

Flash Attention-3: Fast and Accurate Attention with Asynchrony and Low-precision

TIMELINE

Flash Attention 1

Flash Attention 2

Flash Attention 3

2022.10

2023.6

2024.3

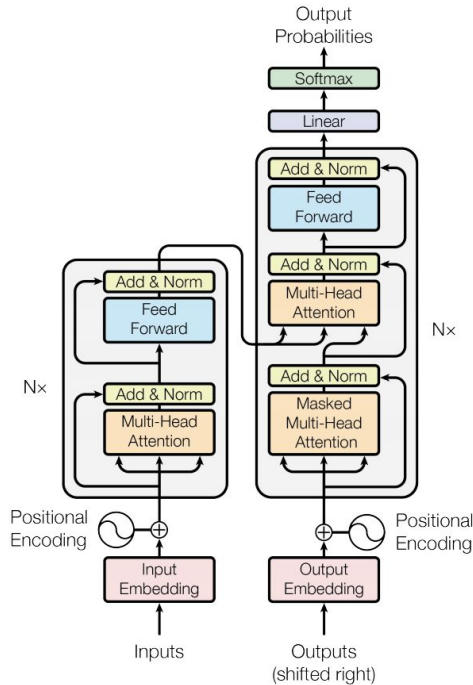
*Flash Attention:
Fast and
Memory-Efficient
Exact Attention
with IO-Awareness*

*Flash Attention-2:
Faster Attention
with Better
Parallelism and
Work Partitioning*

*Flash Attention-3:
Faster Attention
with Better
Parallelism and
Customization*

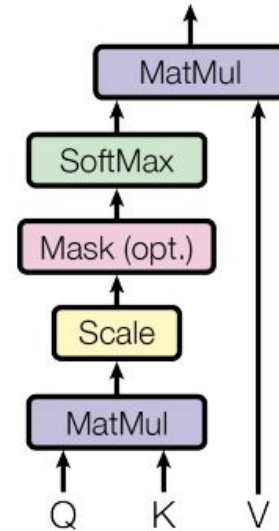
Attention

- Transformer



[Image source: Vaswani et al. (2017)]

