

Diagonal-Tiled Mixed-Precision Attention for Efficient Low-Bit MXFP Inference

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Transformer-based large language models (LLMs) have*
 002 *demonstrated remarkable performance across a wide range*
 003 *of real-world tasks, but their inference cost remains pro-*
 004 *hibitively high due to the quadratic complexity of attention*
 005 *and the memory bandwidth limitations of high-precision*
 006 *operations. In this work, we present a low-bit mixed-*
 007 *precision attention kernel using the microscaling floating-*
 008 *point (MXFP) data format, utilizing the computing capabil-*
 009 *ity on next-generation GPU architectures. Our Diagonal-*
 010 *Tiled Mixed-Precision Attention (DMA) incorporates two*
 011 *kinds of low-bit computation at the tiling-level, and is a*
 012 *delicate fused kernel implemented using Triton, exploit-*
 013 *ing hardware-level parallelism and memory efficiency to*
 014 *enable fast and efficient inference without compromising*
 015 *model performance. Extensive empirical evaluations on*
 016 *NVIDIA B200 GPUs show that our kernel maintains gen-*
 017 *eration quality with negligible degradation, and meanwhile*
 018 *achieves significant speedup by kernel fusion. We release*
 019 *our code at [https://anonymous.4open.science/r/MP-Sparse-](https://anonymous.4open.science/r/MP-Sparse-Attn)*
 020 *Attn.*

021 1. Introduction

022 The rapid development of large language models (LLMs)
 023 has created a growing demand for faster throughput. In-
 024 side the transformer architecture, self-attention becomes the
 025 main inference bottleneck because its cost scales quadrati-
 026 cally with sequence length. To reduce this cost, prior work
 027 has explored several directions, including quantization-
 028 based compression [5], structured sparsity [6], and kernel-
 029 based approximations such as linear attention [10]. At
 030 the same time, modern GPU architectures are adding
 031 stronger support for low-precision computation. In par-
 032 ticular, NVIDIA Blackwell introduces native support for
 033 microscaling number formats, including MXFP8, MXFP4,
 034 and NVFP4, achieves lower bitwidth with better quanti-
 035 zation quality [1, 8]. These advances make low-precision
 036 LLM inference practical in real deployment settings.

037 Despite this promising hardware support, effectively

Table 1. Comparison of representative MXFP data formats, including their block sizes, element formats, and shared scale formats.

Name	Block Size	Element		Scale	
		Format	Bits	Format	Bits
MXFP8	32	FP8 (E4M3 / E5M2)	8	E8M0	8
MXFP4	32	FP4 (E2M1)	4	E8M0	8
NVFP4	16	FP4 (E2M1)	4	E4M3	8

leveraging MX formats to accelerate attention in LLMs
 presents significant challenges. **Challenge 1:** low-bit quan-
 tization can lead to severe accuracy degradation. Directly
 applying 4-bit MXFP formats to attention computation
 causes substantial quantization error (see Tab. 2). While
 recent work [15] introduces smoothing techniques to miti-
 gate quantization errors, it requires additional floating-point
 GEMV operations that reduce kernel throughput. **Challenge**
2: unfused operations in quantization process undermine
 the efficiency of low-bit computation. Our experiment re-
 veals that without kernel fusion, the quantization and format
 conversion brings redundant memory accesses and kernel
 launch cost that cannot be underestimated (refer to Tab. 7).

In this paper, we propose **Diagonal-Tiled Mixed-**
Precision Attention (DMA), the first attention workflow
 designed to operate hybrid MX formats. DMA addresses
 the above challenges via two key techniques: First, we
 adopt a tiling-level mixed-precision design that partitions
 the attention matrix into low- and high-precision regions
 to retain the most salient information for critical tokens in
 high-precision, while leveraging faster low-precision MX
 formats elsewhere to ensure the speedup. Second, we build
 a full-stack attention computation workflow into a fused
 memory-efficient kernel to ensure the end-to-end efficiency.
 It covers quantization, microscaling transformation, low-bit
 encoding and packing, and attention computation. This fu-
 sion allows fine-grained parallel execution on GPU thread
 blocks without storing intermediate results, significantly re-
 ducing memory pressure and synchronization overhead.

We evaluate DMA on Longbench dataset. Experimental
 results show that DMA achieves lossless generation quality
 compared to full-precision attention baseline. Furthermore,

070 our ablation studies demonstrate the effectiveness of mixed-
071 precision tiling, quantization granularity, and diagonal win-
072 dow design in balancing performance and precision.

073 Our main contributions are summarized as follows:

- 074 • We propose Diagonal-Tiled Mixed-Precision Attention
075 (DMA), a new attention workflow that operates with hy-
076 brid MX formats for efficient LLM inference.
- 077 • We develop a fully fused GPU kernel that integrates quan-
078 tization, microscaling transformation, and attention com-
079 putation into one end-to-end workflow, reducing memory
080 access and kernel launch overhead.
- 081 • We conduct extensive experiments and detailed ablation
082 studies to show the lossless generation quality of our ker-
083 nel, and offer practical insights into the trade-off between
084 efficiency and accuracy.

085 2. Related Works

086 2.1. Efficient Attention

087 Under the context of large language models (LLMs), the
088 computational overhead of attention scales quadratically
089 with sequence length, making efficient attention a criti-
090 cal research direction. FlashAttention [4] introduced a
091 highly efficient attention kernel that uses tiling and an on-
092 line softmax approach to eliminate the need to buffer full
093 attention matrices. Quantization-based approaches such
094 as INT-FlashAttention [3] and SageAttention [13] series
095 compress attention operands into low-bit formats. They
096 are plug-and-play quantized attention kernels that achieves
097 2–3× speedup over FlashAttention. Sparse attention meth-
098 ods like SparseAttention [14] limit the number of active to-
099 ken pairs to reduce complexity while preserving important
100 information. As a lightweight sparse attention alternative,
101 TurboAttention [7] combines head-wise quantization and
102 sparsity-aware softmax approximation to deliver 1.2–1.8×
103 attention speed up and 4× KV cache compression.

