

# Performance Enhancement of the Ozaki Scheme on Integer Matrix Multiplication Unit

Journal Title  
XX(X):1–12  
©The Author(s) 2016  
Reprints and permission:  
sagepub.co.uk/journalsPermissions.nav  
DOI: 10.1177/ToBeAssigned  
www.sagepub.com/

SAGE

Yuki Uchino<sup>1</sup>, Katsuhisa Ozaki<sup>2</sup>, and Toshiyuki Imamura<sup>1</sup>

## Abstract

This study was aimed at simultaneously achieving sufficient accuracy and high performance for general matrix multiplications. Recent architectures, such as NVIDIA GPUs, feature high-performance units designed for low-precision matrix multiplications in machine learning models, and next-generation architectures are expected to follow the same design principle. The key to achieving superior performance is to fully leverage such architectures. The Ozaki scheme, a highly accurate matrix multiplication algorithm using error-free transformations, enables higher-precision matrix multiplication to be performed through multiple lower-precision matrix multiplications and higher-precision matrix additions. Ootomo et al. implemented the Ozaki scheme on high-performance matrix multiplication units with the aim of achieving both sufficient accuracy and high performance. This paper proposes alternative approaches to improving performance by reducing the numbers of lower-precision matrix multiplications and higher-precision matrix additions. Numerical experiments demonstrate the accuracy of the results and conduct performance benchmarks of the proposed approaches. These approaches are expected to yield more efficient results in next-generation architectures.

## Keywords

matrix multiplication, fixed-point arithmetic, floating-point arithmetic, Tensor Cores, error-free transformation

## 1 Introduction

The field of machine learning, including AI, is evolving daily, and the scale and complexity of machine learning models are continually increasing. Recent architectures are designed to process these models rapidly and with high energy efficiency. For machine learning models, matrix multiplications using low-precision floating-point systems and integers are essential. Therefore, recent architectures are equipped with high-performance low-precision floating-point and integer matrix multiplication units, and high-performance mixed-precision matrix multiplication units that leverage these capabilities. A prime example of this is the NVIDIA Tensor Cores (see [NVIDIA Corporation \(2024\)](#)). Table 1 shows the specifications of the NVIDIA GPUs equipped with Tensor Core technology. Note that the specification of FP16 TC on RTX 4090 (165 TFLOPS) is for FP16 input and FP32 output, and the specifications for the H100 and H200 are the same as those for the GH200. In the future, numerical computation algorithms that leverage the performance of cutting-edge architectures will be essential. This study focused on researching high-performance matrix multiplication algorithms that maximize the potential of the latest architectures.

A highly accurate matrix multiplication scheme via the error-free transformation of matrix products was proposed in [Ozaki et al. \(2012\)](#). The scheme is called the Ozaki scheme and it converts a matrix product into a sum of multiple matrix products. It is also possible to convert a matrix product into a sum of lower-precision matrix products to take full advantage of the immense computational power of recent architectures. In order to compute  $D \leftarrow AB$  for higher-precision matrices  $A \in \mathbb{R}^{m \times n}$  and  $B \in \mathbb{R}^{n \times p}$ , the Ozaki

**Table 1.** Specifications in TFLOPS/TOPS of NVIDIA GPUs for dense data [NVIDIA Corporation \(2024\)](#)

|         | RTX 4090 | A100 | GH200 | B200 | GB200 |
|---------|----------|------|-------|------|-------|
| FP64    | 1.29     | 9.7  | 34    | 40   | 90    |
| FP64 TC | –        | 19.5 | 67    | 40   | 90    |
| FP32    | 82.6     | 19.5 | 67    | 80   | 180   |
| TF32 TC | 82.6     | 156  | 494   | 1100 | 5000  |
| BF16 TC | 165      | 312  | 989   | 2250 | 5000  |
| FP16 TC | 165      | 312  | 989   | 2250 | 5000  |
| INT8 TC | 661      | 624  | 1979  | 4500 | 10000 |
| FP8 TC  | 661      | –    | 1979  | 4500 | 10000 |
| FP6 TC  | –        | –    | –     | 4500 | 10000 |
| INT4 TC | 1321     | –    | –     | –    | –     |
| FP4 TC  | –        | –    | –     | 9000 | 20000 |

scheme using lower/mixed-precision, provided in [Mukunoki et al. \(2020\)](#), is constructed as the following four steps:

- Extract the lower-precision matrices  $A_1, A_2, \dots, A_k$  from the higher-precision matrix  $A$ , where  $k$  is specified by the user, shifting  $A_i$  to prevent overflow and underflow at lower-precision arithmetic.

<sup>1</sup>RIKEN Center for Computational Science, Japan

<sup>2</sup>Department of Mathematical Sciences, Shibaura Institute of Technology, Japan

## Corresponding author:

Yuki Uchino, 7-1-26 Minatojima-minami-machi, Chuo-ku, Kobe, Hyogo 650-0047, Japan

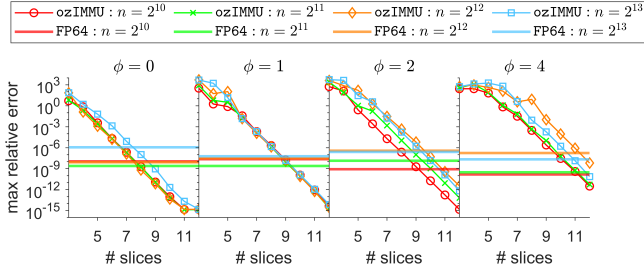
Email: yuki.uchino.fe@riken.jp

- (ii) Extract  $k$  lower-precision matrices  $B_i$  from  $B$  in the similar way as  $A_i$ .
- (iii) Compute  $A_i B_j$  for  $i + j \leq k + 1$  using lower/mixed-precision arithmetic.
- (iv) Reverse (shift)  $A_i B_j$  and accumulate the results into  $D$  using higher-precision arithmetic.

When emulating the GEMM routine, the following fifth step is also performed:

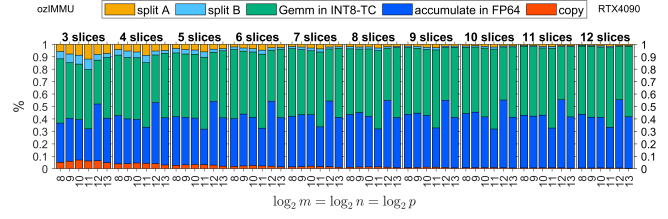
- (v) Compute  $C \leftarrow \alpha D + \beta C$  for scalars  $\alpha, \beta$  and a matrix  $C$ .

Ootomo et al. implemented the Ozaki scheme using the INT8 Tensor Cores and evaluated the performance for emulating DGEMM, an FP64 matrix multiplication routine, as reported in Ootomo et al. (2024). The Ozaki scheme implemented by Ootomo is named the ‘‘Ozaki Scheme on Integer Matrix Multiplication Unit’’ (ozIMMU), and the code is available in Ootomo (2024). Figure 1 shows the accuracy of ozIMMU. Herein,  $\phi$  specifies the tendency of the difficulty in terms of the accuracy of matrix multiplications. At least 7 or 8 slices are required to obtain sufficiently accurate results, even for well-conditioned matrix multiplications. As  $\phi$  increases, more slices are required to obtain sufficient accuracy. Ootomo’s implementation for splitting matrices offers bit masking. Thus, the extracted matrices  $A_i, B_i$  may not be optimal and the splitting method can be improved (see Section 3.1 for details). Improvement of the splitting method should contribute to improving the accuracy of results.

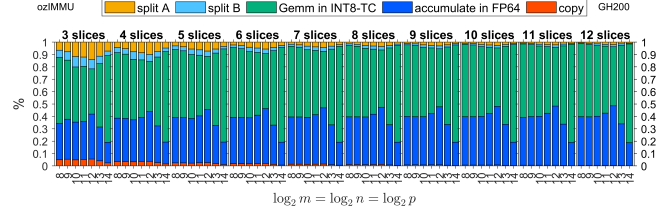


**Figure 1.** Accuracy of ozIMMU. Matrix  $A$  has entries  $a_{ij} := (U_{ij} - 0.5) \cdot \exp(\phi \cdot N_{ij})$ , where  $U_{ij} \in (0, 1)$  are uniformly distributed and  $N_{ij}$  are drawn from standard normal distribution for  $1 \leq i, j \leq m$  and  $m = n = p$ . Matrix  $B$  is composed similarly.

Figures 2 and 3 show the time breakdown of ozIMMU for double-precision matrices  $A, B$ , defined in IEEE Computer Society (2019), on GeForce RTX 4090 and GH200 Grace Hopper Superchip, respectively. Note that ‘‘split A’’, ‘‘split B’’, ‘‘Gemm in INT8-TC’’, ‘‘accumulation in FP64’’, and ‘‘copy’’ correspond to the steps (i), (ii), (iii), (iv), and (v) in the Ozaki scheme for emulating the GEMM routine, respectively. It can be seen that the computing time for the accumulation of matrix products in FP64 is not negligible even though the computation cost is  $\mathcal{O}(mp)$  operations. This is because the performance of the INT8 Tensor Core is nearly 512 times and 60 times higher than that of FP64 on the RTX 4090 and GH200, respectively. Because this ratio increases significantly on the B200 and future architectures, accelerating the accumulation process is critical for the Ozaki scheme.