Scaled Dot-Product Attention



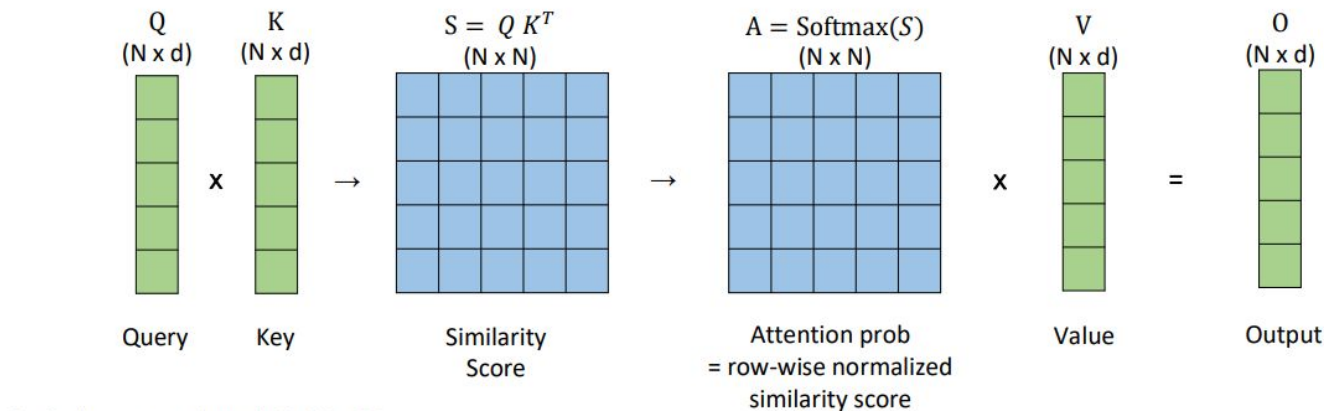
- Scaled Dot-Product Attention**

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Attention

- The entire sequence of length N must be stored for the Softmax

Background: Attention Mechanism



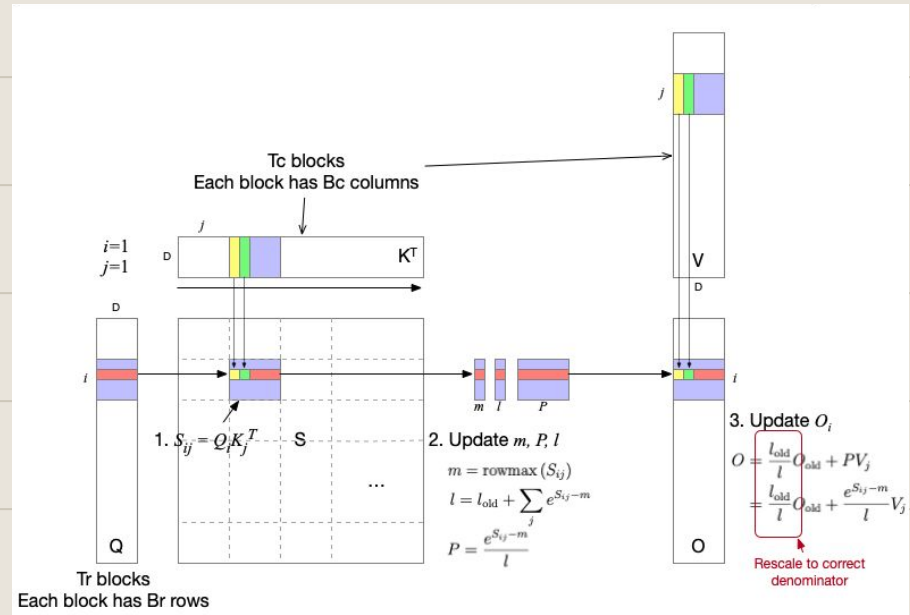
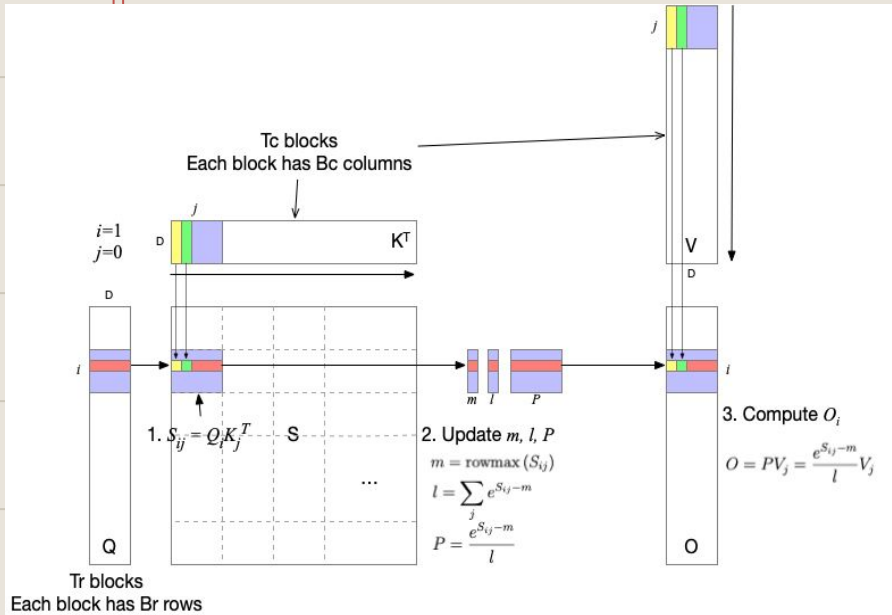
Typical sequence length N : 1K – 8K
Head dimension d : 64 – 128

$$\text{Softmax}([s_1, \dots, s_N]) = \left[\frac{e^{s_1}}{\sum_i e^{s_i}}, \dots, \frac{e^{s_N}}{\sum_i e^{s_i}} \right]$$

$$O = \text{Softmax}(QK^T)V$$

Flash Attention

- Softmax computed using partial scores S , with the max value
- Recomputed and compensated for later.



