

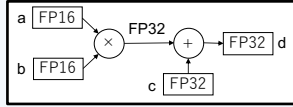
Accurate DGEMM using Tensor Cores

Daichi Mukunoki¹, Katsuhisa Ozaki², Takeshi Ogita³

1 RIKEN Center for Computational Science (JAPAN), 2 Shibaura Institute of Technology (JAPAN), 3 Tokyo Woman's Christian University (JAPAN)

Introduction

- As the demand for deep-learning increases, specialized hardware dedicated to low-precision matrix-multiplication, which is the kernel of those tasks, has been developing (e.g., NVIDIA Tensor Core, Google TPU, ARM Matherhorn's Matmul)
- NVIDIA Tensor Cores** are a special processing unit that enables 4×4 matrix multiplication operations on FP16 inputs with FP32 precision and return the result on FP32. Besides, cuBLAS provides the cublasGemmEx routine using the above Tensor Core operation.
- Markidis et al. [1] proposed a method to achieve nearly the SGEMM equivalent accuracy on the matrices with the dynamic range supported on FP16 using the cublasGemmEx routine.
- In this poster, we present an accurate DGEMM implementation, which has the DGEMM compatible interface and can be computed using cublasGemmEx on Tensor Cores
- Our implementation returns the **correctly-rounded** result: it supports bit-level **reproducibility** and is compatible with other correctly-rounded GEMMs such as ExBLAS [2] and OzBLAS [3]



Tensor Core operation with FP32 accuracy

Data and computation precisions on GEMM routines

	Mat. A	Mat. B	Mat. C	Computation
cublasDgemm	FP64	FP64	FP64	FP64
cublasGemmEx (with Tensor Cores†)	FP16	FP16	FP32	FP32
our DGEMM	FP64	FP64	FP64	Correctly-rounded

†Other operations and data formats are also supported through the same interface without Tensor Cores

Method (Ozaki scheme)

- Our method is based on the Ozaki scheme [4], which is an accurate matrix-multiplication algorithm based on the error-free transformation for dot-product/matrix-multiplication.
- In this study, we modified the original version of the Ozaki scheme, which was designed to be computed using DGEMM, to be computed using Tensor Cores on FP64 input/output.

Ozaki scheme for Tensor Cores

```
function (C ← DGEMM-TC-CR(m, n, k, A, B)) // A ∈ ℝm×k, B ∈ ℝk×n
    (Asplit[sA], cA[sA]) ← SplitA(m, k, A) // Asplit is obtained on FP16
    (Bsplit[sB], cB[sB]) ← SplitB(k, n, B) // Bsplit is obtained on FP16
    r = 1
    for (q = 1 : sB) do
        for (p = 1 : sA) do
            Ctmp1 = GEMMFP32(m, n, k, Asplit[p], Bsplit[q])
            Ctmp2[r] = ScaleUp(Ctmp1, cA[p] × cB[q])
            r = r + 1
        end for
    end for
    C = NearSum(Ctmp2, r)
end function
```

Spitting for a vector (for matrix multiply, it is performed for inner-product wise direction)

```
function ((xsplit[sx], c[sx]) ← Split(n, x))
    ρ = ceil(log2(uFP64-1) - (log2(uFP32-1) - log2(n)) / 2)
    μ = max0 ≤ i ≤ n-1 (|xi|)
    j = 0
    while (μ ≠ 0) do
        j = j + 1
        τ = ceil(log2(μ))
        σ = 2ρ+τ
        xtmp = f1FP64((x + σ) - σ) // xtmp is the split vector on FP64
        x = f1FP64(x - xsplit[j])
        θ = floor(log2(65504) - τ) // 65504 is maximum of the FP16 values
        c[j] = 2θ // this must be hold for scaling-up later
        xsplit[j] = f1FP16(xtmp[c[j]]) // Scaling-down and casting from FP64 to FP16
        μ = max0 ≤ i ≤ n-1 (|xi|)
    end while
    sx = j
end function
```

NearSum [5] computes the correctly-rounded summation

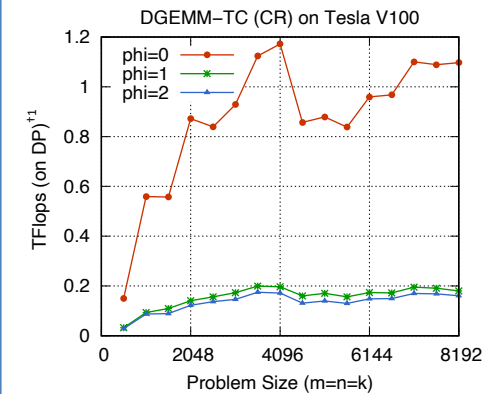
Note: In this method, the number of split matrices depends on the range of the absolute values of the input matrices