104 2.2. MXFP Quantization

105 The recent introduction of Microscaling Floating-Point
106 (MXFP) formats is an advancement over traditional Block
107 Floating Point (BFP), designed to improve the numerical
108 efficiency and deployment flexibility of low-precision com-
109 putations in AI workloads [9]. As supported by NVIDIA’s
110 latest Blackwell GPU architecture, MXFP4 and MXFP8 al-
111 low for significantly higher theoretical throughput, up to
112 2× to 4× compared to FP16 while still maintaining com-
113 petitive numerical accuracy in mixed-precision settings [8].
114 However, despite these theoretical advantages, systematic
115 and practical software support for MXFP-based attention
116 kernels remains limited, especially for end-to-end pipelines
117 that operate below the standard IEEE precision formats.
118 This motivates our first contribution: a full-stack low-bit at-
119 tention pipeline for MXFP formats (e.g., MXFP8, MXFP4,

NVFP4), including efficient quantization and packing, low-
bit computation combined with OnlineSoftmax technique.

120 3. Preliminaries

121 3.1. MXFP Data Format

122 Microscaling (MX) formats decompose a tensor into low-
123 precision elements and a shared exponent scale per block
124 (typically 32 or 16 elements). This approach allows
125 dynamic range coverage much larger than conventional
126 floating-point formats while significantly reducing storage
127 and compute cost. Table 1 summarizes the MX formats
128 used in this work. MX formats adopt a shared scale in
129 E8M0 for every 32 elements per block and are named by
130 prefixing “MX” to the element data format, i.e., MXFP8
131 (with FP8 data of E5M2/E4M3) or MXFP4 (FP4 data of
132 E2M1). Since the exponent of FP32 uses 8-bit, the rep-
133 resentable range of these formats covers that of FP32, en-
134 suring compatibility with high-dynamic-range input distri-
135 butions. While NVFP4 has FP8 (E4M3) shared scale for
136 every 16 elements in each block. The finer-grained scal-
137 ing and quantization granularity significantly reduces quan-
138 tization error compared to MXFP4, enhancing downstream
139 accuracy. 140

141 3.2. Online Softmax

142 To enable faster attention computation by dividing them
143 into blocks, many efficient attention kernels uses Online-
144 Softmax [4, 13] to ensure the equivalence of the results. For
145 example, FlashAttention avoids materializing and storing
146 the full attention matrix by computing the attention output
147 in a tile-wise fashion. It fuses the attention score compu-
148 tation, softmax normalization, and value aggregation into a
149 single memory-efficient kernel, using online row-wise soft-
150 max. 151

152 Each attention head is processed in parallel across
153 blocks. L_Q, L_K are the sequence lengths of query and
154 key/value matrices, respectively. D is the head dimen-
155 sion. We divide \mathbf{Q} into tiles of size B_M . This results in
156 $\lceil L_Q/B_M \rceil$ thread blocks executing in parallel per batch per
157 head, each responsible for a tile of query $\mathbf{Q}_i \in \mathbb{R}^{B_M \times D}$.
158 Within each thread block, it iterate over the key/value for
159 $\lceil L_k/B_N \rceil$ tiles. Let $\mathbf{Q}_i \in \mathbb{R}^{B_M \times D}$ denote the i^{th} query
160 tile, and $\mathbf{K}_j, \mathbf{V}_j \in \mathbb{R}^{B_N \times D}$ denote the j^{th} key and value
161 tiles respectively. For each (i, j) tile pair, the attention score
162 matrix is computed as:

$$163 \mathbf{S}_{ij} = \tilde{\mathbf{Q}}_i \times \mathbf{K}_j^\top, \quad \tilde{\mathbf{Q}}_i = \frac{\mathbf{Q}_i}{\sqrt{D}}. \quad (1)$$

164 It then applies OnlineSoftmax across key/value tiles to elim-
165 inate the need to store the full matrices of intermediate re-
166 sults. It tracks the running maximum of \mathbf{S}_i in \mathbf{m}_i and accu-
167 mulates the intermediate attention outputs \mathbf{O}_{i-1} across all

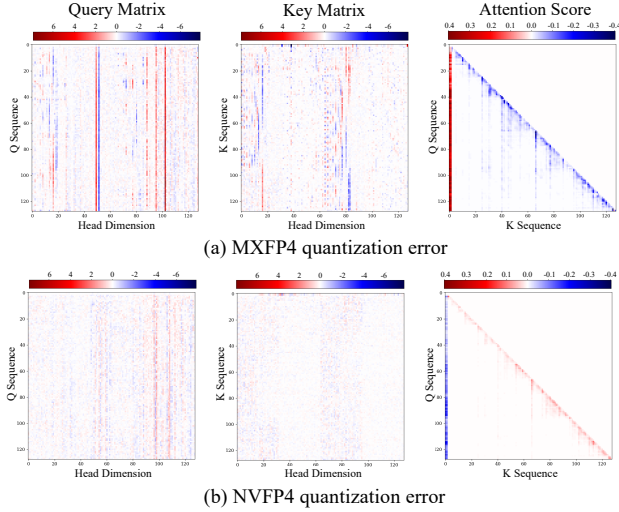


Figure 1. Visualization of quantization error of MXFP4 and NVFP4 format for query, key and attention score.

previous tiles, which is normalized by the ratio of the adjacent maximum values. The accumulated output \mathbf{O}_i and normalization factor l_i are updated incrementally, tile by tile. We finalize the results of the i^{th} query tile by $\text{diag}(l_i)^{-1}\mathbf{O}_i$.