[Image source: <https://insujang.github.io/2024-01-21/flash-attention/>]

Flash Attention

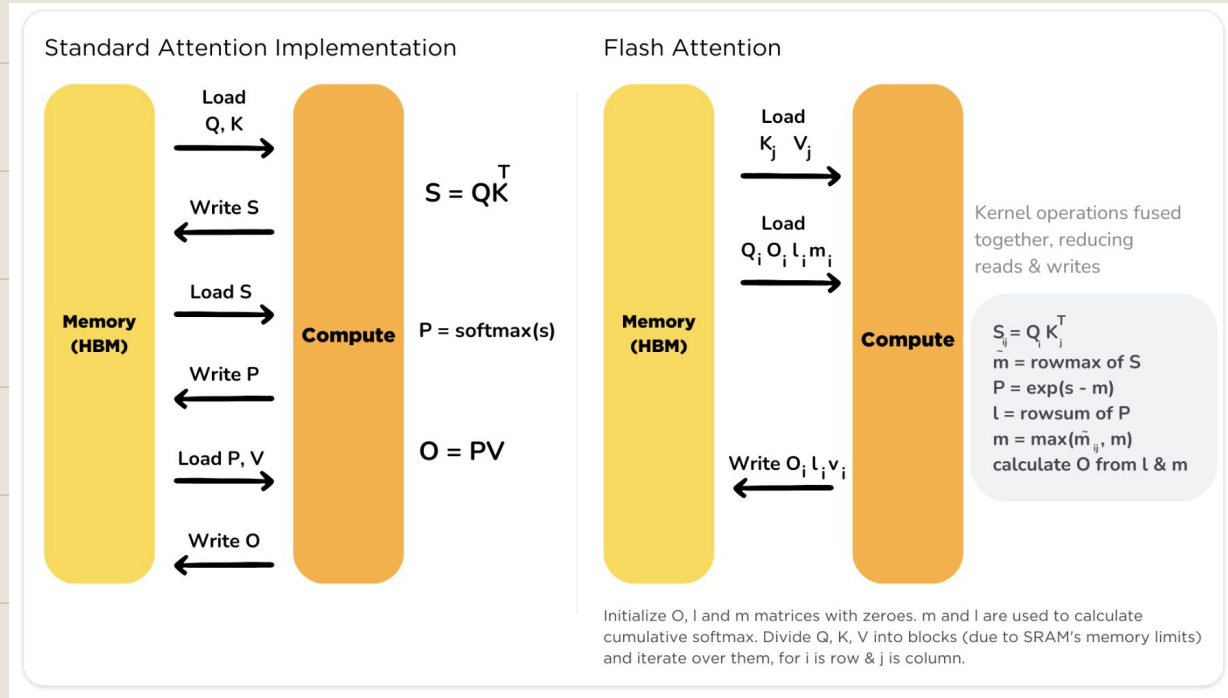
- Recompute O1 with new max M1

Loop 0	Loop 1
$S_0 = QK_0$	$S_1 = QK_1$
$m_0 = \text{rowmax}(S_0)$	$m_1 = \max(\text{rowmax}(S_0), m_0)$
$l_0 = \text{rowsum}(e^{s_0 - m_0})$	$l_1 = l_0 * e^{m_0 - m_1} + e^{s_1 - m_1}$
$P_0 = \frac{e^{s_0 - m_0}}{l_0}$	$P_1 = \frac{e^{s_1 - m_1}}{l_1}$
$O_0 = P_0 V_0$	$O_1 = O_0 * \frac{l_0}{l_1} + P_1 V_1$

[Image source: <https://zhuanlan.zhihu.com/p/607364156>]

