

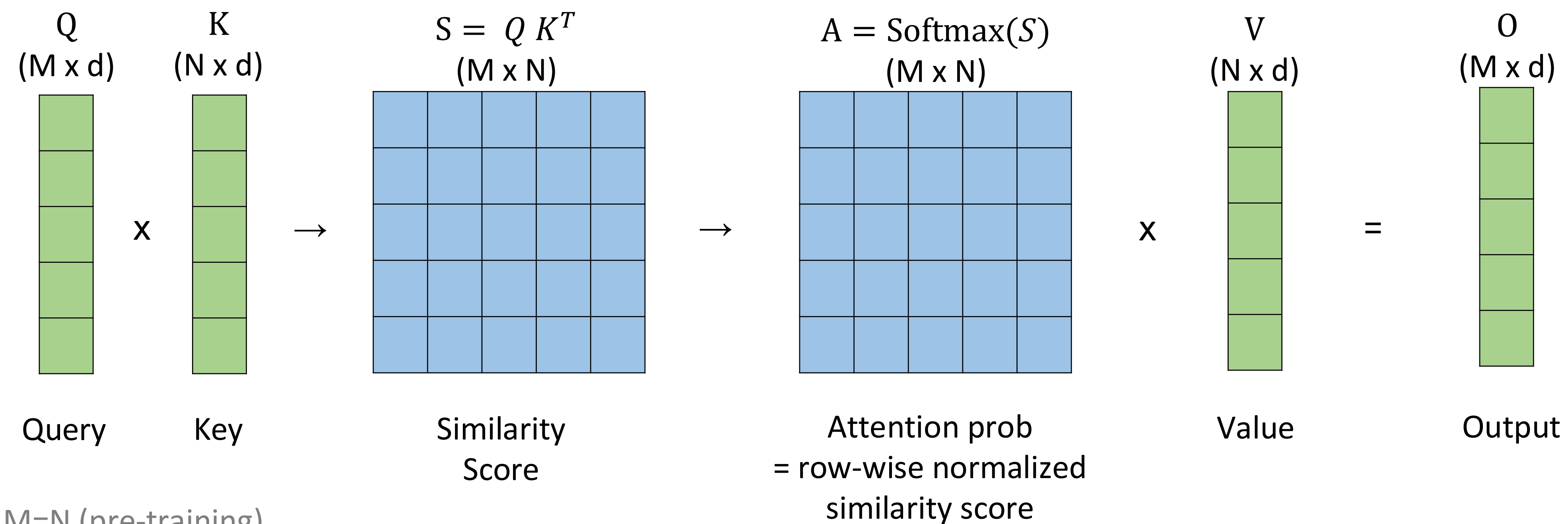
# CUTLASS and FlashAttention-3

Jay Shah, Colfax Research  
<https://research.colfax-intl.com/>

## Outline of talk:

1. Recap of attention and Flash Attention.
2. High-level overview of the FlashAttention-3 algorithm.
3. Translating the algorithm into working code built on CUTLASS.

# Background: Attention Mechanism



Suppose  $M=N$  (pre-training).  
Typical sequence length  $N$ : 1K – 16K  
Head dimension  $d$ : 64 – 256