4. Challenges in Low-bit MXFP Attention

Although full 4-bit quantization can reduce computation and memory cost, directly applying MXFP4 to the QK^T leads to clear degradation in attention quality. As shown in Fig. 1, the quantization error of MXFP4 is much pronounced than that of NVFP4. Table 2 confirms this trend quantitatively. Lower-bitwidth MX formats introduce larger quantization errors under all metrics. Direct MXFP4 causes a clear drop in attention-score fidelity, with cosine similarity decreasing from 0.988 to 0.714, PSNR from 71.70 to 60.82, and L1 error increasing from 0.246 to 0.924. NVFP4 is notably more stable, and adding token-wise quantization on top of it brings only marginal changes.

We observe that the error exhibits a clear channel-wise structure in query and key matrices, indicating that some feature dimensions are consistently more sensitive to low-bit quantization than others. However, this channel dimension is exactly the reduction axis in the QK^T multiplication. Therefore, it is unable to use channel-wise scaling in a fused low-bit attention kernel without introducing substantial implementation complexity.

5. Diagonal-Tiled Mixed-Precision Attention

In this section, we present the overall design of our method, which is built on the typical tiling-style of FlashAttention, but extends it with two components tailored for low-bit

Table 2. Quantization error of attention score using different data formats. “+” means combining with tokenwise quantization.

Format	Cos Sim (\uparrow)	PSNR (\uparrow)	L1 (\downarrow)	RMSE (\downarrow)
MXFP8	0.988	71.70	0.246	0.003
MXFP4	0.714	60.82	0.924	0.009
NVFP4	0.982	69.37	0.309	0.003
NVFP4+	0.983	69.63	0.312	0.003
Ours	0.988	71.70	0.248	0.003

MXFP inference: a diagonal-tiled mixed-precision computation scheme to preserve the most sensitive attention regions (Sec. 5.2), and a fused quantization kernel to avoid the overhead of separate low-bit pre-processing steps (Sec. 5.3).

5.1. Overall Workflow

We follow the tiled execution pattern of FlashAttention to partition attention into sub-tensors for memory-efficient and parallel computation. Built on this workflow, DMA introduces a quantization-aware attention kernel that supports tile-level mixed precision on GPU tensor cores with native MXFP support. Meanwhile, we fuse the low-bit pre-processing pipeline into a single Triton implementation. Instead of executing quantization, low-bit encoding, packing, and scale conversion as separate operators, DMA performs them within the kernel, which reduces intermediate memory traffic, kernel launch overhead, and synchronization cost. This fused design is important for maintaining the end-to-end efficiency of mixed-precision attention. The overall framework is illustrated in Fig. 2.

5.2. Diagonal-Tiled Attention Workflow

We propose Diagonal-Tiled Mixed-Precision strategy to mitigate the accuracy degradation. As detailed in Sec. 4, low-bit quantization introduces significant quantization error, and the most influential scores typically concentrate along the diagonal of the attention matrix. Prior solutions, such as SageAttention2 [12], compensates for quantization error by performing additional GEMV operations between full-precision K and the mean-pooled query vector, which brings extra full-precision computation overhead. To mitigate this, we propose a mixed-precision strategy in tiling-level that selectively retains higher-precision tiles near diagonal, while aggressively quantizing peripheral regions to ensure the overall throughput.

Based on the tiling-wise attention paradigm, we introduce a diagonal window with size T , determining the amount of tokens for higher-precision computation along the token dimension. For each query tile indexed by n_Q , only the last T tokens are computed using higher-precision representations of query and key, while all preceding tiles are aggressively quantized. Taking causal attention as an

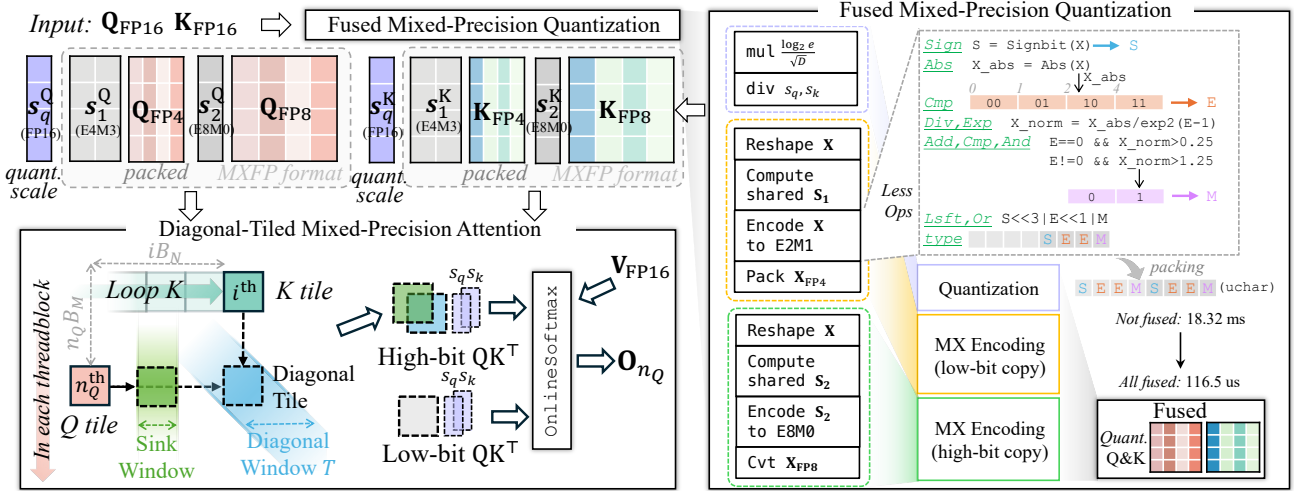


Figure 2. Overview workflow of our Diagonal-Tiled Mixed-Precision Attention. It first applies fused mixed-precision quantization to produce low-bit and high-bit Q and K representations, and then performs diagonal-tiled mixed-precision attention, using higher precision near the diagonal and lower precision elsewhere to balance accuracy and efficiency.