Flash Attention

- Benefits : Reducing Read & Write cycle



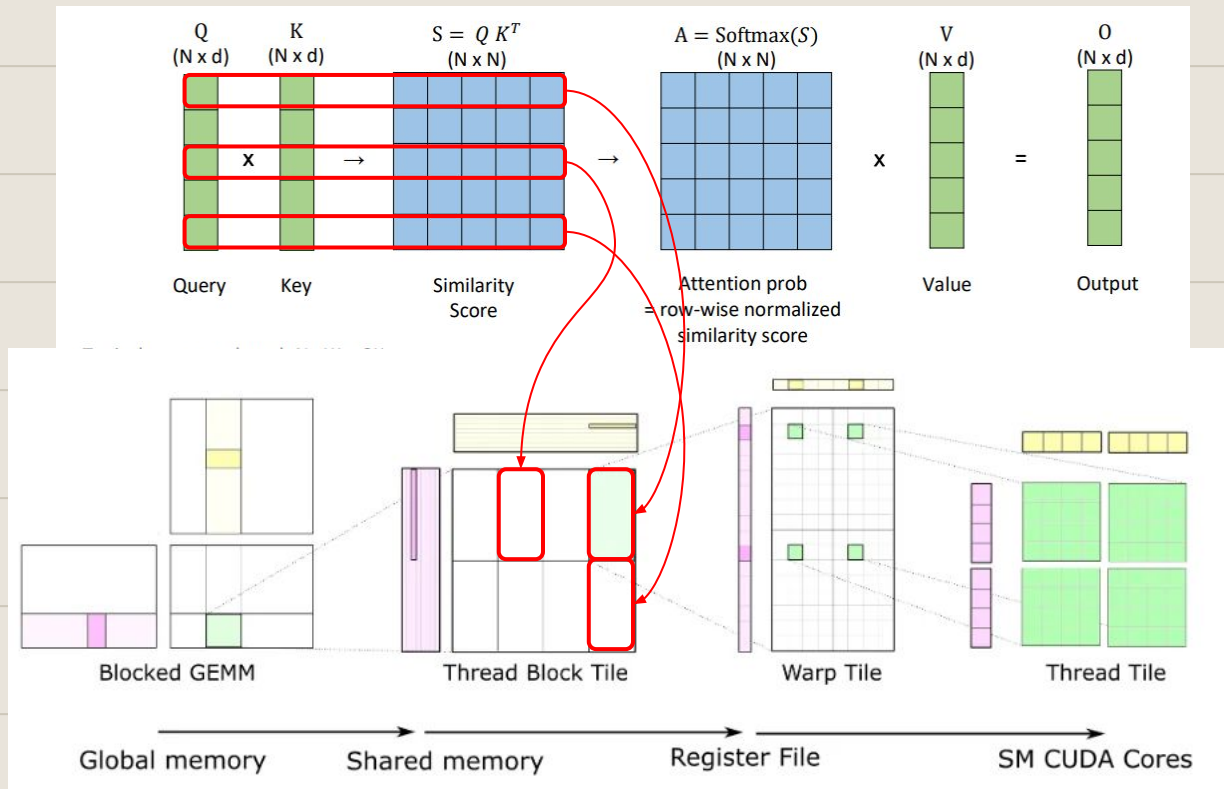
Flash Attention

- Limitation

- **Limited inter-thread block parallelism:** For long sequences, a single CTA (Cooperative Thread Array) handles the computation, which restricts parallelization.
- **Rigid memory layout:** The fixed structure limits flexibility and adaptability to different models or hardware.
- **Performance benefits mainly for long sequences:** Gains are significant only for long sequences; performance is limited for short sequences or small batch sizes.