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

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

Attention scales quadratically in sequence length  $N$

# Naïve (Standard) Attention Algorithm

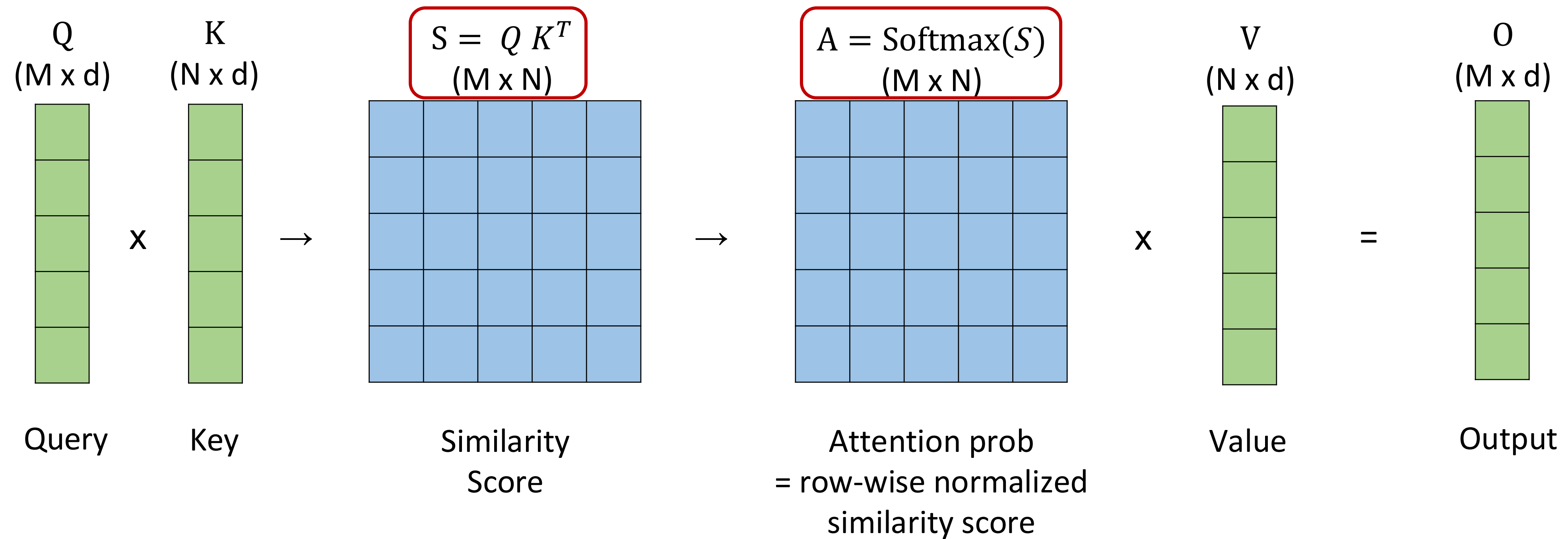
---

**Algorithm 1** Standard Attention

---

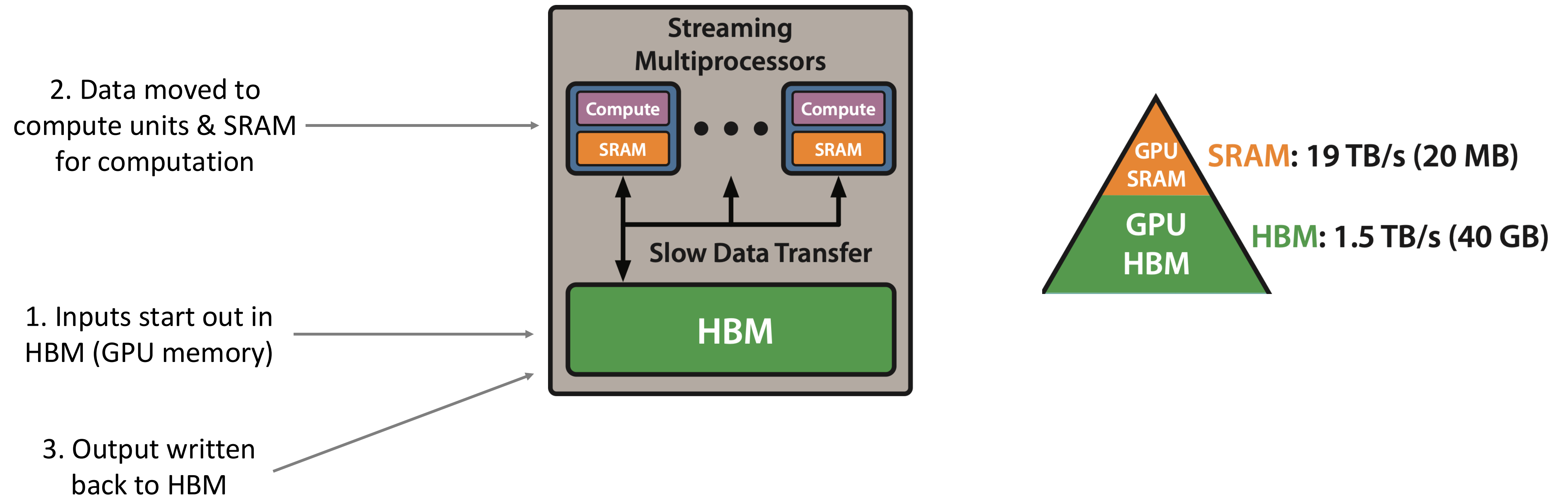
- 1: Load  $Q$  and  $K$  by blocks from HBM.
  - 2: Compute  $S = (1/\sqrt{d})QK^T$  (GEMM-I).
  - 3: Write  $S$  to HBM.
  - 4: Read  $S$  from HBM.
  - 5: Compute  $S = S - \text{rowmax}(S)$ .
  - 6: Compute  $P = \text{softmax}(S)$ .
  - 7: Write  $P$  to HBM.
  - 8: Load  $P$  and  $V$  by blocks from HBM.
  - 9: Compute  $O = PV$  (GEMM-II).
  - 10: Write  $O$  to HBM.
-