Algorithm 1 Diagonal-Tiled Mixed-Precision Attention

Input: $\tilde{Q}^l \in \mathbb{R}_{FP4}^{B_M \times D}$, $\tilde{Q}^h \in \mathbb{R}_{FP8}^{B_M \times D}$ are pre-processed with softmax scale. $\mathbf{K}^l \in \mathbb{R}_{FP4}^{B_N \times D}$. $\mathbf{K}^h \in \mathbb{R}_{FP8}^{B_N \times D}$. $V \in \mathbb{R}_{FP16}^{B_N \times D}$. B_M, B_N : tile sizes in sequence dimension for query and key, respectively. n_Q : thread block index. T : diagonal window size for high-precision.

1: Initialize: $\mathbf{m}_0 \leftarrow -\infty, \mathbf{l}_0 \leftarrow 1, \mathbf{O}_0 \leftarrow \mathbf{0}$

Phase 1: Low-precision tiles

2: **for** $i \leftarrow 0$ **to** $\lceil (n_Q B_M - T) / B_N \rceil$ **do**

3: $\mathbf{S}_i \leftarrow \tilde{Q}_{n_Q}^l \times \mathbf{K}_i^{lT}$

4: $\mathbf{O}_i, \mathbf{l}_i, \mathbf{m}_i \leftarrow \text{OnlineSoftmax}(\mathbf{S}_i, \mathbf{V}_i, \mathbf{O}_{i-1}, \mathbf{l}_{i-1}, \mathbf{m}_{i-1})$

5: **end for**

Phase 2: High-precision tiles in diagonal window

6: **for** $i \leftarrow \lceil (n_Q B_M - T) / B_N \rceil$ **to** $\lceil n_Q B_M / B_N \rceil$ **do**

▷ Skip upper triangular part for *causal attention*

7: $\mathbf{S}_i \leftarrow \tilde{Q}_{n_Q}^h \times \mathbf{K}_i^{hT}$

8: $\Delta_i \leftarrow (n_Q B_M + \mathcal{I}_{B_M}) < (i B_N + \mathcal{I}_{B_N})$

9: $\mathbf{S}_i \leftarrow \mathbf{S}_i - \text{inf} \cdot \Delta_i$

▷ Mask upper triangular of the global attention

10: $\mathbf{O}_i, \mathbf{l}_i, \mathbf{m}_i \leftarrow \text{OnlineSoftmax}(\mathbf{S}_i, \mathbf{V}_i, \mathbf{O}_{i-1}, \mathbf{l}_{i-1}, \mathbf{m}_{i-1})$

11: **end for**

12: $\mathbf{O}_{n_Q} \leftarrow \text{diag}(\mathbf{l}_i)^{-1} \mathbf{O}_i$

Output: Causal attention output $\mathbf{O}_{n_Q} \in \mathbb{R}^{B_M \times D}$

237 example, we detail the workflow in Algorithm 1 in two
238 phases.

239 **Phase 1.** Starting from the left columns in head dimen-
240 sion. Low-precision attention tiles are computed over the
241 lower triangular region using low-bit quantized query and

key matrices. We follow the standard OnlineSoftmax to in- 242
crementally accumulate the intermediate output \mathbf{O}_i . The 243
first phase terminates after completing the computation for 244
the first $\lceil (n_Q B_M - T) / B_N \rceil$ tokens. 245

Phase 2. Key/value tiles between $\lceil (n_Q B_M - T) / B_N \rceil$ 246
and $\lceil n_Q B_M / B_N \rceil$ will be processed in this phase. To ob- 247
tain a more precised attention weight, $\tilde{Q}_{n_Q}^h, \mathbf{K}_i^h$ of the high- 248
precision copy will be used. In causal attention, Δ_i masks 249
scores within the current tile that correspond to the upper 250
triangular region of the global attention matrix, where \mathcal{I}_{B_M} 251
and \mathcal{I}_{B_N} are indices within each tile. It ensures that no 252
query attends to future positions globally. Final outputs are 253
normalized by the accumulated softmax scaling factors \mathbf{l}_i . 254

This diagonal-tiled execution maintains the performance 255
and memory benefits of quantized attention while preserv- 256
ing the numerical precision of the attention mechanism to 257
keep the generation quality. 258

Compatibility with Non-Causal Attention While Algo- 259
rithm 1 illustrates the workflow of causal attention, our 260
method is also compatible with non-causal attention. In 261
the non-causal attention workflow, the low-bit computation 262
covers both the lower and upper triangular regions of the 263
attention matrix, excluding the diagonal window. 264

The structure of the workflow remains similar to that of 265
the causal case, but the iteration range of Phase 1 differs: the 266
tile index i spans two disjoint intervals. Specifically, tiles 267
from 1 to $\lceil (n_Q B_M - T) / B_N \rceil$ corresponds to the lower 268
triangular region, and tiles from $\lceil (n_Q B_M + T) / B_N \rceil$ 269
to $\lceil n_Q B_M / B_N \rceil$ covers the upper triangular. During this 270

271 phase, the attention weight \mathbf{S}_i is computed using the low-
 272 precision copies of the query and key matrices. In Phase
 273 2, tokens located within the diagonal window $[\lceil (n_Q B_M -$
 274 $T/2)/B_N \rceil, \lceil (n_Q B_M + T/2)/B_N \rceil]$ are reprocessed using
 275 high-precision copies of the inputs. The results are then ac-
 276 cumulated into \mathbf{O}_i computed in Phase 1. Once all key/value
 277 tiles have been iterated over, the final attention output is cal-
 278 culated by $\text{diag}(l_i)^{-1} \mathbf{O}_i$.

279 This partitioning is valid due to the equivalence under
 280 column-wise transformations in matrix multiplication.
 281 Specifically, the dot-product attention is a linear operation
 282 with respect to the key-value axis, and the output remains
 283 mathematically correct as long as all tiles of the attention
 284 matrix are covered and accumulated appropriately. Hence,
 285 splitting the computation into lower, upper triangular and
 286 diagonal regions does not compromise correctness, pro-
 287 vided the transformations are applied consistently across
 288 tiles.

