unit.
- More research is required to put the Ozaki scheme into practice.
 - because the Ozaki scheme does not emulate the floating-point arithmetic, whereas most HPC applications are conducted by floating-point.