# Naïve Attention is Bottlenecked by Memory Reads/Writes



**The biggest cost is in moving the bits!**  
Standard implementation requires repeated R/W  
from slow GPU memory

# Background: GPU Compute Model & Memory Hierarchy



[Blogpost](#): Horace He, Making Deep Learning Go Brrrr From First Principles.

Can we exploit the memory asymmetry to get speed up?  
With IO-awareness (accounting for R/W to different levels of memory)

# How to Reduce HBM Reads/Writes: Tiling and Recomputation

## Challenges:

- (1) Don't materialize scores matrix to HBM, and compute softmax normalization without access to full input.
- (2) Backward without the large attention matrix from forward.

## Approaches:

- (1) Kernel fusion with tiling and online softmax: for a given tile of Q, load KV block by block and accumulate values of corresponding tile of O.
- (2) Recomputation: Don't store attn. matrix from forward, recompute it in the backward.

# FlashAttention-2 Algorithm

---

**Algorithm 2** FlashAttention-2 (FMHA)

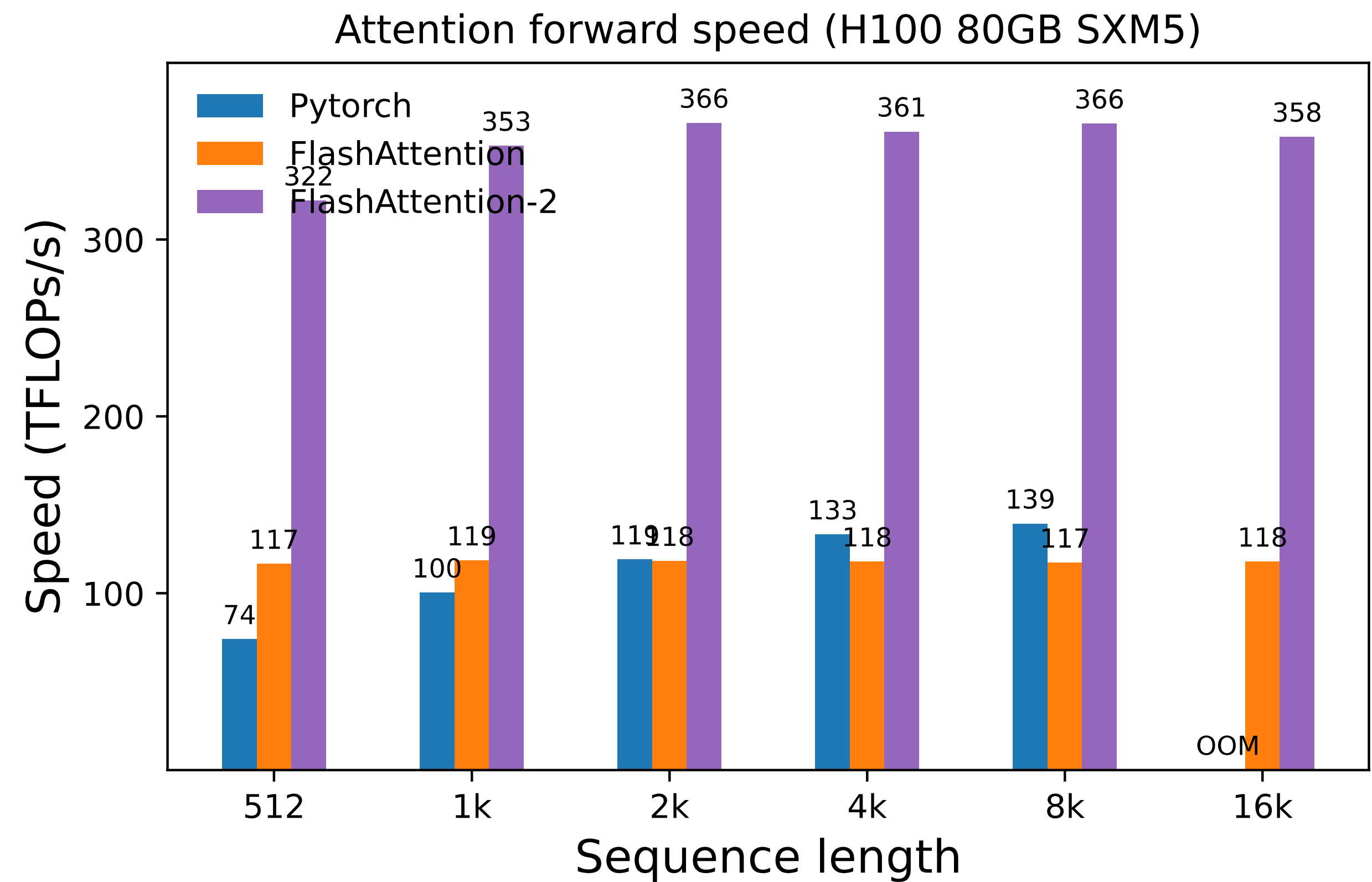
---

```
1: for  $i$  in range(tiles of  $Q$ ) do