289 5.3. Dual MXFP Quantization Kernel Fusion

290 In this section, we describe the fused mixed-precision quan-
 291 tization kernel used in DMA. The goal is to convert the in-
 292 put tensor into both low-bit and high-bit MXFP representa-
 293 tions within a single fused pipeline, so that the subsequent
 294 mixed-precision attention kernel can directly consume the
 295 quantized outputs without launching another kernels. The
 296 overall procedure is summarized in Algorithm 2.

297 **Pre-process softmax scale (Step 1).** Before quantization,
 298 we first apply the standard softmax scaling factor to the
 299 query tensor. Specifically, when the input is Q , we multi-
 300 ply it by $\log_2 e / \sqrt{D}$, where D is the head dimension.
 301 This incorporates the softmax normalization factor into the
 302 quantized computation in advance, so that the subsequent
 303 QK^\top accumulation is already aligned with the scaled at-
 304 tention score formulation. Since our kernel computes the
 305 matrix product in base-2 arithmetic, folding this factor into
 306 the query in advance avoids an extra scaling step after accu-
 307 mulation and simplifies the fused implementation.

308 **Quantization scale (Step 2).** We fuse quantization with
 309 MXFP number format conversion at the kernel level. Fol-
 310 lowing the two-level scaling strategy in SageAttention3, we
 311 note that the shared scale in NVFP4 uses the FP8 (E4M3)
 312 format, which ranges from $[-448, 448]$, while each element
 313 is represented in FP4, with a dynamic range of $[-6, 6]$.
 314 Therefore, the representable range of NVFP4 is the product
 315 of these two bounds. To fully utilize the available dynamic
 316 range and accommodate potential outliers without resorting
 317 to clipping, we compute the quantization scale in Step 2
 318 (line 4) of Algorithm 2, scaling the original tensor into this
 319 representable range prior to quantization.

Algorithm 2 Fused Mixed-Precision Quantization

Input: $\mathbf{X} \in \mathbb{R}^{B \times D}$: input FP16 tensor. B : size of a tile
 in sequence dimension, D : size of head dimension. V_1 ,
 V_2 : sizes of a block for low-/high-precision MXFP for-
 mats. l_1, l_2, u_1, u_2 : lower/upper bound for both ele-
 ments formats. e^{\max} : the exponent of the largest nor-
 mal number of the element format.

Step 1: Pre-process softmax scale
 1: **if** is Query **then**
 2: $\mathbf{X}_{\text{sm-scaled}} \leftarrow \mathbf{X} \cdot \frac{\log_2 e}{\sqrt{D}}$
 3: **end if**
Step 2: Compute the quantization scale
 4: $\mathbf{S}_q \leftarrow \max_D (|\mathbf{X}_{\text{sm-scaled}}|) / (448 \times 6)$
 5: $\mathbf{X}_{\text{scaled}} \leftarrow \mathbf{X}_{\text{sm-scaled}} / \mathbf{S}_q$
Step 3: Compute shared scale for low precision format
 6: $\mathbf{X}'_{\text{scaled}} \leftarrow \mathbf{X}_{\text{scaled}} \triangleright$ reshape to $[B, D // V_1, V_1]$
 7: $\mathbf{S}_{\text{FP4}} \leftarrow \max_D (|\mathbf{X}'_{\text{scaled}}|) / u_1$
 8: $\mathbf{X}_{\text{clamped}} \leftarrow \text{clamp}(\mathbf{X}'_{\text{scaled}} / \mathbf{S}_{\text{FP4}}, l_1, u_1)$
Step 4: Encode $\mathbf{X}_{\text{clamped}}$ to E2M1
 9: $\mathbf{X}_{\text{FP4}} \leftarrow \text{Algorithm 3}(\mathbf{X}_{\text{clamped}})$
Step 5: Pack two FP4 into one UINT8 along D
 10: **if** pack along last dimension **then**
 11: $\mathbf{X}'_{\text{FP4}} \leftarrow \mathbf{X}_{\text{FP4}} \triangleright$ reshape to $[B, (D + 1) // 2, 2]$
 12: $\mathbf{L}, \mathbf{H} \leftarrow \mathbf{X}'_{\text{FP4}}[:, :, 0], \mathbf{X}'_{\text{FP4}}[:, :, 1]$
 13: $\mathbf{X}_{\text{packed}} \leftarrow (\mathbf{H} \ll 4) \mid \mathbf{L}$
 14: **end if**
Step 6: Compute shared scale for high precision format
 15: $\mathbf{X}''_{\text{scaled}} \leftarrow \mathbf{X}_{\text{scaled}} \triangleright$ reshape to $[B, D // V_2, V_2]$
 16: $\mathbf{S}_{\text{shared}} \leftarrow \lfloor \log_2 (\max_D |\mathbf{X}''_{\text{scaled}}|) \rfloor - e^{\max}$
 17: $\mathbf{X}_{\text{FP8}} \leftarrow \text{clamp}(\mathbf{X}''_{\text{scaled}} / 2^{\mathbf{S}_{\text{shared}}}, l_2, u_2)$
Step 7: Convert shared scale into E8M0
 18: $\mathbf{S}_{\text{FP8}} \leftarrow \text{clamp}(\mathbf{S}_{\text{shared}} + 127, 0, 254)$
Output: Packed FP4 tensor $\mathbf{X}_{\text{packed}}$, FP8 tensor \mathbf{X}_{FP8} ,
 shared scale of NVFP4 \mathbf{S}_{FP4} , shared scale of MXFP8
 \mathbf{S}_{FP8} , quantization scale \mathbf{S}_q