**Figure 2.** Time breakdown of ozIMMU on NVIDIA GeForce RTX 4090



**Figure 3.** Time breakdown of ozIMMU on NVIDIA GH200 Grace Hopper Superchip

In this paper, we propose new implementation methods for accelerating the accumulation of the Ozaki scheme using the INT8 Tensor Core. In addition, an alternative splitting method is also applied to improve the accuracy of results. For the number of slices, the splitting method contributes to obtaining more accurate results than those of Ootomo’s ozIMMU. The remainder of this paper is organized as follows: Section 2 overviews previous studies about the Ozaki scheme using an integer matrix multiplication unit; Section 3 presents the proposed method for accelerating the accumulation and optimizing the splitting method; Section 4 shows numerical results to illustrate the efficiency of the proposed methods; Section 5 provides a rounding error analysis of the Ozaki scheme with the proposed methods; and Section 6 presents final remarks.

## 2 Previous study

In this section, we briefly summarize previous studies. Let  $u$  be the relative error unit, e.g.,  $u = 2^{-53}$  for double-precision floating-point numbers defined in IEEE 754. Define  $\mathbb{F}$  as a set of binary floating-point numbers with  $u$ . Let  $O_{m,n}$  be an  $m \times n$  zero matrix. Suppose that  $A \in \mathbb{F}^{m \times n}$  and  $B \in \mathbb{F}^{n \times p}$ . Ozaki et al. proposed an error-free transformation of a matrix product  $AB$  in Ozaki et al. (2012). For a user-specified constant  $k \in \mathbb{N}$ , the Ozaki scheme transforms each of  $A$  and  $B$  into  $k + 1$  matrices

$$\begin{aligned} A &:= A_1 + A_2 + \cdots + A_k + V_k, \\ B &:= B_1 + B_2 + \cdots + B_k + W_k, \end{aligned}$$

where  $A_i \in \mathbb{F}^{m \times n}$ ,  $B_i \in \mathbb{F}^{n \times p}$ ,

$$V_i := A - \sum_{j=1}^i A_j, \quad \text{and} \quad W_i := B - \sum_{j=1}^i B_j, \quad (1)$$

such that  $V_i \in \mathbb{F}^{m \times n}$  and  $W_i \in \mathbb{F}^{n \times p}$  for  $i \leq k$ , and  $|(A_s)_{ij}| \geq |(A_{s+1})_{ij}|$  holds if  $(A_s)_{ij} \neq 0$  for  $s = 1, \dots, k-2$ , and  $B_s$  satisfies the corresponding relation. If the technique provided in Minamihata et al.

(2016) is used, then  $|(A_s)_{ij}| > |(A_{s+1})_{ij}|$  holds when  $(A_s)_{ij} \neq 0$ .

Let  $\mathbb{F}_N$  be a set of  $N$ -bit binary floating-point numbers. For the Ozaki scheme using the FP16 Tensor Core with accumulation in FP32 and outputs in FP32, the mixed-precision splitting method with shifting to prevent overflow and underflow was utilized in Mukunoki et al. (2020). For  $A \in \mathbb{F}_{64}^{m \times n}$  and  $B \in \mathbb{F}_{64}^{n \times p}$ , the splitting method produces

$$A =: \text{diag}(\mu'^{(1)})A'_1 + \cdots + \text{diag}(\mu'^{(k)})A'_k + V_k, \quad (2)$$

$$B =: B'_1 \text{diag}(\nu'^{(1)}) + \cdots + B'_k \text{diag}(\nu'^{(k)}) + W_k \quad (3)$$

with double-precision shift values  $\mu'^{(i)} \in \mathbb{F}_{64}^m$ ,  $\nu'^{(i)} \in \mathbb{F}_{64}^p$  and half-precision matrices  $A'_i \in \mathbb{F}_{16}^{m \times n}$ ,  $B'_i \in \mathbb{F}_{16}^{n \times p}$  for  $i = 1, \dots, k$ . Algorithm 1 represents the splitting method for obtaining (2) using Minamihata's technique. Equation (3) can be obtained by transposing the results of Algorithm 1 executed for  $B^T$ . Algorithm 1 can be described as a loop that is executed until  $s = k$  for simplicity; however, the loop terminates when  $A = O_{m,n}$  in practice. Note that the binary logarithm is used at the 3rd line in Algorithm 1; however, the calculation using the binary logarithm occasionally returns erroneous results. Therefore, it is better to use a calculation method without the binary logarithm, such as a bitwise operation or a technique leveraging rounding error in floating-point arithmetic as follows:

$$\mu_i'^{(s)} \leftarrow u^{-1}\alpha_i + (1 - u^{-1})\alpha_i,$$

where  $\alpha_i := \max_j |a_{ij}|$ . This technique was developed by Rump and provided in Ozaki et al. (2013).

---

**Algorithm 1** Mixed-precision splitting method from Mukunoki et al. (2020) for Ozaki scheme between FP64 and FP16 using floating-point arithmetic in round-to-nearest-even mode with Minamihata's technique from Minamihata et al. (2016)

---

**Input:**  $A \in \mathbb{F}_{64}^{m \times n}$ ,  $k \in \mathbb{N}$

**Output:**  $A'_s \in \mathbb{F}_{16}^{m \times n}$ ,  $\mu'^{(s)} \in \mathbb{F}_{64}^m$ ,  $s = 1, \dots, k$

```

1:  $\beta \leftarrow \lceil (29 - \log_2 n)/2 \rceil$ 
2: for  $s = 1, \dots, k$  do
3:    $\mu_i'^{(s)} \leftarrow 2^{\lceil \log_2 \max_j |a_{ij}| \rceil}$  for  $i = 1, \dots, m$ 
4:    $\sigma_i \leftarrow 0.75 \cdot \mu_i'^{(s)} \cdot 2^\beta$  for  $i = 1, \dots, m$ 
5:    $(A_s)_{ij} \leftarrow (a_{ij} + \sigma_i) - \sigma_i \ \forall i, j$  {Extract in
      round-to-nearest-even mode}
6:    $(A'_s)_{ij} \leftarrow \text{FP16}((\mu_i'^{(s)})^{-1} \cdot (A_s)_{ij}) \ \forall i, j$  {Convert to
      FP16}
7:    $a_{ij} \leftarrow a_{ij} - (A_s)_{ij} \ \forall i, j$ 
8: end for
```

---

After splitting matrices as in (2) and (3),

$$\begin{aligned}
AB &= \sum_{s+t \leq k+1} \text{diag}(\mu'^{(s)})A'_s B'_t \text{diag}(\nu'^{(t)}) \\
&+ \sum_{s+t > k+1} \text{diag}(\mu'^{(s)})A'_s B'_t \text{diag}(\nu'^{(t)}) \\
&+ \sum_{s=1}^k \text{diag}(\mu'^{(s)})A'_s W_k \\
&+ \sum_{s=1}^k V_k B'_s \text{diag}(\nu'^{(s)}) \\
&+ V_k W_k
\end{aligned}$$

holds and there is no rounding error in  $A'_s B'_t$  for  $1 \leq s, t \leq k$  on the FP16 Tensor Core. In Mukunoki et al. (2020), the approximation of  $AB$  was computed only using matrix multiplications with the FP16 Tensor Core, as

$$AB \approx \sum_{s+t \leq k+1} \text{diag}(\mu'^{(s)})A'_s B'_t \text{diag}(\nu'^{(t)})$$

shown in Algorithm 2. This method was referred to as the “Fast Mode” in Mukunoki et al. (2020). For larger  $k$ ,  $|(W_k)_{ij}|$  and  $|(V_k)_{ij}|$  are smaller and the approximation is more accurate.

---

**Algorithm 2** Mixed-precision matrix multiplication method for Ozaki scheme using FP16 Tensor Core in fast mode from Mukunoki et al. (2020)

---

**Input:**  $A'_s \in \mathbb{F}_{16}^{m \times n}$ ,  $\mu'^{(s)} \in \mathbb{F}_{64}^m$ ,  $B'_s \in \mathbb{F}_{16}^{n \times p}$ ,  $\nu'^{(s)} \in \mathbb{F}_{64}^n$ ,  $s = 1, \dots, k$

**Output:**  $C \in \mathbb{F}_{64}^{m \times p}$

```

1:  $C \leftarrow O_{m,p}$  { $C' \in \mathbb{F}_{32}^{m \times p}$ }
2: for  $g = 2, \dots, k+1$  do
3:   for  $s = 1, \dots, g-1$  do
4:      $C' \leftarrow A'_s B'_{g-s}$  {Compute using GEMM on FP16
      Tensor Core}
5:      $C \leftarrow C + \text{diag}(\mu'^{(s)})\text{FP64}(C')\text{diag}(\nu'^{(t)})$ 
      {Compute in FP64}
6:   end for
7: end for
```

---

Let  $\mathbb{I}_N$  be a set of  $N$ -bit signed integers. It holds that  $-2^{N-1} \leq i \leq 2^{N-1} - 1$  for all  $i \in \mathbb{I}_N$ . Remember that no error occurs in integer arithmetic, barring overflow. For the Ozaki scheme using the INT8 Tensor Core with accumulation in INT32, the mixed-precision splitting method via bit masking with shifting is offered in Ootomo et al. (2024) and provided in Ootomo (2024). Algorithm 3 shows the splitting method. Let

$$\beta := \min \left( 7, \left\lfloor \frac{31 - \log_2 n}{2} \right\rfloor \right). \quad (4)$$

Assume that  $\beta \geq 1$ , i.e.,  $n \leq 2^{29}$ . That splitting method produces

$$A =: \text{diag}(\mu'')(2^{-\beta+1}A''_1 + \cdots + 2^{-k\beta+1}A''_k) + V_k, \quad (5)$$

$$B =: (2^{-\beta+1}B''_1 + \cdots + 2^{-k\beta+1}B''_k)\text{diag}(\nu'') + W_k \quad (6)$$

**Algorithm 3** Mixed-precision splitting method for Ozaki scheme via bit masking from [Ootomo et al. \(2024\)](#). This algorithm is specialized for the Ozaki scheme using the INT8 Tensor Core.