2:   Load  $bM \times d$  tile  $Q_i$  from HBM to SMEM.
3:   Initialize  $bM \times d$  accumulator  $O_i = (0)$ .
4:   Initialize  $bM \times 2$  rowmax  $m_i = (-\infty)$  and  $bM \times 1$  rowsum  $\Sigma_i = (0)$ .
5:   for  $j$  in range(tiles of  $K$ ) do
6:     Load  $bN \times d$  tile  $K_j$  from HBM to SMEM.
7:     Compute  $S_{ij} = (1/\sqrt{d})(Q_i K_j^T)$  (SS-GEMM-I).
8:     Update rowmax  $m_i = (m_i^{\text{new}}, m_i^{\text{old}})$ , tracking rowmax at steps  $j$  and  $j - 1$ .
9:     Compute  $\tilde{P}_{ij} = \exp(S_{ij} - m_i^{\text{new}})$ .
10:    Update rowsum  $\Sigma_i = \exp(m_i^{\text{old}} - m_i^{\text{new}})\Sigma_i + \text{rowsum}(\tilde{P}_{ij})$ .
11:    Load  $bN \times d$  tile  $V_j$  from HBM to SMEM.
12:    Compute  $O_i = \exp(m_i^{\text{old}} - m_i^{\text{new}})O_i + \tilde{P}_{ij}V_j$  (RS-GEMM-II).
13:   end for
14:   Compute  $O_i = (1/\Sigma_i)O_i$ .
15:   Write  $O_i$  to HBM.
16: end for
```

---



# Challenge: Optimizing FlashAttention for Modern Hardware - H100



FA2 only gets to 35-40% utilization (no WGMMMA, no TMA)