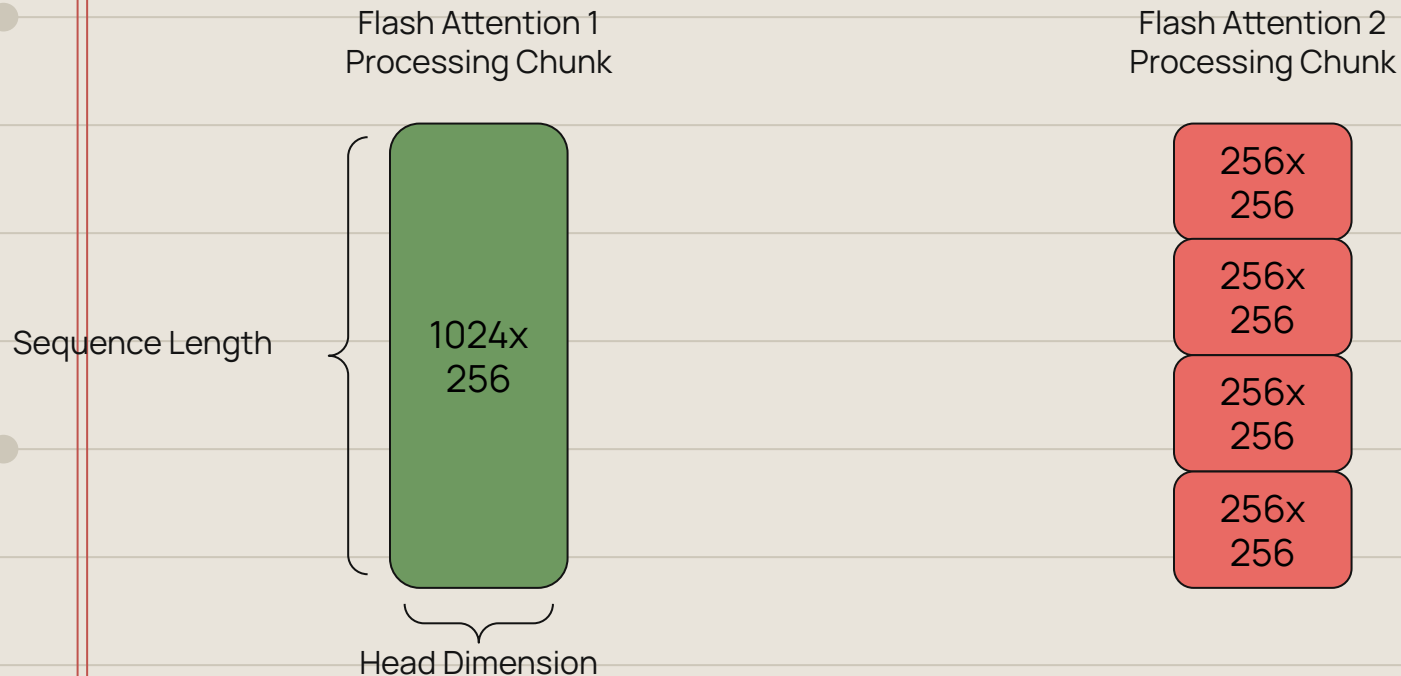
Flash Attention 2

- Main idea : Thread Block parallelism (Global max and update later)



Flash Attention 2

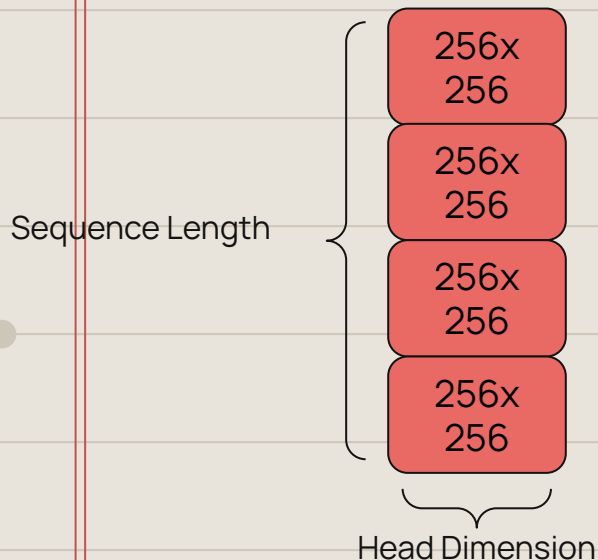
- Main idea : Thread Block parallelism (Global max and update later)



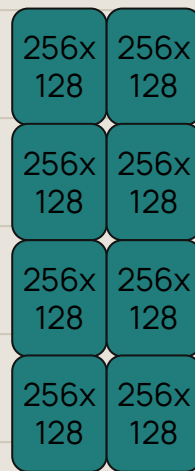
Flash Attention 3

- Main idea : $M(\text{Head dim})$ parallelism + Re-ordering + Block Quantization
- $M(\text{Head dim})$ parallelism