**Input:**  $A \in \mathbb{F}_{64}^{m \times n}$ ,  $k \in \mathbb{N}$

**Output:**  $A''_s \in \mathbb{F}_8^{m \times n}$ ,  $\mu'' \in \mathbb{F}_{64}^m$ ,  $s = 1, \dots, k$

- 1:  $\beta \leftarrow \min(7, \lfloor (31 - \log_2 n)/2 \rfloor)$
- 2:  $\mu''_i \leftarrow 2^{\lfloor \log_2 \max_j |a_{ij}| \rfloor}$  for  $i = 1, \dots, m$
- 3: **for**  $s = 1, \dots, k$  **do**
- 4: Extract sth  $\beta$  bits of mantissa of  $a_{ij} \forall i, j$  via bit masking and hold those as INT8 variable  $(A''_s)_{ij}$
- 5: **end for**

with double-precision shift values  $\mu'' \in \mathbb{F}_{64}^m$ ,  $\nu'' \in \mathbb{F}^p$  and 8-bit integer matrices  $A''_i \in \mathbb{F}_8^{m \times n}$ ,  $B''_i \in \mathbb{F}_8^{n \times p}$  for  $i = 1, \dots, k$ .

After splitting matrices as in (5) and (6),

$$\begin{aligned}
 AB &= \sum_{s+t \leq k+1} \text{diag}(\mu'') 2^{-\beta s+1} 2^{-\beta t+1} A''_s B''_t \text{diag}(\nu'') \\
 &+ \sum_{s+t > k+1} \text{diag}(\mu'') 2^{-\beta s+1} 2^{-\beta t+1} A''_s B''_t \text{diag}(\nu'') \\
 &+ \sum_{s=1}^k \text{diag}(\mu'') 2^{-\beta s+1} A''_s W_k \\
 &+ \sum_{s=1}^k V_k 2^{-\beta s+1} B''_s \text{diag}(\nu'') \\
 &+ V_k W_k
 \end{aligned}$$

holds and there is no error in  $A''_s B''_t$  for  $1 \leq s, t \leq k$  on the INT8 Tensor Core barring overflow. In [Ootomo et al. \(2024\)](#), the approximation of  $AB$  is computed using only matrix multiplications with the INT8 Tensor Core, as

$$AB \approx \sum_{s+t \leq k+1} \text{diag}(\mu'') 2^{-\beta s+1} 2^{-\beta t+1} A''_s B''_t \text{diag}(\nu''), \quad (7)$$

shown in Algorithm 4. For larger  $k$ ,  $|(W_k)_{ij}|$  and  $|(V_k)_{ij}|$  are smaller and the approximation is more accurate.

**Algorithm 4** Mixed-precision matrix multiplication method for Ozaki scheme using INT8 Tensor Core from [Ootomo et al. \(2024\)](#)

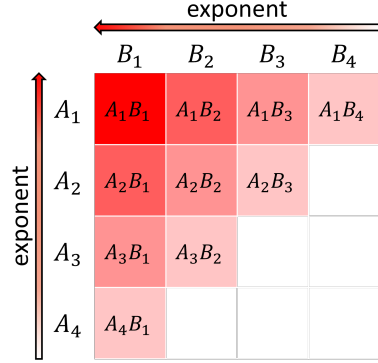
**Input:**  $A''_s \in \mathbb{F}_8^{m \times n}$ ,  $\mu''^{(s)} \in \mathbb{F}_{64}^m$ ,  $B''_s \in \mathbb{F}_8^{n \times p}$ ,  $\nu''^{(s)} \in \mathbb{F}_{64}^p$ ,  $s = 1, \dots, k$

**Output:**  $C \in \mathbb{F}_{64}^{m \times p}$

- 1:  $\beta \leftarrow \min(7, \lfloor (31 - \log_2 n)/2 \rfloor)$
- 2:  $C \leftarrow O_{m,p}$   $\{C'' \in \mathbb{F}_{32}^{m \times p}\}$
- 3: **for**  $g = 2, \dots, k+1$  **do**
- 4: **for**  $s = 1, \dots, g-1$  **do**
- 5:  $C'' \leftarrow A''_s B''_{g-s}$  {Compute using GEMM on INT8 Tensor Core}
- 6:  $C \leftarrow C + \text{diag}(\mu'') 2^{-\beta s+1} 2^{-\beta(g-s)+1} \text{FP64}(C'') \text{diag}(\nu'')$  {Compute in FP64}
- 7: **end for**
- 8: **end for**

Figure 4 represents images of Algorithms 2 and 4. In Figure 4,  $A_i := \text{diag}(\mu^{(i)}) A'_i$  and  $B_i := B'_i \text{diag}(\nu^{(i)})$  for

Algorithms 1 and 2, and  $A_i := \text{diag}(\mu'') 2^{-i\beta+1} A'_i$  and  $B_i := 2^{-i\beta+1} B'_i \text{diag}(\nu'')$  for Algorithms 3 and 4. The exponent of  $A_i B_j$  is larger for smaller indices of  $i, j$ , so those terms have stronger effects on the accuracy of the final result.



**Figure 4.** Images of matrix multiplications in Ozaki scheme for  $k = 4$

### 3 Proposed methods

#### 3.1 Splitting with rounding to nearest

We begin by explaining our motivation with a simple example. Algorithm 3 described in Section 2 used bitmask for the splitting. We assume that each number keeps three bits. Let a floating-point number be  $x := 377$ , represented by  $(10101111)_2$  in binary format. The bitmask strategy divides  $x$  into three numbers  $x = x_1 + x_2 + x_3$ :

$$\begin{cases}
 x_1 = (\underline{1}01000000)_2, \\
 x_2 = (0000\underline{1}1000)_2, \\
 x_3 = (000000\underline{1}11)_2.
 \end{cases} \quad (8)$$

In decimal format,

$$x_1 = 256 + 0 + 64, \quad x_2 = 0 + 16 + 8, \quad x_3 = 4 + 2 + 1.$$

Here  $x_1$  is not the nearest number to  $x$ .

On the other hand, let  $x_1$  and  $x_2$  be the nearest number to  $x$  and  $x - x_1$ , respectively. Then,

$$\begin{cases}
 x_1 = (\underline{1}01000000)_2, \\
 x_2 = (000\underline{1}00000)_2, \\
 x_3 = -(00000000\underline{1})_2.
 \end{cases} \quad (9)$$

In decimal format,

$$x_1 = 256 + 0 + 64, \quad x_2 = 32, \quad x_3 = -1.$$

There are two drawbacks in the bitmask strategy. One drawback is that even if the leading bit is zero, we keep it, so that the number of significant figures decreases as in  $x_2$  in (8). The other drawback is that the truncation error can be reduced. For example, if we set two slices, there is a difference in the truncation error  $|x_3|$  in (8) and (9). It is expected that the nearest strategy makes the truncation error small. Summarizing, even if we obtain a computed result with  $k$  slices using the rounding to nearest strategy, the



accuracy may be comparable to that with  $k + 1$  slices using the bitmask strategy described in Section 2 as the previous study. In this case, we must carefully monitor the following via numerical examples:

- the increase in the splitting cost because the cost for the bitmask is low;
- the reduction in the number of matrix multiplications; and
- the reduction in the term for the accumulation.

Note that in the worst case, the results of the splitting with bitmask and the nearest strategy are the same. It can be explained using the following example:

$$x := (100010001)_2.$$

For the bit mask strategy, we have

$$\begin{cases} x_1 = (\underline{100000000})_2, \\ x_2 = (0000\underline{10000})_2, \\ x_3 = (00000000\underline{1})_2. \end{cases}$$

For the rounding to the nearest strategy, we obtain

$$\begin{cases} x_1 = (\underline{100000000})_2, \\ x_2 = (0000\underline{10000})_2, \\ x_3 = (00000000\underline{1})_2. \end{cases}$$

The following is an algorithm for the splitting with the round to nearest strategy.

---

**Algorithm 5** Proposed method for mixed-precision splitting for Ozaki scheme between FP64 and INT8 using floating-point arithmetic in round-to-nearest-even mode with Minamihata's technique from [Minamihata et al. \(2016\)](#)

---

**Input:**  $A \in \mathbb{F}_{64}^{m \times n}$ ,  $k \in \mathbb{N}$

**Output:**  $A'_s \in \mathbb{F}_8^{m \times n}$ ,  $\mu'^{(s)} \in \mathbb{F}_{64}^m$ ,  $s = 1, \dots, k$

- 1:  $\beta \leftarrow \min(7, \lfloor (31 - \log_2 n)/2 \rfloor)$
  - 2: **for**  $s = 1, \dots, k$  **do**
  - 3:  $\mu_i'^{(s)} \leftarrow 2^{\lceil \log_2 \max_j |a_{ij}| \rceil} \cdot 2^{1-\beta}$  for  $i = 1, \dots, m$
  - 4:  $\sigma_i \leftarrow 0.75 \cdot 2^{53} \cdot \mu_i'^{(s)}$  for  $i = 1, \dots, m$
  - 5:  $(A_s)_{ij} \leftarrow (a_{ij} + \sigma_i) - \sigma_i$  for  $i = 1, \dots, m$  {Extract in round-to-nearest-even mode}
  - 6:  $(A'_s)_{ij} \leftarrow \text{INT8}((\mu_i'^{(s)})^{-1} \cdot (A_s)_{ij}) \forall i, j$  {Convert to INT8}
  - 7:  $a_{ij} \leftarrow a_{ij} - (A_s)_{ij} \forall i, j$
  - 8: **end for**
- 

### 3.2 Group-wise error-free accumulation

Next, we propose a method for accelerating the accumulation in FP64 for ozIMMU. Algorithm 4 requires the computation of a sum of  $k(k + 1)/2$  FP64 matrices at the 6th line. The accumulation accounts for a large ratio of the total computation time of ozIMMU as shown in Figures 2 and 3. For this challenge, we propose a method for accelerating ozIMMU by reducing the number of additions in FP64.