ScaleUp computes element-wise scalar-multiplication of two matrices

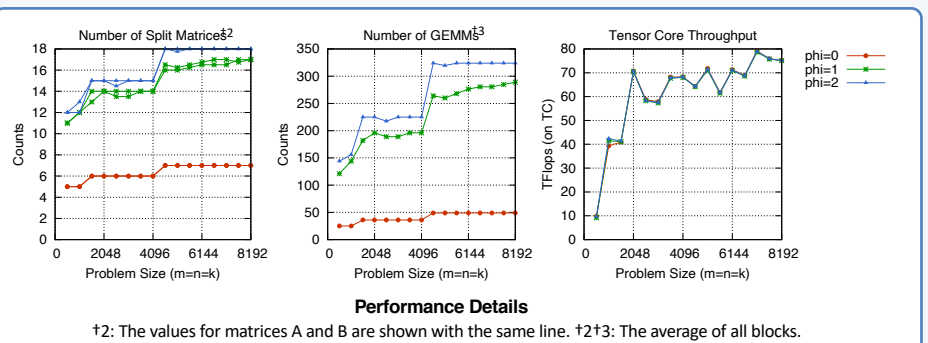
u_{FP64} denotes the unit round-off of FP64: u_{FP64} = 2⁻⁵³
u_{FP32} denotes the unit round-off of FP32: u_{FP32} = 2⁻²⁴

Performance on Tesla V100

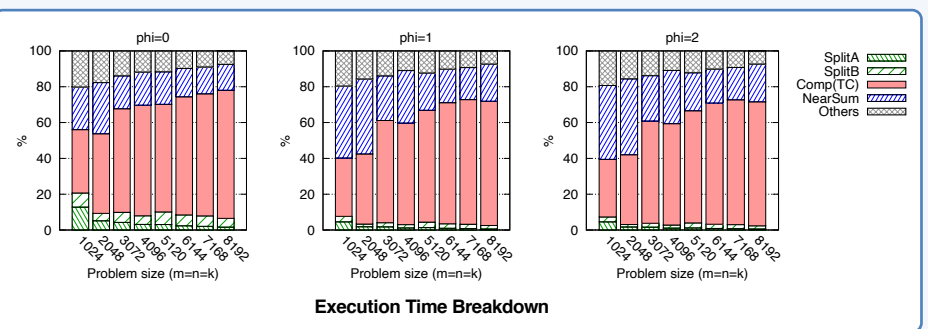
- Environment: NVIDIA Tesla V100 GPU (PCIe-32GB): the theoretical peak performance (with 1.38 GHz) is 7065.6 GFlops on FP64 and 113049.6 GFlops on Tensor Cores. CUDA 10.1, GPU driver version 418.39.
- Performance depends on the input (the number of split matrices depends on the range of the absolute values of the input matrices): matrices are initialized with (rand-0.5) × exp(φ × randn) to vary the floating-point range: 9.3E-10 – 5.0E-01 on φ = 0, 1.5E-09 – 1.6E+02 on φ = 1, and 1.3E-09 – 4.8E+04 on φ = 2.



†1: "Flops (on DP)" is the number of FP64 floating-point operations, corresponding to the standard DGEMM, per second



†2: The values for matrices A and B are shown with the same line. †3: The average of all blocks.



Conclusion

- We presented an accurate GEMM implementation with the DGEMM compatible interface using cublasGemmEx performed on Tensor Cores, which returns the correctly-rounded result on FP64
- The performance depends on the input matrices. For example, for matrices initialized with random numbers having the dynamic range of 1E+9, we achieved approx. 1.2 TFlops.
- Although our implementation brings no performance advantage against cuBLAS DGEMM on GPUs that support fast FP64 (1/16 of Tensor Cores), it can be beneficial on hardware with limited FP64 support such as NVIDIA Tesla T4, whose FP64 performance is 1/256 of Tensor Cores.
- Our approach opens the way to utilize AI-oriented hardware, which supports limited FP64 performance, for more general purposes and may impact the system co-design.
- The latest work including some extensions will be published soon.

Reference and Acknowledgement

- Markidis, S., Chien, S.W.D., Laure, E., Peng, I.B., Vetter, J.S.: NVIDIA Tensor Core Programmability, Performance Precision. In: 2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), pp. 522–531 (2018)
- S. Collange, D. Defour, S. Graillat, R. Iakymchuk: Numerical reproducibility for the parallel reduction on multi- and many-core architectures, Parallel Computing, Vol. 49, pp. 83–97, 2015.
- Mukunoki, D., Ogita, T., Ozaki, K.: Reproducible BLAS Routines with Tunable Accuracy Using Ozaki Scheme for Many-core Architectures. In: 13th International Conference on Parallel Processing and Applied Mathematics (PPAN2019) (2019), (to appear)
- K. Ozaki, T. Ogita, S. Oishi, S. M. Rump: "Error-free transformations of matrix multiplication by using fast routines of matrix multiplication and its applications," Numerical Algorithms, Vol. 59, No. 1, pp. 95–118, 2012.
- S. Rump, T. Ogita, S. Oishi, "Accurate floating-point summation part II: Sign, k-fold faithful and rounding to nearest," SIAM Journal on Scientific Computing, Vol. 31, No. 2, pp. 1269–1302, 2008.

* This research was partially supported by MEXT as "Exploratory Issue on Post-K computer" (Development of verified numerical computations and super-high-performance computing environment for extreme researches) and the Japan Society for the Promotion of Science (JSPS) KAKENHI Grant Number 19K20286. This research used computational resources of Cygnus provided by Multidisciplinary Cooperative Research Program in Center for Computational Sciences, University of Tsukuba.