Compute shared scale (Step 3 and Step 6). We first re-
 shape the scaled input tensor along the packing dimension
 into a new matrix, where each row groups the V_1 or V_2 ele-
 ments that share a common scale in the MXFP format. For
 NVFP4, the shared scale is computed directly as the ab-
 solute maximum of each group. For MXFP4 and MXFP8
 formats, however, the scale is represented in integer E8M0
 format. To maximize the representable range, we need to
 normalize the input exponent so that the largest exponent in
 the data aligns with the maximum representable exponent
 (denoted as e^{\max}) of the low-bit element format. Conse-
 quently, the shared exponent $\mathbf{S}_{\text{shared}}$ stores this offset from
 the input's exponent to e^{\max} . In E5M2, the maximum ex-
 ponent e^{\max} is 15 (i.e., $(11110)_2 = 30$ with a bias of 15,
 excluding reserved patterns with all 1 bits to represent in-

Algorithm 3 Encode FP16 tensor into E2M1 format**Input:** $\mathbf{X} \in \mathbb{R}_{\text{FP16}}^{B \times D} \in [-6.0, 6.0]$: clamped FP16 tensor.**Step 4.1: Extract signbit**1: $\mathbf{S} \leftarrow \text{sign}(\mathbf{X})$ **Step 4.2: Compute 2-bit exponent**2: $\mathbf{X}_{\text{abs}} \leftarrow |\mathbf{X}|$ 3: $\mathbf{E} \leftarrow \mathbb{I}[\mathbf{X}_{\text{abs}} \geq 1] + \mathbb{I}[\mathbf{X}_{\text{abs}} \geq 2] + \mathbb{I}[\mathbf{X}_{\text{abs}} \geq 4]$ **Step 4.3: Compute 1-bit mantissa**4: $\text{bias} \leftarrow 1$ 5: $\mathbf{X}_{\text{norm}} \leftarrow \mathbf{X}_{\text{abs}} / 2^{\mathbf{E} - \text{bias}}$ 6: $\mathbf{M} \leftarrow \mathbb{I}[\mathbf{E} = 0] \cdot \mathbb{I}[\mathbf{X}_{\text{norm}} > 0.25] + \mathbb{I}[\mathbf{E} \neq 0] \cdot \mathbb{I}[\mathbf{X}_{\text{norm}} > 1.25]$ **Step 4.4: Assemble quantized FP4 integer**7: $\mathbf{X}_{\text{FP4}} \leftarrow (\mathbf{S} \lll 3) \mid (\mathbf{E} \lll 1) \mid \mathbf{M}$ **Output:** Quantized tensor $\mathbf{X}_{\text{FP4}} \in \mathbb{R}_{\text{FP4}}^{B \times D}$

335 finite and NaN). In E4M3, $e^{\max} = 8$ (i.e., $(1111)_2 = 15$
 336 with a bias of 7). Notably, E4M3 does not strictly follow
 337 IEEE-754, its normal numbers allows exponent bits set to
 338 all 1. The division of exponent shared scale allows full uti-
 339 lization of the exponent range within the limited bit budget
 340 of MXFP format.

341 **Encoding of FP4 format (Step 4).** Algorithm 3 shows
 342 the steps to quantize an FP16 tensor into E2M1 format (1-
 343 bit sign, 2-bit exponent, 1-bit mantissa). We begin by ex-
 344 tracting the sign bit (Step 4.1). The 2-bit exponent is de-
 345 termined by thresholding the absolute value of the input
 346 against $\{1, 2, 4\}$, yielding exponent values in $\{0, 1, 2, 3\}$,
 347 which can be represented by 2 bits (Step 4.2). We employ
 348 the Kronecker delta function $\mathbb{I}[\cdot]$ in line 3 and 6, which re-
 349 turns 1 if the condition is true and 0 otherwise. In this way,
 350 subnormal values are assigned $E = 0$. To compute man-
 351 tassas (Step 4.3), We normalize the input by the implied
 352 scale factor and assign the mantissa bit to 1 if \mathbf{X}_{norm} ex-
 353 ceeds 0.25 (the midpoint of 0 and 0.5). For normal values
 354 ($E \neq 0$), the mantissa is set to 1 if \mathbf{X}_{norm} exceeds the mid-
 355 point of the two representable values under that exponent.
 356 For example, when $E = 3$, E2M1 can represent 4 and 6,
 357 which correspond to normalized values of 1 and 1.5, making
 358 1.25 the comparison threshold for mantissa. To implement
 359 `roundTiesToEven` according to the standard, we prefer
 360 rounding to even mantissas (i.e., $M = 0$) in tie-breaking
 361 scenarios. For example, for input value 5, we prefer round-
 362 ing to 4 to have mantissa bit as 0, rather than 6. So the
 363 comparison to midpoint threshold should be greater but not
 364 equal. In final, step 4 constructs the 4-bit representation us-
 365 ing bitwise shift and OR operations.

366 **Packing (Step 5).** We encode two FP4 values into a single
 367 byte, assigning the value with the higher index to the most

significant 4 bits and the other to the least significant 4 bits.
 This compact representation improves memory bandwidth
 utilization.

Scale Conversion into E8M0 (Step 7). The shared scale
 for MXFP8 should be in E8M0 (unsigned 8-bit integer) for-
 mat according to the official document of MXFP comput-
 ing. Therefore, we add 127 before clamping it into 0 to
 254, ensuring it lies within the valid exponent range of E8.

6. Experiments

We conduct a comprehensive evaluation of the proposed
 DMA operator in terms of both accuracy and efficiency. We
 report results across a range of tasks and datasets to assess
 the overall performance. In addition, we perform detailed
 ablation studies to analyze the impact of key design choices,
 including numerical precision, quantization granularity, di-
 agonal window size, and kernel fusion strategies.

6.1. Settings

Models and Datasets. We evaluate our method using
 LLaMA-3.1-8B and LLaMA-3.2-3B. Performance is mea-
 sured on LongBench [2], which focuses on long-context
 understanding with sequence lengths ranging from 2.5K to
 30K tokens. We report results across various subtasks to as-
 sess general long-context capabilities. For all reported met-
 rics, higher values indicate better performance.