Define  $\mathbb{G}_g \subset \mathbb{R}^2$  for  $g = 3, \dots, k + 1$  as

$$\mathbb{G}_g := \{(i, j) \mid i + j = g\}.$$

Recall that Algorithm 4 performs

$$C \leftarrow \sum_{t=2}^{k+1} \sum_{s=1}^{g-1} \text{diag}(\mu'') 2^{-\beta s+1} 2^{-\beta(g-s)+1} A''_s B''_{g-s} \text{diag}(\nu''). \quad (10)$$

Let  $(i_1, j_1), \dots, (i_{g-1}, j_{g-1}) \in \mathbb{G}_g$  such that  $i_s \neq i_t$  and  $j_s \neq j_t$  implies  $s \neq t$  for  $1 \leq s, t \leq g - 1$ , where  $g \in \{3, \dots, k + 1\}$ . Then, the inner sum of (10) can be expressed as follows:

$$\begin{aligned} & \sum_{s=1}^{g-1} \text{diag}(\mu'') 2^{-\beta i_s+1} 2^{-\beta j_s+1} A''_{i_s} B''_{j_s} \text{diag}(\nu'') \\ &= 2^{-\beta g+2} \text{diag}(\mu'') \left( \sum_{s=1}^{g-1} A''_{i_s} B''_{j_s} \right) \text{diag}(\nu'') \end{aligned} \quad (11)$$

which uses that  $i_s + j_s = g$  for all  $s = 1, \dots, g - 1$ . If no overflow occurs in  $A''_{i_s} B''_{j_s} + A''_{i_t} B''_{j_t}$  ( $1 \leq s, t \leq g - 1$ ) in INT32, the summation can be performed on the accumulator of GEMM on the INT8 Tensor Core. Therefore, applying the above transformation reduces the numbers of slow conversions from INT32 to FP64, scalings, and slow summations in FP64. Let  $r \in \mathbb{N}$  be

$$r := \max(1, 2^{31-2\beta-\lceil \log_2 n \rceil}). \quad (12)$$

Then, the summation

$$\sum_{s=1}^{\min(r, g-1)} A''_{i_s} B''_{j_s}$$

can be computed without error using GEMM on the INT8 Tensor Core. The summation of  $r$  instances of  $A''_{i_s} B''_{j_s}$  can be computed without error; however, the summation of all  $|\mathbb{G}_g|$  instances of  $A''_{i_s} B''_{j_s}$  cannot always be computed without error. The proof validating error-free summation is provided in Section 5. Finally, we present Algorithm 6, which has a reduced number of summands in FP64 accumulation.

Assuming  $r \geq k$  for simplify, all  $A''_i B''_j$  for  $(i, j) \in \mathbb{G}_g$  can be accumulated without overflow for  $g = 3, \dots, k + 1$ . Algorithm 7 shows the Ozaki scheme with groupwise error-free accumulation for  $r \geq k$ .

### 3.3 Combination of proposals

Next, we accelerate ozIMMU by reducing the number of summands in FP64 accumulation and improving the splitting method. For this purpose, we combine the methods proposed in Sections 3.1 and 3.2. In order to use group-wise error-free accumulation as in Algorithm 6,  $A$  and  $B$  need to be expressed as in (5) and (6). Thus, we determine

$$\begin{aligned} \mu_i'' & \leftarrow 2^{\lceil \log_2 \max_j |a_{ij}| \rceil} \cdot 2^{1-\beta}, \\ \sigma_i & \leftarrow 0.75 \cdot 2^{53-\beta(s-1)} \cdot \mu_i'' \end{aligned}$$

as the 3rd and 4th lines in Algorithm 5. Then,  $A_s$  are extracted using rounding to nearest floating-point arithmetic as

$$(A_s)_{ij} \leftarrow (a_{ij} + 2^{\beta(s-1)} \sigma_i) - 2^{\beta(s-1)} \sigma_i$$

for the constant  $2^{\beta(s-1)} \sigma_i$  with a common ratio of  $2^\beta$ . Finally, we obtain Algorithm 8. The shift values for  $A_s$

**Algorithm 6** Proposed method for groupwise error-free accumulation for Ozaki scheme using INT8 Tensor Core

**Input:**  $A_s'' \in \mathbb{I}_8^{m \times n}$ ,  $\mu''^{(s)} \in \mathbb{F}_{64}^m$ ,  $B_s'' \in \mathbb{I}_8^{n \times p}$ ,  $\nu''^{(s)} \in \mathbb{F}_{64}^p$ ,  $s = 1, \dots, k$   
**Output:**  $C \in \mathbb{F}_{64}^{m \times p}$

- 1:  $\beta \leftarrow \min(7, \lfloor (31 - \log_2 n)/2 \rfloor)$
- 2:  $r \leftarrow \max(1, 2^{31-2\beta-\lceil \log_2 n \rceil})$  {#addends for error-free accum.}
- 3:  $C \leftarrow O_{m,p}$
- 4: **for**  $g = 2, \dots, k+1$  **do**
- 5:    $q \leftarrow 0$
- 6:    $C'' \leftarrow O_{m,p}$   $\{C'' \in \mathbb{I}_{32}^{m \times p}\}$
- 7:   **for**  $s = 1, \dots, g-1$  **do**
- 8:      $q \leftarrow q+1$
- 9:      $C'' \leftarrow C'' + A_s'' B_{g-s}''$  {Compute using GEMM on INT8 Tensor Core}
- 10:    **if**  $q == r \parallel s == g-1$  **then**
- 11:      $C \leftarrow C + 2^{-\beta g+2} \text{diag}(\mu'') \text{FP64}(C'') \text{diag}(\nu'')$  {Compute in FP64}
- 12:      $q \leftarrow 0$
- 13:      $C'' \leftarrow O$
- 14:    **end if**
- 15: **end for**
- 16: **end for**

**Algorithm 7** Simple version of proposed method for groupwise error-free accumulation for Ozaki scheme using INT8 Tensor Core for  $r \geq k$

**Input:**  $A_s'' \in \mathbb{I}_8^{m \times n}$ ,  $\mu''^{(s)} \in \mathbb{F}_{64}^m$ ,  $B_s'' \in \mathbb{I}_8^{n \times p}$ ,  $\nu''^{(s)} \in \mathbb{F}_{64}^p$ ,  $s = 1, \dots, k$   
**Output:**  $C \in \mathbb{F}_{64}^{m \times p}$

- 1:  $\beta \leftarrow \min(7, \lfloor (31 - \log_2 n)/2 \rfloor) = 7$
- 2:  $C \leftarrow O_{m,p}$
- 3: **for**  $g = 2, \dots, k+1$  **do**
- 4:    $C'' \leftarrow O_{m,p}$   $\{C'' \in \mathbb{I}_{32}^{m \times p}\}$
- 5:   **for**  $s = 1, \dots, g-1$  **do**
- 6:      $C'' \leftarrow C'' + A_s'' B_{g-s}''$  {Compute using GEMM on INT8 Tensor Core}
- 7:   **end for**
- 8:    $C \leftarrow C + \text{diag}(\mu'') 2^{-\beta g+2} \text{FP64}(C'') \text{diag}(\nu'')$  {Compute in FP64}
- 9: **end for**

are  $\mu_i''^{(s)} \cdot 2^{\beta(s-1)}$ ; thus,  $\mu_i''^{(s)}$  are determined before the “for” statement and the shift values for  $A_s$  for  $s \geq 1$  are calculated by shifting  $\mu_i''^{(s)}$  by a constant ratio  $2^\beta$ . Hence, the algorithm finds the maximum absolute values  $\max_j |a_{ij}|$  for  $i = 1, \dots, m$  only once before the “for” statement. On the other hand, the shift values  $\mu_i''^{(s)}$  are determined inside the “for” statement in Algorithm 5. Therefore, the algorithm finds the maximum absolute values  $\max_j |a_{ij}|$  for  $i = 1, \dots, m$  at most  $k$  times.

## 4 Numerical examples

We have implemented the proposed methods in Uchino (2024), replacing specific code in Ootomo’s ozIMMU library with our code. All numerical experiments were conducted on NVIDIA GH200 Grace Hopper Superchip and NVIDIA

**Algorithm 8** Another proposed method for mixed-precision splitting for Ozaki scheme between FP64 and INT8 using floating-point arithmetic in round-to-nearest-even mode with Minamihata’s technique from Minamihata et al. (2016). The results can be expressed as in (5) for error-free group-wise accumulation.