# FlashAttention-3: Optimizing FlashAttention for H100 GPU

Jay Shah\*, Ganesh Bikshandi\*, Ying Zhang, Vijay Thakkar, Pradeep Ramani, Tri Dao

## 1. New instructions on H100:

- **WGMMMA**: higher throughput MMA primitive, async, collectively executed by a warpgroup (= 4 contiguous warps)
- **TMA**: faster loading from gmem  $\leftrightarrow$  smem, async, saves registers

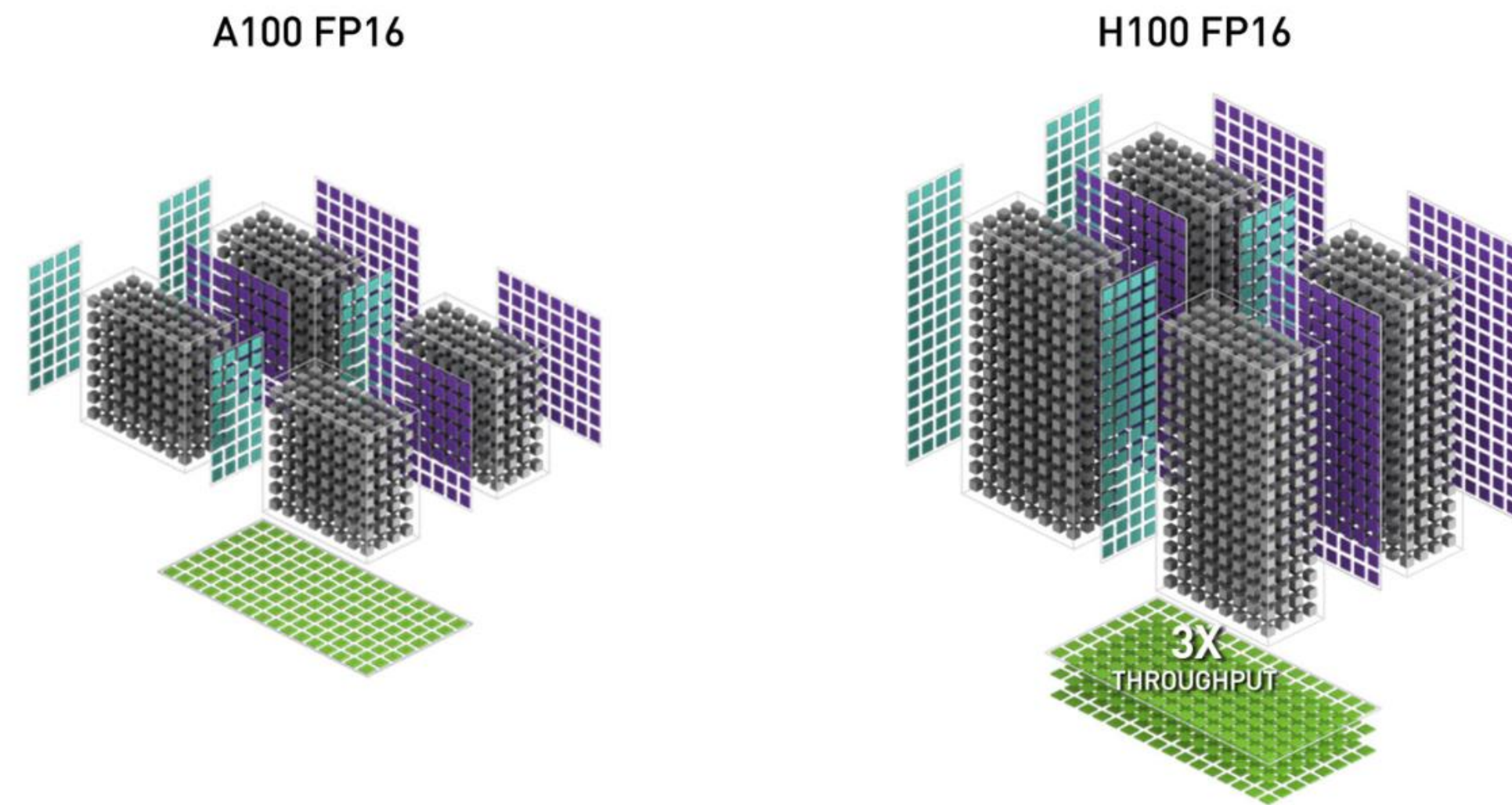
## 2. Asynchrony

- Builds on asynchronous wgmma, TMA, transaction barrier
- Inter-warpgroup overlapping: warp-specialization, pingpong