Implementation. We implement DMA in Triton. We
 compare our method with the native attention kernel, which
 refers to the SDPA kernel originally supported in `torch`
 and computed in BF16 format [11]. All experiments are
 conducted on a single NVIDIA B200 GPU.

6.2. Accuracy

Performance on LongBench We evaluate our method
 on LongBench to assess long-context language understand-
 ing performance. Table 3 shows that our method im-
 proves the average score over the native attention base-
 line on both LLaMA3.1-8B and LLaMA3.2-3B. Specifi-
 cally, for LLaMA3.1-8B, the average score improves from
 44.11 to 46.43. For LLaMA3.2-3B, the average score im-
 proves from 35.84 to 37.20. The gains are broad across
 many tasks, including `repobench-p`, `samsun`, `trec`,
 and `triviaqa`. In particular, `repobench-p` improves
 from 42.23 to 54.00 on LLaMA3.1-8B and from 44.14 to
 50.15 on LLaMA3.2-3B. The average results across both
 models indicate that our method preserves, and in most
 cases improves, long-context accuracy compared with the
 native implementation.

Table 3. Comparison of attention implementations on LLaMA3.2-3B and LLaMA3.1-8B. “Native” refers to SDPA implementation supported by PyTorch.

Task	LLaMA3.1-8B		LLaMA3.2-3B	
	Native	Ours	Native	Ours
2wikimqa	39.15	32.54	30.75	36.33
dureader	28.40	32.93	28.11	33.04
gov_report	33.87	34.89	32.26	33.04
hotpotqa	49.89	50.49	46.18	48.59
lcc	49.46	56.96	42.98	45.54
lsht	41.50	44.50	26.75	30.75
multi_news	26.16	27.07	25.56	21.78
multifieldqa_en	52.26	52.65	46.84	47.87
multifieldqa_zh	57.57	58.81	48.06	53.62
musique	25.82	26.81	20.35	25.77
narrativeqa	25.28	26.70	19.76	25.44
passage_count	4.02	6.22	3.40	1.50
passage_retrieval_en	99.00	96.00	80.00	37.00
passage_retrieval_zh	92.14	85.50	8.25	10.50
qasper	42.64	43.05	32.13	39.62
qmsum	23.47	24.67	22.52	23.63
repobench-p	42.23	54.00	44.14	50.15
samsum	36.60	43.99	34.79	42.74
trec	65.00	71.50	62.50	69.50
triviaqa	77.01	87.89	83.68	88.13
vcsum	14.80	17.92	13.65	16.69
Avg.	44.11	46.43	35.84	37.20

413 6.3. Efficiency

414 Table 4 reports the latency breakdown of different formats
415 and block-scale configurations. We compare our imple-
416 mentation with several fixed-format baselines, including
417 MXFP4, NVFP4, and MXFP8. For each setting, we report
418 the attention time, quantization overhead, and total runtime.

419 Among all evaluated settings, our configuration with di-
420 agonal and sink sizes set to 128 achieves the lowest total
421 latency, at 7.776 ms. This is lower than MXFP4 (12.980
422 ms), NVFP4 (13.404 ms), and MXFP8 (16.771 ms). In par-
423 ticular, the main reduction comes from the attention kernel
424 time, which is 7.110 ms in our 128/128 setting, compared
425 with 12.491 ms, 12.941 ms, and 16.480 ms for the three
426 baselines, respectively. We also evaluate a larger block-
427 scale configuration with diagonal and sink sizes set to 256.
428 In this case, the total latency increases to 15.720 ms. Com-
429 pared with the 128/128 setting, this suggests that a larger
430 block size is less efficient in our current implementation.

431 6.4. Ablation Study

432 **Mixed-Precision Block Tile Sizes.** We conduct ablation
433 experiments to evaluate the impact of different block tile

Table 4. Latency breakdown of different block-scale types and configurations. “MP Size” denotes the mixed-precision block size used for the higher-bit diagonal and sink blocks.

Format	MP Size	Attn (ms)	Quant (ms)	Total (ms)
MXFP4	–	12.491	0.242	12.980
NVFP4	–	12.941	0.204	13.404
MXFP8	–	16.480	0.044	16.771
Ours	128	7.110	0.382	7.776
Ours	256	15.056	0.382	15.720

Table 5. Similarity metrics under different token numbers for diagonal and sink windows.

Diag.	Sink	Bit _{high} (%)	Cos Sim ↑	Rel. L1 ↓	RMSE ↓	PSNR ↑
-	-	0.0	0.778	0.620	0.065	43.715
-	-	100.0	0.819	0.547	0.059	44.568
0	128	1.15	0.781	0.780	0.072	42.817
128	0	1.15	0.782	0.644	0.066	43.635
128	128	2.30	0.822	0.539	0.059	44.657
512	512	9.22	0.826	0.542	0.058	44.731
2048	2048	36.87	0.852	0.521	0.054	45.352

Table 6. Ablation study of kernel fusion components to analyze their effect on throughput (TOPS) and latency. **Encode**: encoding FP16 to FP4/FP8 MX format. **Pack**: packing two FP4 values to one UINT8. **Scale Cvt.**: converting microscaling scalar to E8M0 format. **MP**: fusing quantization for mixed bitwidth to a single kernel.

Encode	Pack	Scale Cvt.	MP	$L = 2k$ (μs)	$L = 8k$ (μs)
\times	\times	\times	\times	7262.41	22628.96
\checkmark	\times	\times	\times	802.90	1113.77
\checkmark	\checkmark	\times	\times	740.64	942.67
\checkmark	\checkmark	\checkmark	\times	179.97	299.69
\checkmark	\checkmark	\checkmark	\checkmark	97.87	282.46

434 sizes on the similarity between the quantized attention ma-
435 trix and its full-precision counterpart. As shown in Tab. 5,
436 we report the similarity before and after quantization to
437 show the representation errors under varying diagonal and
438 sink tile configurations, including Cosine Similarity, Rel-
439 ative L1 Distance, Root Mean Square Error (RMSE), and
440 Peak Signal-to-Noise Ratio (PSNR). The Bit_{high}% column
441 denotes the percentage of values in the attention matrix that
442 compute in high-precision.