Flash Attention 2
Processing Chunk



Flash Attention 3
Processing Chunk



Flash Attention 3

- Main idea : M parallelism + Re-ordering + Block Quantization
- Re-ordering



Flash attention 1&2

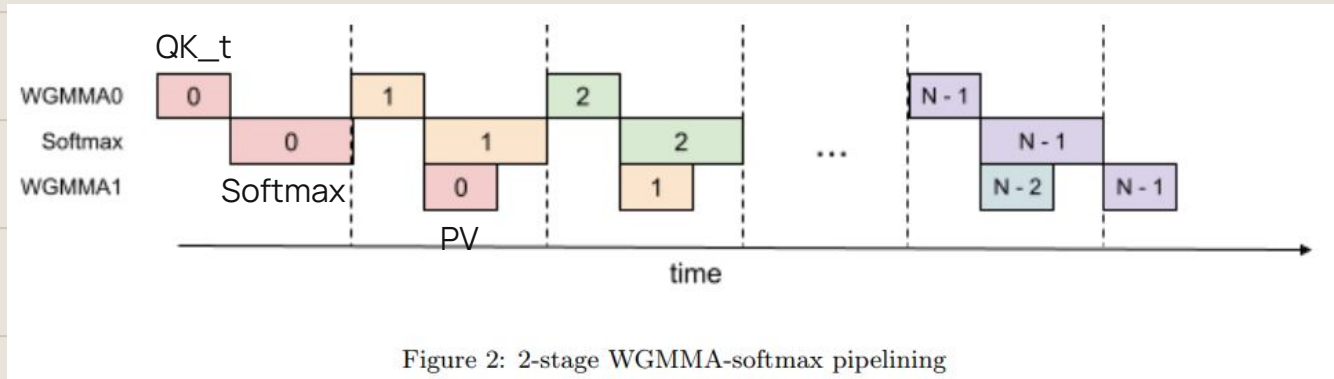
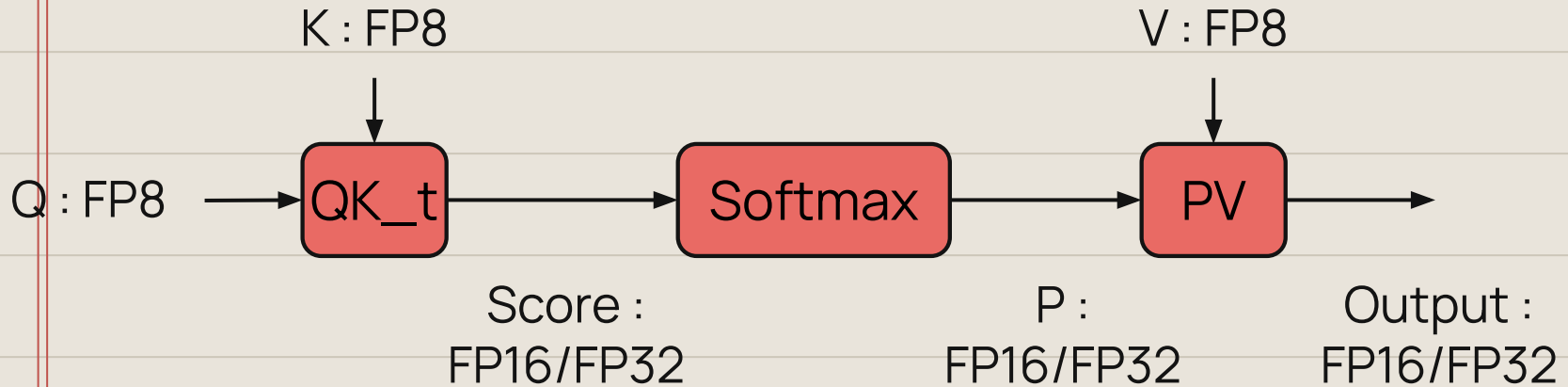


Figure 2: 2-stage WGMMA-softmax pipelining
[Image source: Flash Attention 3]