- Intra-warpgroup overlapping: softmax and async matmul

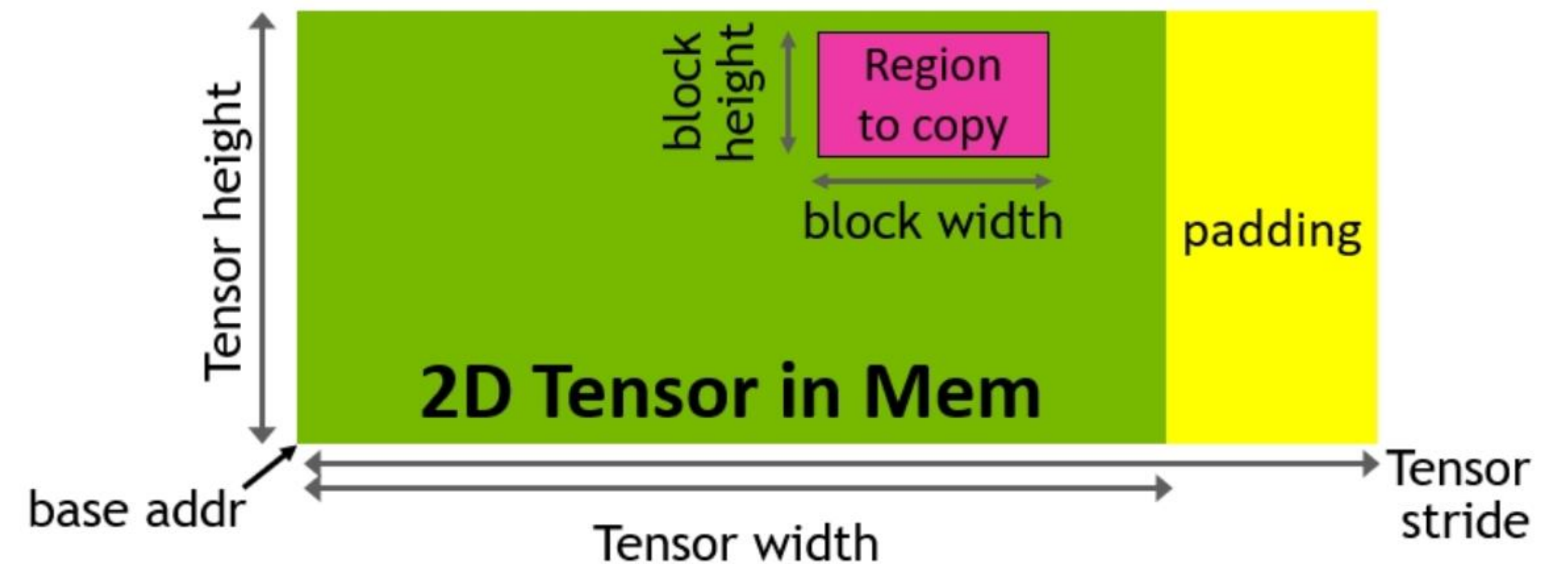
## 3. Low-precision – FP8: layout conformance, in-kernel V transpose

Upshot: **1.6-3x** speedup

# New Instructions: WGMMA & TMA



wgmma necessary, mma.sync can only reach 2/3 peak throughput



TMA: accelerate gmem -> smem, saves registers

WGMMA and TMA integrate into a producer-consumer warp-specialized pipelined design. Use warpgroup-wide register reallocation to give consumers greater share of registers.