**Input:**  $A \in \mathbb{F}_{64}^{m \times n}$ ,  $k \in \mathbb{N}$   
**Output:**  $A_s'' \in \mathbb{I}_8^{m \times n}$ ,  $\mu'' \in \mathbb{F}_{64}^m$ ,  $s = 1, \dots, k$

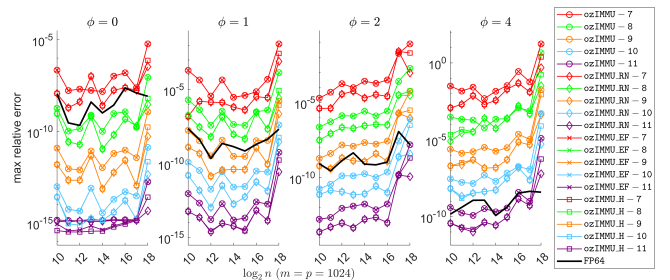
- 1:  $\beta \leftarrow \min(7, \lfloor (31 - \log_2 n)/2 \rfloor)$
- 2:  $\mu_i'' \leftarrow 2^{\lceil \log_2 \max_j |a_{ij}| \rceil} \cdot 2^{1-\beta}$  for  $i = 1, \dots, m$
- 3: **for**  $s = 1, \dots, k$  **do**
- 4:    $\sigma_i \leftarrow 0.75 \cdot 2^{53-\beta(s-1)} \cdot \mu_i''$  for  $j = 1, \dots, m$
- 5:    $(A_s)_{ij} \leftarrow (a_{ij} + \sigma_i) - \sigma_i \forall i, j$  {Extract in round-to-nearest-even mode}
- 6:    $(A_s'')_{ij} \leftarrow \text{INT8}((\mu_i'')^{(s)-1} \cdot 2^{-\beta(s-1)} \cdot (A_s)_{ij}) \forall i, j$  {Convert to INT8}
- 7:    $a_{ij} \leftarrow a_{ij} - (A_s)_{ij} \forall i, j$
- 8: **end for**

GeForce RTX 4090 with the GNU C++ Compiler 11.4.1 and NVIDIA CUDA Toolkit 12.5.82. The tested methods will be denoted as follows:

- ozIMMU- $k$  : Ootomo’s implementation with  $k$  slices
- ozIMMU\_RN- $k$ : Proposed method in Section 3.1 with  $k$  slices
- ozIMMU\_EF- $k$ : Proposed method in Section 3.2 with  $k$  slices
- ozIMMU\_H- $k$  : Proposed method in Section 3.3 with  $k$  slices

### 4.1 Accuracy

Figure 5 shows the accuracies of ozIMMU- $k$ , ozIMMU\_RN- $k$ , ozIMMU\_EF- $k$ , and ozIMMU\_H- $k$  for  $k = 7, \dots, 11$ . We set  $A \in \mathbb{F}_{64}^{n \times n}$  as  $a_{ij} := (U_{ij} - 0.5) \cdot \exp(\phi \cdot N_{ij})$ , where  $U_{ij} \in (0, 1)$  are uniformly distributed random numbers and  $N_{ij}$  are drawn from the standard normal distribution, for  $1 \leq i, j \leq n$ . The constant  $\phi$  specifies the tendency of the difficulty in terms of matrix multiplication accuracy. A matrix  $B \in \mathbb{F}_{64}^{n \times n}$  is set similarly to  $A$ .

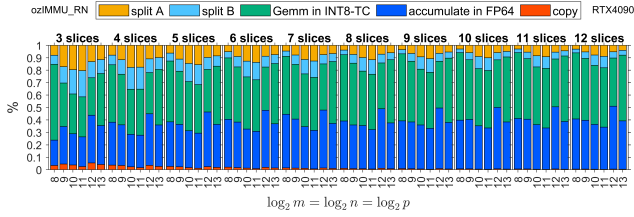


**Figure 5.** Comparison of accuracy between ozIMMU and proposed methods. Matrix  $A$  has entries  $a_{ij} := (U_{ij} - 0.5) \cdot \exp(\phi \cdot N_{ij})$ , where  $U_{ij} \in (0, 1)$  are uniformly distributed random numbers and  $N_{ij}$  are drawn from the standard normal distribution, for  $1 \leq i, j \leq m$  and  $m = n = p$ . Matrix  $B$  has a similar composition.

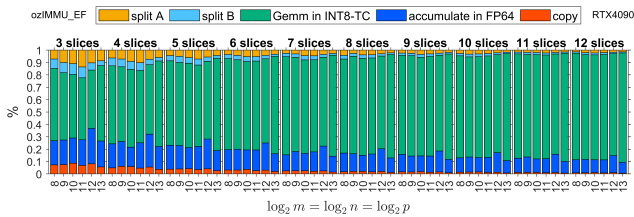
The accuracy deteriorated for  $n > 2^{17}$  because the maximum number of significant bits  $\beta = \min(7, \lfloor (31 - \log_2 n)/2 \rfloor)$  of the elements of the split matrices is less than 7. The accuracies of  $\text{ozIMMU}$  and  $\text{ozIMMU\_EF}$  are comparable, as are those of,  $\text{ozIMMU\_RN}$  and  $\text{ozIMMU\_H}$ .  $\text{ozIMMU\_RN}$  and  $\text{ozIMMU\_H}$ , which use a splitting method via rounding to nearest floating-point arithmetic, are more accurate than  $\text{ozIMMU}$  and  $\text{ozIMMU\_EF}$ , which use a splitting method via bit masking. Occasionally, to obtain results comparable with that of FP64,  $\text{ozIMMU\_RN-}k$ ,  $\text{ozIMMU\_H-}k$ ,  $\text{ozIMMU-}(k+1)$ , and  $\text{ozIMMU\_EF-}(k+1)$  are required; i.e.,  $\text{ozIMMU\_RN}$  and  $\text{ozIMMU\_H}$  require fewer splits than  $\text{ozIMMU}$  and  $\text{ozIMMU\_EF}$ . In particular, for  $\phi = 2$ ,  $\text{ozIMMU\_RN-}9$  and  $\text{ozIMMU\_H-}9$  produce comparable results to FP64; however, the accuracies of the results of  $\text{ozIMMU-}9$  and  $\text{ozIMMU\_EF-}9$  are worse than that of FP64. In such cases,  $k = 10$  is required for  $\text{ozIMMU}$  and  $\text{ozIMMU\_EF}$  to attain more accurate results than FP64.

## 4.2 Time breakdown

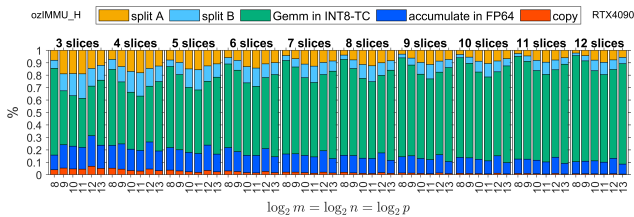
Figures 6–8 and Figures 9–11 show time breakdowns of the proposed methods on RTX 4090 and GH200, respectively. Note that “split A”, “split B”, “Gemm in INT8-TC”, “accumulation in FP64”, and “copy” correspond to the steps (i), (ii), (iii), (iv), and (v) in the Ozaki scheme for emulating the GEMM routine as shown in Section 1.



**Figure 6.** Time breakdown of  $\text{ozIMMU\_RN}$  on NVIDIA GeForce RTX 4090

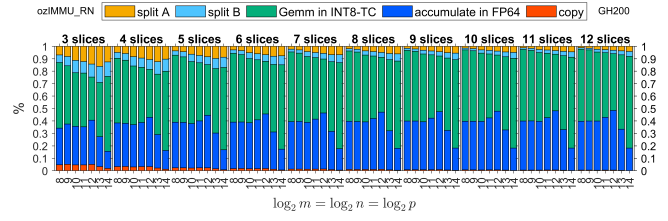


**Figure 7.** Time breakdown of  $\text{ozIMMU\_EF}$  on NVIDIA GeForce RTX 4090

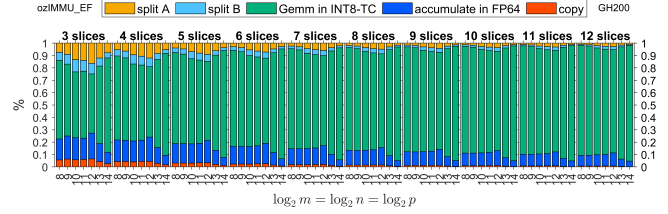


**Figure 8.** Time breakdown of  $\text{ozIMMU\_H}$  on NVIDIA GeForce RTX 4090

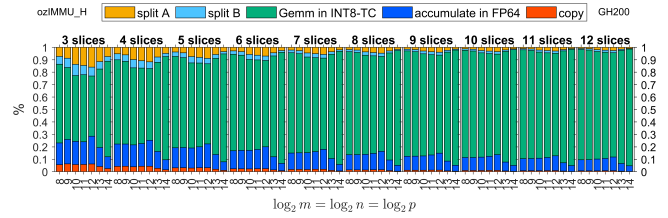
The execution time of FP64 accumulation in  $\text{ozIMMU\_RN}$  is not negligible, because the number of additions in FP64 is