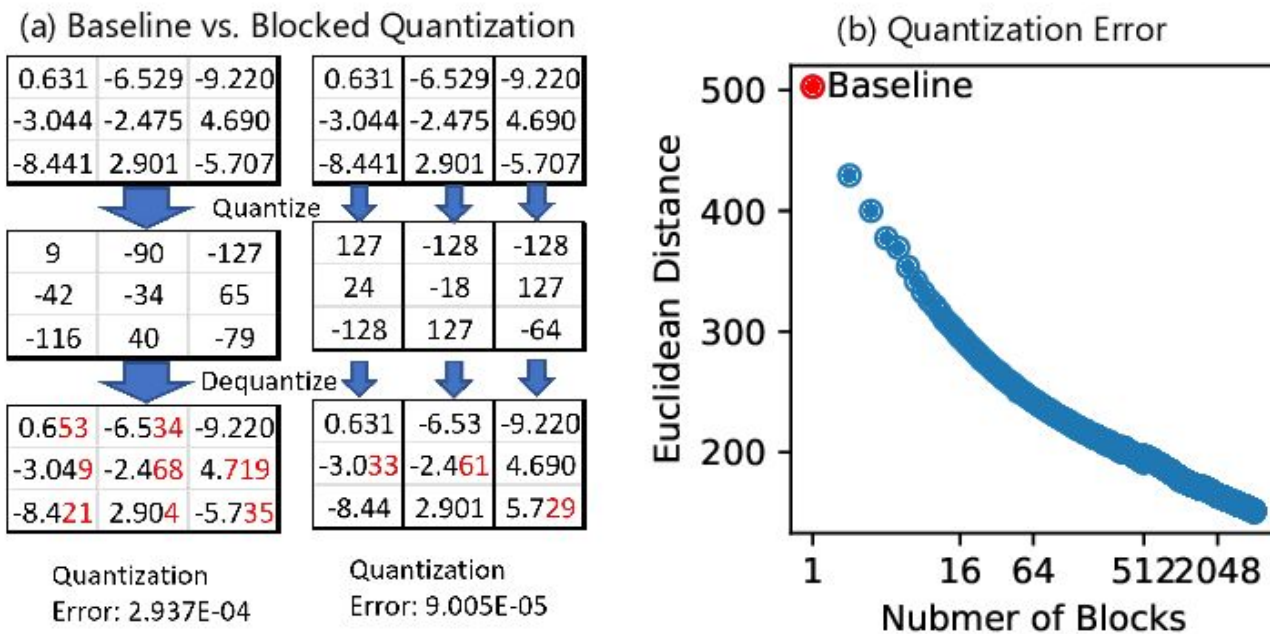
Flash Attention 3

- Main idea : M parallelism + Re-ordering + Block Quantization
- FP8 Block Quantization



Flash Attention 3

- Main idea : M parallelism + Re-ordering + Block Quantization
- FP8 Block Quantization



Flash Attention 3

- Validation
- GPU Optimize library
 - cuDNN : CUDA Deep Neural Network
 - Triton : GPU Programming for Neural Networks

[Image source: Flash Attention - 3]

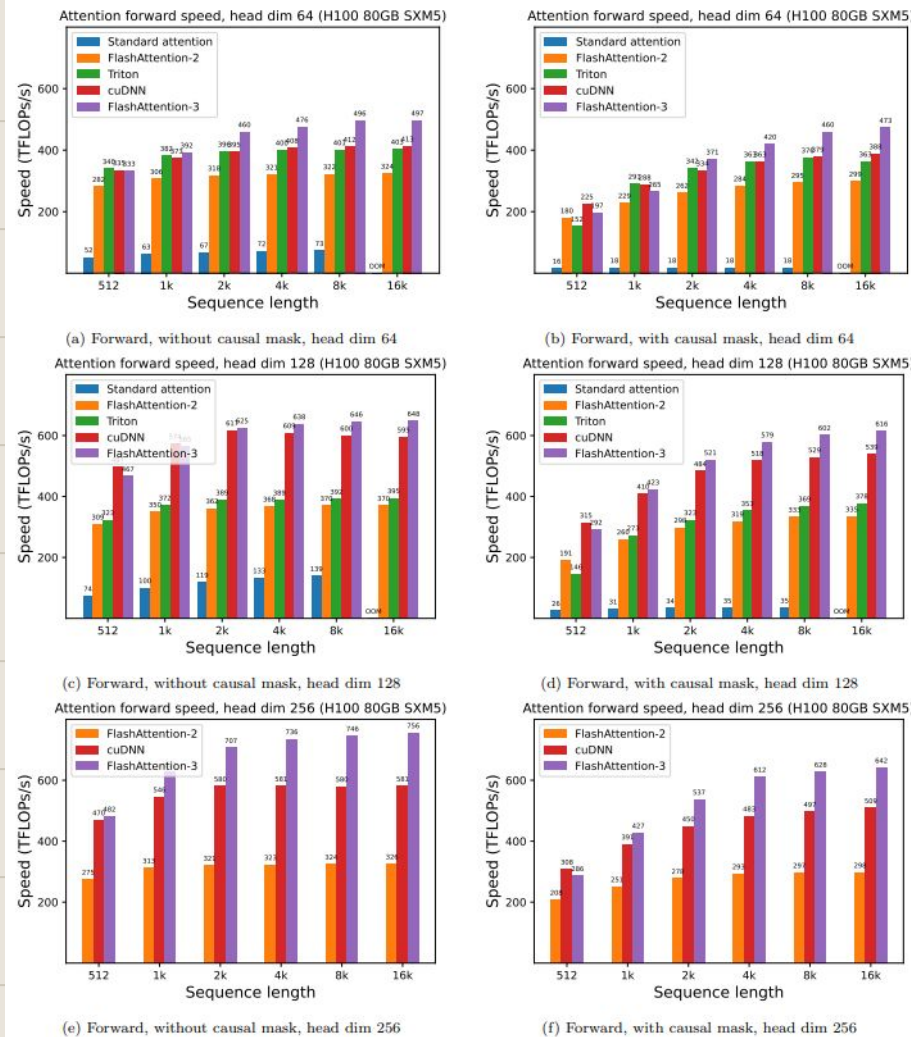


Figure 5: Attention forward speed (FP16/BF16) on H100 GPU

**THANK
YOU!**