443 We observe that increasing the tile size from 128 to 512
444 improves similarity metrics marginally. But it leads to a
445 significant degradation in throughput due to reduced paral-
446 lelism and slower computation as shown in Tab. 4. Based
447 on this trade-off, we use the 128/128 configuration as the
448 default setting in the following experiments.

449 **Kernel Fusion** We evaluate the contribution of each ker-
450 nel fusion component in Tab. 6. The fully unfused baseline,
451 which executes all quantization-related steps separately, in-

Table 7. Latency breakdown of the non-fused MX encoding pipeline and our fused implementation.

Operator	Time	Time (%)
Not fused	18.320 ms	–
- Element encoding	17.742 ms	100.0%
MinOps	2.110 ms	11.89%
ArgMinOps	2.054 ms	11.58%
Direct_Copy	1.256 ms	7.08%
CompareEq	890.18 us	5.02%
AddOps	810.56 us	4.57%
MulFuncator	755.10 us	4.26%
Memcpy / Memset	206.85 us	1.17%
- Element packing	81.22 us	100.0%
BitwiseOr	23.74 us	29.24%
lshift	22.63 us	27.86%
- Scalar Convert	521.47 us	100.0%
IndexOps	35.75 us	6.85%
DeviceSelectSweep	25.73 us	4.93%
Write_Indices	20.42 us	3.91%
Direct_Copy	17.12 us	3.28%
Memcpy	13.57 us	2.60%
Kernel Fusion (Ours)	116.480 us	–

Table 8. Performance comparison under different quantization granularities. Latency are reported in milliseconds.

Granu.	Latency ↓	Cos Sim ↑	Rel. L1 ↓	RMSE ↓	PSNR ↑
Per-Tensor	6.276 ms	0.732	0.560	0.067	43.479
Per-Block	6.366 ms	0.736	0.558	0.067	43.531
Per-Token	7.131 ms	0.822	0.539	0.059	44.657

452 curs a latency of 7262.41 μ s for sequence length $L = 2k$
 453 and 22628.96 μ s for $L = 8k$. When only in-kernel FP16-to-
 454 MX encoding is enabled, the latency is reduced to 802.90 μ s
 455 for $L = 2k$ and 1113.77 μ s for $L = 8k$. After further
 456 enabling FP4 packing, which packs two FP4 values into
 457 one `uint8`, the latency is further reduced to 740.64 μ s for
 458 $L = 2k$ and 942.67 μ s for $L = 8k$. When scale conversion
 459 of microscaling factors to the `e8m0` format is also fused into
 460 the kernel, the latency drops substantially to 179.97 μ s for
 461 $L = 2k$ and 299.69 μ s for $L = 8k$. Finally, after integrating
 462 mixed-precision quantization into a single fused kernel, de-
 463 noted as **MP**, the latency is further reduced to 97.87 μ s for
 464 $L = 2k$ and 282.46 μ s for $L = 8k$, achieving the best over-
 465 all performance. Compared with the fully unfused baseline,
 466 this corresponds to a 74.2 \times speedup for $L = 2k$ and an
 467 80.1 \times speedup for $L = 8k$.

468 Overall, these results demonstrate that each fusion com-
 469 ponent contributes to latency reduction, and a fully fused
 470 kernel is critical for achieving an efficient mixed-precision
 471 attention.

Quantization Granularity. We further investigate the
 472 impact of different quantization granularities on both la-
 473 tency and kernel output precision. As shown in Tab. 8, we
 474 evaluate four granularity settings: Per-Tensor, Per-Block,
 475 and Per-Token. Latency is measured using the 5 warmups
 476 and average of 10 runs. The results show that finer granu-
 477 larity, such as Per-Token, achieves the highest output simi-
 478 larity before and after quantization, as indicated by higher
 479 Cosine Similarity (0.822), lower RMSE (0.059), and higher
 480 PSNR (44.657), but it also incurs the highest latency (7.131
 481 ms). In contrast, Per-Tensor and Per-Block offer lower la-
 482 tency but exhibit larger quantization errors. Therefore, the
 483 choice of quantization granularity depends on the task re-
 484 quirement, specifically whether latency or output quality is
 485 prioritized.
 486

7. Conclusion 487

In this paper, we present DMA, a diagonal-tiled mixed-
 488 precision attention method for low-bit MXFP attention
 489 computation for efficient LLM inference. DMA improves
 490 the accuracy of low-bit attention by preserving more sensi-
 491 tive diagonal regions in higher precision, while keeping the
 492 remaining regions in efficient low-bit formats. Combined
 493 with a fused dual-MXFP quantization kernel, our design re-
 494 duces quantization preprocessing overhead and maintains
 495 practical efficiency. Experimental results show that DMA
 496 provides a better balance between accuracy and latency.
 497

Limitation 498

Our current study is limited to a small set of model sizes and
 499 benchmark settings, and the evaluation is mainly focused on
 500 long-context text workloads. In particular, we mainly vali-
 501 date DMA on text-based tasks and do not extend the experi-
 502 ments on vision or vision-language settings. In addition, the
 503 mixed-precision tiling policy is not validated on extremely
 504 long sequence lengths, diverse hardware settings, or other
 505 model architectures and attention variants. Extensive ex-
 506 periments and validations are important directions for our
 507 future work.
 508

509

References

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

[15] Jintao Zhang, Jia Wei, Pengle Zhang, Xiaoming Xu, 566

Haofeng Huang, Haoxu Wang, Kai Jiang, Jianfei Chen, and 567

Jun Zhu. Sageattention3: Microscaling fp4 attention for in- 568

ference and an exploration of 8-bit training, 2026. 1 569