**Figure 9.** Time breakdown of  $\text{ozIMMU\_RN}$  on NVIDIA GH200 Grace Hopper Superchip



**Figure 10.** Time breakdown of  $\text{ozIMMU\_EF}$  on NVIDIA GH200 Grace Hopper Superchip

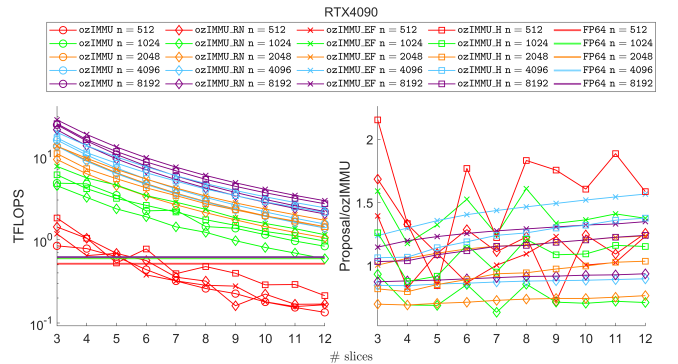


**Figure 11.** Time breakdown of  $\text{ozIMMU\_H}$  on NVIDIA GH200 Grace Hopper Superchip

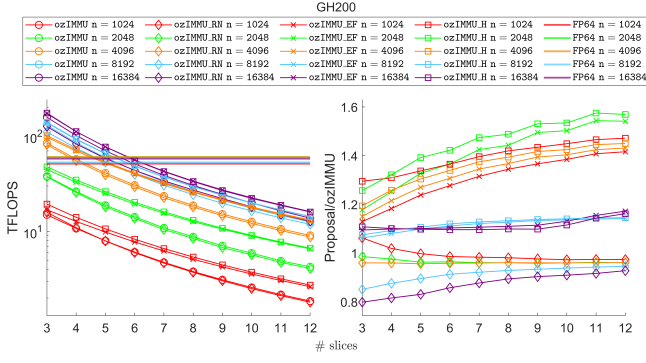
not reduced. The execution times of FP64 accumulation in  $\text{ozIMMU\_EF}$  and  $\text{ozIMMU\_H}$  are less than those in  $\text{ozIMMU}$  and  $\text{ozIMMU\_RN}$ . From Figures 6–8, the computation time of splitting via rounding to nearest floating-point arithmetic is not so fast because FP64 is much slower than the lower-precision arithmetic on RTX 4090. From Figures 10 and 11, the ratios of the computation time of splitting in  $\text{ozIMMU\_EF}$  and  $\text{ozIMMU\_H}$  are comparable on GH200.

## 4.3 Performance

Figures 12 and 13 show throughput in TFLOPS and ratio to  $\text{ozIMMU}$ . A smaller number of slices corresponds to better performance because the total computation cost is less.



**Figure 12.** Throughput comparison on NVIDIA GeForce RTX 4090



**Figure 13.** Throughput comparison on NVIDIA GH200 Grace Hopper Superchip

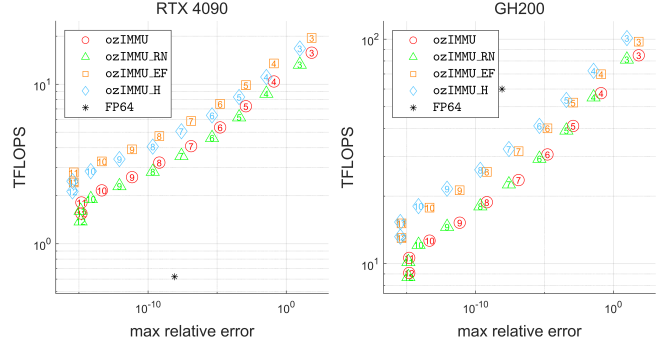
On RTX 4090, all methods are faster than or comparable to matrix multiplication in FP64 for  $n \geq 1024$  because FP64 is much slower than the lower-precision arithmetic.  $\text{ozIMMU\_EF}$  and  $\text{ozIMMU\_H}$  are faster than  $\text{ozIMMU}$  almost everywhere. In particular,  $\text{ozIMMU\_EF-12}$  and  $\text{ozIMMU\_H-12}$  are respectively 1.6 and 1.4 times faster than  $\text{ozIMMU}$  for  $n = 4096$ .  $\text{ozIMMU\_RN}$  is slower than  $\text{ozIMMU}$ . For  $k = 3$ ,  $\text{ozIMMU}$ ,  $\text{ozIMMU\_RN}$ ,  $\text{ozIMMU\_EF}$ , and  $\text{ozIMMU\_H}$  are 38.8, 33.8, 44.4, and 40.0 times faster than matrix multiplication in FP64 for  $n = 16384$ , respectively. For  $k = 12$ ,  $\text{ozIMMU}$ ,  $\text{ozIMMU\_RN}$ ,  $\text{ozIMMU\_EF}$ , and  $\text{ozIMMU\_H}$  are 1.4, 1.0, 1.9, and 1.6 times faster than matrix multiplication in FP64 for  $n = 1024$ , respectively. In addition,  $\text{ozIMMU}$ ,  $\text{ozIMMU\_RN}$ ,  $\text{ozIMMU\_EF}$ , and  $\text{ozIMMU\_H}$  are 9.4, 8.5, 12.0, and 10.9 times faster than FP64 for  $k = 7$  and  $n = 16384$ , and 7.4, 6.8, 9.5, and 8.6 times faster than FP64 for  $k = 8$  and  $n = 16384$ .

On GH200, the methods with small  $k$  are faster than matrix multiplication in FP64 for  $n \geq 4096$ ; however, the methods are much slower than FP64 otherwise.  $\text{ozIMMU\_EF}$  and  $\text{ozIMMU\_H}$  are faster than  $\text{ozIMMU}$ . In particular,  $\text{ozIMMU\_EF-11}$  and  $\text{ozIMMU\_H-11}$  are respectively 1.5 and 1.6 times faster than  $\text{ozIMMU}$  for  $n = 2048$ .  $\text{ozIMMU\_RN}$  is not as slow on GH200 as it is on RTX 4090. For  $k = 3$ ,  $\text{ozIMMU}$ ,  $\text{ozIMMU\_RN}$ ,  $\text{ozIMMU\_EF}$ , and  $\text{ozIMMU\_H}$  are 2.7, 2.2, 3.0, and 3.0 times faster than matrix multiplication in FP64 for  $n = 16384$ , respectively.  $\text{ozIMMU-5}$ ,  $\text{ozIMMU\_RN-5}$ ,  $\text{ozIMMU\_EF-6}$ , and  $\text{ozIMMU\_H-6}$  have comparable computation times to matrix multiplication in FP64 for  $n = 16384$ ; however,  $\text{ozIMMU-k}$ ,  $\text{ozIMMU\_RN-k}$ ,  $\text{ozIMMU\_EF-(k+1)}$ , and  $\text{ozIMMU\_H-(k+1)}$  are slower than FP64 for  $n < 16384$  or  $k > 5$ .

#### 4.4 Performance vs. Accuracy

Figure 14 shows throughput in a TFLOPS vs. Accuracy plot, illustrating relationships between performance and accuracy for the Ozaki scheme and for matrix multiplication in FP64 for  $n = 4096$  and  $\phi = 0$ . The numbers inside the symbols are the numbers of splices for the Ozaki scheme.

On RTX 4090, all methods are better than matrix multiplication in FP64 with respect to both performance and accuracy for  $k \geq 8$ .  $\text{ozIMMU\_H}$  and  $\text{ozIMMU\_EF}$  produce results with comparable performances that are more accurate than those of  $\text{ozIMMU}$ . In particular,



**(a)** NVIDIA GeForce RTX 4090 **(b)** NVIDIA GH200 Grace Hopper Superchip

**Figure 14.** Relationship between performance and accuracy for  $n = 4096$  and  $\phi = 0$

$\text{ozIMMU\_H-8}$  and  $\text{ozIMMU\_EF-9}$  produce results with comparable performances that are more accurate than those of  $\text{ozIMMU-7}$ . On GH200 as well,  $\text{ozIMMU\_H}$  and  $\text{ozIMMU\_EF}$  produce results with comparable performances that are more accurate than those of  $\text{ozIMMU}$ . In particular,  $\text{ozIMMU\_H-7}$  and  $\text{ozIMMU\_EF-7}$  produce results with comparable performances that are more accurate than those of  $\text{ozIMMU-6}$ . These results also indicate that  $\text{ozIMMU\_H}$  and  $\text{ozIMMU\_EF}$  can be computed with less memory consumption than  $\text{ozIMMU}$ .

### 5 Rounding error analysis

In this section, we give a rounding analysis for a fixed number of slices. We have described Algorithms 3 and 4 ( $\text{ozIMMU}$ ), and proposed Algorithms 5 and 6, Algorithms 3 and 6, and Algorithms 8 and 6. Even if we used the rounding to nearest strategy as in Algorithms 5 and 8, their error bounds would be the same as that of the bitmask strategy used in Algorithm 5. Therefore, we focus on Algorithms 3 and 4 ( $\text{ozIMMU}$ ) and Algorithms 3 and 6. Below, absolute value notation applied to a matrix means the matrix from applying absolute value element-wise. We assume that neither overflow nor underflow occurs. For  $A \in \mathbb{F}^{m \times n}$  and  $B \in \mathbb{F}^{n \times p}$ , the following deterministic error bound is given by Jeannerod and Rump (2013):

$$|AB - \text{fl}(AB)| \leq nu|A||B|.$$

The following alternative probabilistic error bound comes from Higham and Mary (2019) (which has the details on the assumptions and probabilities):

$$|AB - \text{fl}(AB)| \lesssim \sqrt{nu}|A||B|.$$

In this section, our aim is to derive an error bound on the computed result. Let  $T \in \mathbb{F}^{m \times p}$  be a computed result using  $k$  slices, such as

$$A \approx A_1 + A_2 + \dots + A_k, \quad B \approx B_1 + B_2 + \dots + B_k.$$

Note that for  $i = 1, 2, \dots, k$ ,

$$A_i := \text{diag}(\mu'^{(i)})A'_i, \quad B_i := B'_i \text{diag}(\nu'^{(i)})$$

for Algorithms 1, 2, and 5, and

$$A_i := \text{diag}(\mu'')2^{-i\beta+1}A''_i, \quad B_i := 2^{-i\beta+1}B''_i \text{diag}(\nu'')$$



for Algorithms 3, 4, 6, 7, and 8. Let  $T_k$  be an approximation of  $AB$  with  $k$  slices. Our aim is to obtain the upper bound on  $|AB - T|$  as follows:

$$\begin{aligned} & |AB - T_k| \\ & \leq \left| AB - \sum_{i=1}^k \sum_{j=1}^{k-i+1} A_i B_j \right| + \left| \sum_{i=1}^k \sum_{j=1}^{k-i+1} A_i B_j - T_k \right|. \end{aligned} \quad (13)$$

Note that we can obtain the exact product  $A_i B_j$  in the above by using GEMM in the INT8 Tensor Core and scaling by powers of two. The first term of the bound in (13) indicates the truncation error, and the second term shows an error arising in the accumulation process. As matrix scaling does not affect rounding errors, a discussion on it is omitted.

In the following subsections, we use an upper bound provided in Jeannerod and Rump (2013) for a sum:

$$\left| \sum_{i=1}^n p - \text{fl} \left( \sum_{i=1}^n p \right) \right| \leq (n-1)u \sum_{i=1}^n |p_i|, \quad p \in \mathbb{F}^n. \quad (14)$$

Let  $\text{ufp}(c)$  for  $c \in \mathbb{R}$  be defined as

$$\text{ufp}(c) := \begin{cases} 0 & \text{if } c = 0, \\ 2^{\lfloor \log_2 |c| \rfloor} & \text{otherwise.} \end{cases}$$

The following condition is used to check whether a real number is a floating-point number. For  $c \in \mathbb{R}$ , if  $|c|$  is smaller than the maximum floating-point number and  $c$  is an integral multiple of the minimum positive floating-point number, then it holds that

$$c \in 2u \cdot \text{ufp}(c)\mathbb{Z} \Rightarrow c \in \mathbb{F}. \quad (15)$$

### 5.1 Error bound for Algorithms 3 and 4 (ozIMMU)

Let matrices  $V_k \in \mathbb{F}^{m \times n}$  and  $W_k \in \mathbb{F}^{n \times p}$  be defined as in (1) where these show the truncation errors of  $A$  and  $B$  for the case of  $k$  slices, respectively. Notation  $e$  indicates the vector  $e = (1, 1, \dots, 1)^T \in \mathbb{F}^n$ . Define two vectors  $g \in \mathbb{F}^m$  and  $h \in \mathbb{F}^p$  related to row-wise maximums in  $A$  and column-wise maximums in  $B$  in the sense of the unit in the first place as

$$g_i := \text{ufp} \left( \max_j |a_{ij}| \right), \quad f_j := \text{ufp} \left( \max_i |b_{ij}| \right).$$

Then, from Figure 15 and (1), we have

$$|V_i| \leq 2^{-\beta i+1} g e^T, \quad |W_i| \leq 2^{-\beta i+1} e f^T, \quad (16)$$

and matrices  $|A_i|$  and  $|B|$  are bounded as

$$|A_i| \leq 2^{-\beta(i-1)+1} g e^T, \quad |B| \leq 2 e f^T. \quad (17)$$

Because

$$AB = (A_1 + \dots + A_k + V_k)(B_1 + \dots + B_k + W_k),$$

we have

$$AB - \sum_{i=1}^k \sum_{j=1}^{k-i+1} A_i B_j = \sum_{i=1}^k A_i W_{k-i+1} + V_k B.$$

From the last equation, (16), and (17), the truncation error is bounded as

$$\begin{aligned} & \left| AB - \sum_{i=1}^k \sum_{j=1}^{k-i+1} A_i B_j \right| \\ & \leq \sum_{i=1}^k |A_i| |W_{k-i+1}| + |V_k| |B| \\ & \leq \sum_{i=1}^k 2^{-\beta(i-1)+1} g e^T \cdot 2^{-\beta(k-i+1)+1} e f^T \\ & \quad + 2 \cdot 2^{-\beta k+1} g e^T \cdot e f^T \\ & = 4n \sum_{i=1}^k 2^{-\beta k} g f^T + 4n \cdot 2^{-\beta k} g f^T \\ & = 4(k+1)n 2^{-\beta k} g f^T. \end{aligned} \quad (18)$$

If  $|g_i - |a_{ij}||$  for  $1 \leq j \leq n$  and  $|f_j - |b_{ij}||$  for  $1 \leq i \leq n$  are very small, we can assume that

$$4n g f^T \approx |A| |B|. \quad (19)$$

In this case, we have

$$\left| AB - \sum_{i=1}^k \sum_{j=1}^{k-i+1} A_i B_j \right| \lesssim (k+1) 2^{-\beta k} |A| |B|. \quad (20)$$

We next focus on an error in the accumulation, the second term in (13). If  $A_i$  and  $B_i$  are obtained by Algorithm 3,

$$|A| = \sum_{i=1}^k |A_i|, \quad |B| = \sum_{i=1}^k |B_i| \quad (21)$$

are satisfied. Because we compute the sum of  $\frac{1}{2}k(k+1)$  floating-point matrices, the following bound is immediately obtained from (14) and (21):

$$\begin{aligned} & \left| \sum_{i=1}^k \sum_{j=1}^{k-i+1} A_i B_j - T_k \right| \\ & \leq \left( \frac{1}{2}k(k+1) - 1 \right) u \sum_{i=1}^k \sum_{j=1}^{k-i+1} |A_i B_j| \\ & \leq \left( \frac{1}{2}k(k+1) - 1 \right) u |A| |B|. \end{aligned} \quad (22)$$

We will show that  $k' (\leq k)$  exists such that

$$\sum_{s=1}^{k'} \sum_{t=1}^{k'-s+1} A_s B_t \in \mathbb{F}^{m \times n}. \quad (23)$$

Then, the bound (22) is improved. From Figure 15, we have

$$(A_s)_{ij} \in 2^{-s\beta+1} g_i \mathbb{Z}, \quad (B_t)_{ij} \in 2^{-t\beta+1} f_j \mathbb{Z}, \quad (24)$$

and these derive

$$(A_s)_{i\ell} (B_t)_{\ell j} \in 2^{-\beta(s+t)+2} g_i f_j \mathbb{Z}. \quad (25)$$

Figure 15 shows an image of (24). It follows from this that

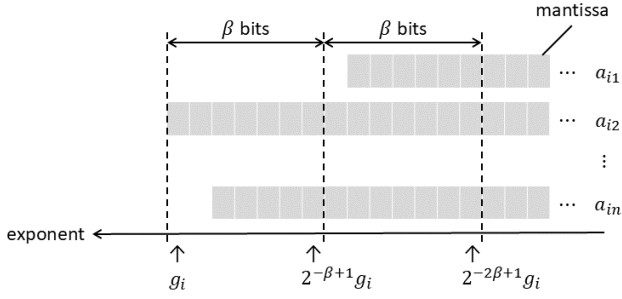


Figure 15. An image of (24)

$$\left( \sum_{s=1}^{k'} \sum_{t=1}^{k'-s+1} A_s B_t \right)_{ij} \in 2^{-\beta(k'+1)+2} g_i f_j \mathbb{Z}. \quad (26)$$

From Figure 15, the elements of the matrices are bounded as

$$|A_s|_{ij} \leq 2^{-\beta(s-1)+1} g_i, \quad |B_t|_{ij} \leq 2^{-\beta(t-1)+1} f_j,$$

so that we have

$$|A_s B_t|_{ij} \leq n 2^{-\beta(s+t-2)+2} g_i f_j. \quad (27)$$

Therefore, from (27), we have the following bound on a dot product:

$$\begin{aligned} & \left| \sum_{s=1}^{k'} \sum_{t=1}^{k'-s+1} A_s B_t \right|_{ij} \\ &= \left| \sum_{s=2}^{k'+1} \sum_{i_2+j_2=s} A_{i_2} B_{j_2} \right|_{ij} \\ &\leq \left( \sum_{s=2}^{k'+1} \sum_{i_2+j_2=s} |A_{i_2}| |B_{j_2}| \right)_{ij} \\ &\leq \sum_{s=2}^{k'+1} \sum_{i_2+j_2=s} n 2^{-(i_2+j_2-2)\beta+2} g_i f_j \\ &= \sum_{s=2}^{k'+1} \sum_{i_2+j_2=s} n 2^{-(s-2)\beta+2} g_i f_j \\ &= \sum_{s=2}^{k'+1} (s-1) n 2^{-(s-2)\beta+2} g_i f_j. \end{aligned}$$

If  $\beta \geq 3$ ,

$$\sum_{s=3}^{k'+1} (s-1) n 2^{-(s-2)\beta+2} g_i f_j \leq n 2^2 g_i f_j$$

Finally, for  $\beta \geq 3$ , we have

$$\left| \sum_{s=1}^{k'} \sum_{t=1}^{k'-s+1} A_s B_t \right|_{ij} \leq 8n \cdot g_i f_j \quad (28)$$

From (26), (28), and (15), if

$$2u \cdot \text{ufp}(8n) \leq 2^{-\beta(k'+1)+2} \quad (29)$$

is satisfied, then (23) holds.

Let  $k'_{\max}$  be the maximum  $k'$  satisfying (29). No rounding error occurs for  $i = 2, 3, \dots, k'_{\max}$ . Therefore, we have

$$\begin{aligned} & \left| \sum_{k=1}^k \sum_{j=1}^{k-i+1} A_i B_j - T_k \right| \\ &\leq \left( \frac{1}{2} k(k+1) - \frac{1}{2} k'_{\max}(k'_{\max}+1) - 1 \right) u |A| |B|. \end{aligned} \quad (30)$$

Summarizing, from (18) and (30), we have

$$\begin{aligned} & |AB - T_k| \\ &\leq 4(k+1) n 2^{-\beta k} g f^T \\ &\quad + \left( \frac{1}{2} k(k+1) - \frac{1}{2} k'_{\max}(k'_{\max}+1) - 1 \right) u |A| |B|. \end{aligned}$$

If (19) is satisfied, then we have

$$\begin{aligned} & |AB - T_k| \\ &\lesssim (k+1) 2^{-\beta k} |A| |B| \\ &\quad + \left( \frac{1}{2} k(k+1) - \frac{1}{2} k'_{\max}(k'_{\max}+1) - 1 \right) u |A| |B|. \end{aligned}$$

## 5.2 Error bound for Algorithms 3 and 6

We focus on

$$\sum_{i_2+j_2=s} A''_{i_2} B''_{j_2}, \quad A''_{i_2} \in \mathbb{I}_{\beta+1}^{m \times n}, \quad B''_{j_2} \in \mathbb{I}_{\beta+1}^{n \times p}. \quad (31)$$

We consider  $r$  such that the following holds

$$\sum_{\substack{i_2+j_2=s \\ s \leq \min(r, k+1)}} A''_{i_2} B''_{j_2} \in \mathbb{I}_{32}^{m \times p}. \quad (32)$$

Let  $E_{m,n}$  be the matrix  $\mathbb{I}_2^{m \times n}$  whose elements are all ones. From the definitions in (31),

$$|A''_{i_2}| \leq (2^\beta - 1) E_{m,n}, \quad |B''_{j_2}| \leq (2^\beta - 1) E_{n,p},$$

and we have

$$\begin{aligned} & \left| \sum_{i_2+j_2=s} A''_{i_2} B''_{j_2} \right| \leq \sum_{i_2+j_2=s} |A''_{i_2}| |B''_{j_2}| \\ &\leq (s-1) n (2^\beta - 1)^2 E_{m,p}. \end{aligned}$$

Since the largest number in  $\mathbb{I}_{32}$  is  $2^{31} - 1$ , if  $s \leq r$  for  $r := \max(1, 2^{31-2\beta-\lceil \log_2 n \rceil})$  in (12), we have

$$\begin{aligned} & (s-1) n (2^\beta - 1)^2 \leq (r-1) n (2^\beta - 1)^2 \\ &\leq (r-1) \cdot 2^{\lceil \log_2 n \rceil} \cdot 2^{2\beta} \\ &= r \cdot 2^{\lceil \log_2 n \rceil} \cdot 2^{2\beta} - 2^{\lceil \log_2 n \rceil} \cdot 2^{2\beta} \\ &\leq 2^{31} - 1, \end{aligned}$$

which implies that (32) holds. Note that

$$(s-1) n (2^\beta - 1)^2 < 2^{31} - 1$$

for  $\beta \geq 1$ . This is why constant  $r$  was set as in (12).

If we define the number of the terms,  $w$ , for the accumulation as

$$w := \left\lceil \frac{k}{r} \right\rceil \left( k - \frac{r}{2} \left\lfloor \frac{k-1}{r} \right\rfloor \right),$$

we have

$$|AB - T| \leq 4(k+1)n2^{-\beta k}gf^T + (w-1)u|A||B|.$$

If (19) is satisfied, then we have

$$|AB - T| \lesssim (k+1)2^{-\beta k}|A||B| + (w-1)u|A||B|.$$

## 6 Conclusion

We proposed three implementation methods for accelerating the Ozaki scheme using the INT8 Tensor Core. In the original implementation, `ozIMMU`, provided by Ootomo et al., the ratio of the accumulation in FP64 is not negligible, as it accounts for approximately 40–50 % of the total computation time. The proposed methods `ozIMMU_EF` and `ozIMMU_H` reduce the computation time ratio of the accumulation to approximately 10–20 % and achieve a 1.2- to 1.6-fold speedup. With future architectures expected to achieve high speed in lower-precision matrix multiplication, the computation time ratio of the accumulation in FP64 is expected to increase. Thus, these proposed methods contribute to leveraging the performance of future architectures more effectively than `ozIMMU`. The proposed methods `ozIMMU_RN` and `ozIMMU_H` offer alternative splitting methods using floating-point arithmetic in round-to-nearest-even mode and produce more accurate results than `ozIMMU` for the same number of slices.

## Acknowledgements

We thank Dr. Hiroyuki Ootomo from NVIDIA for his helpful comments on the implementation of `ozIMMU`.

## Declaration of conflicting interests

The authors declare no competing interests.

## Funding

This study was partially supported by JSPS Grant-in-Aid for JSPS Fellows No. 22KJ2741, JSPS Grant-in-Aid for Research Activity Start-up No. 24K23874, and JSPS KAKENHI Grant No. 23K28100.

## Supplemental material

Not applicable.

## References

- Higham NJ and Mary T (2019) A new approach to probabilistic rounding error analysis. *SIAM journal on scientific computing* 41(5): A2815–A2835.
- IEEE Computer Society (2019) IEEE standard for floating-point arithmetic. *IEEE Std 754-2019 (Revision of IEEE 754-2008)* DOI:10.1109/IEEESTD.2019.8766229.
- Jeannerod CP and Rump SM (2013) Improved error bounds for inner products in floating-point arithmetic. *SIAM Journal on Matrix Analysis and Applications (SIMAX)* 34(2): 338–344.

Minamihata A, Ozaki K, Ogita T and Oishi S (2016) Improved extraction scheme for accurate floating-point summation. In: *The 35th JSST Annual Conference International Conference on Simulation Technology*.

Mukunoki D, Ozaki K, Ogita T and Imamura T (2020) Dgemm using tensor cores, and its accurate and reproducible versions. In: Sadayappan P, Chamberlain BL, Juckeland G and Ltaief H (eds.) *High Performance Computing*. Cham: Springer International Publishing, pp. 230–248.

NVIDIA Corporation (2024) NVIDIA Tensor Cores. URL <https://www.nvidia.com/en-us/data-center/tensor-cores/>.

Ootomo H (2024) ozIMMU - DGEMM on Int8 Tensor Core. URL <https://github.com/enpls0/ozIMMU>.

Ootomo H, Ozaki K and Yokota R (2024) Dgemm on integer matrix multiplication unit. *The International Journal of High Performance Computing Applications* in Press. DOI:10.1177/10943420241239588.

Ozaki K, Ogita T, Oishi S and Rump SM (2012) Error-free transformations of matrix multiplication by using fast routines of matrix multiplication and its applications. *Numerical Algorithms* 59(1): 95–118.

Ozaki K, Ogita T, Oishi S and Rump SM (2013) Generalization of error-free transformation for matrix multiplication and its application. *Nonlinear Theory and Its Applications, IEICE* 4(1): 2–11.

Uchino Y (2024) Accelerator for ozIMMU. R-CCS github repository. URL [https://github.com/RIKEN-RCCS/accelerator\\_for\\_ozIMMU](https://github.com/RIKEN-RCCS/accelerator_for_ozIMMU).

## Author Biographies

Yuki Uchino is a postdoctoral researcher at RIKEN R-CCS. He received his Ph.D. in engineering from Shibaura Institute of Technology in 2024. His research interests include reliable computing, numerical linear algebra, and highly accurate algorithms. He won the student poster presentation award at the 38th JSST Annual International Conference on Simulation Technology (JSST 2019) and the student presentation awards at JSST 2021 and International Workshop on Reliable Computing and Computer-Assisted Proofs (ReCAP 2022).

Katsuhisa Ozaki is a full professor in the Department of Mathematical Sciences at Shibaura Institute of Technology. He received his Ph.D. in engineering from Waseda University in 2007. He was an Assistant Professor (2007–2008) and a Visiting Lecturer (2008–2009) at Waseda University. At Shibaura Institute of Technology, he has served as an Assistant Professor (2010–2013) and an Associate Professor (2013–2019), and has currently been a Professor since 2019. His research interests include reliable computing, particularly addressing rounding error problems in finite-precision arithmetic. He mainly focuses on numerical linear algebra and develops fast and accurate algorithms.

Toshiyuki Imamura is a team leader of Large-scale Parallel Numerical Computing Technology Team at RIKEN R-CCS, and is responsible for developing numerical libraries on Fugaku. He received his Diploma and Doctorate in Applied Systems and Sciences from Kyoto University in 1993 and 2000. He was a Researcher at CCSE, JAERI (1996–2003), a visiting scientist at HLRS (2002), and an associate professor at the University of Electro-Communications (2003–2012). His research interests include HPC, auto-tuning technology, and parallel eigenvalue

computation. His research group won the HPL-MpX ranking (2020-2021) and was nominated as the Gordon Bell Prize finalist in SC05, SC06, and SC20.