# FlashAttention-3

## Algorithm:

### CTA View

---

**Algorithm 1** FLASHATTENTION-3 forward pass **without** intra-consumer overlapping – CTA view

---

**Require:** Matrices  $\mathbf{Q}_i \in \mathbb{R}^{B_r \times d}$  and  $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$  in HBM, key block size  $B_c$  with  $T_c = \lceil \frac{N}{B_c} \rceil$ .

- 1: Initialize pipeline object to manage barrier synchronization with  $s$ -stage circular SMEM buffer.
- 2: **if** in producer warpgroup **then**
- 3:   Deallocate predetermined number of registers.
- 4:   Issue load  $\mathbf{Q}_i$  from HBM to shared memory.
- 5:   Upon completion, commit to notify consumer of the load of  $\mathbf{Q}_i$ .
- 6:   **for**  $0 \leq j < T_c$  **do**
- 7:     Wait for the  $(j \% s)$ th stage of the buffer to be consumed.
- 8:     Issue loads of  $\mathbf{K}_j, \mathbf{V}_j$  from HBM to shared memory at the  $(j \% s)$ th stage of the buffer.
- 9:     Upon completion, commit to notify consumers of the loads of  $\mathbf{K}_j, \mathbf{V}_j$ .
- 10:   **end for**
- 11: **else**
- 12:   Reallocate predetermined number of registers as function of number of consumer warps.
- 13:   On-chip, initialize  $\mathbf{O}_i = (0) \in \mathbb{R}^{B_r \times d}$  and  $\ell_i, m_i = (0), (-\infty) \in \mathbb{R}^{B_r}$ .
- 14:   Wait for  $\mathbf{Q}_i$  to be loaded in shared memory.
- 15:   **for**  $0 \leq j < T_c$  **do**
- 16:     Wait for  $\mathbf{K}_j$  to be loaded in shared memory.
- 17:     Compute  $\mathbf{S}_i^{(j)} = \mathbf{Q}_i \mathbf{K}_j^T$  (SS-GEMM). Commit and wait.
- 18:     Store  $m_i^{\text{old}} = m_i$  and compute  $m_i = \max(m_i^{\text{old}}, \text{rowmax}(\mathbf{S}_i^{(j)}))$ .
- 19:     Compute  $\tilde{\mathbf{P}}_i^{(j)} = \exp(\mathbf{S}_i^{(j)} - m_i)$  and  $\ell_i = \exp(m_i^{\text{old}} - m_i) \ell_i + \text{rowsum}(\tilde{\mathbf{P}}_i^{(j)})$ .
- 20:     Wait for  $\mathbf{V}_j$  to be loaded in shared memory.
- 21:     Compute  $\mathbf{O}_i = \text{diag}(\exp(m_i^{\text{old}} - m_i)) \mathbf{O}_i + \tilde{\mathbf{P}}_i^{(j)} \mathbf{V}_j$  (RS-GEMM). Commit and wait.
- 22:     Release the  $(j \% s)$ th stage of the buffer for the producer.
- 23:   **end for**
- 24:   Compute  $\mathbf{O}_i = \text{diag}(\ell_i)^{-1} \mathbf{O}_i$  and  $L_i = m_i + \log(\ell_i)$ .
- 25:   Write  $\mathbf{O}_i$  and  $L_i$  to HBM as the  $i$ th block of  $\mathbf{O}$  and  $L$ .
- 26: **end if**

---



# Asynchrony: Overlapping GEMM and Softmax

Why overlapping?

**Example:** headdim 128, block size 128 x 192

FP16 WGMMMA:  $2 \times 2 \times 128 \times 192 \times 128 = 12.6$  MFLOPS, 4096 FLOPS/cycle -> **3072 cycles**

MUFU.EX2:  $128 \times 192 = 24.6\text{k}$  OPS, 16 OPS/cycle -> **1536 cycles**

MUFU.EX2 takes 50% the cycles of WGMMMA!

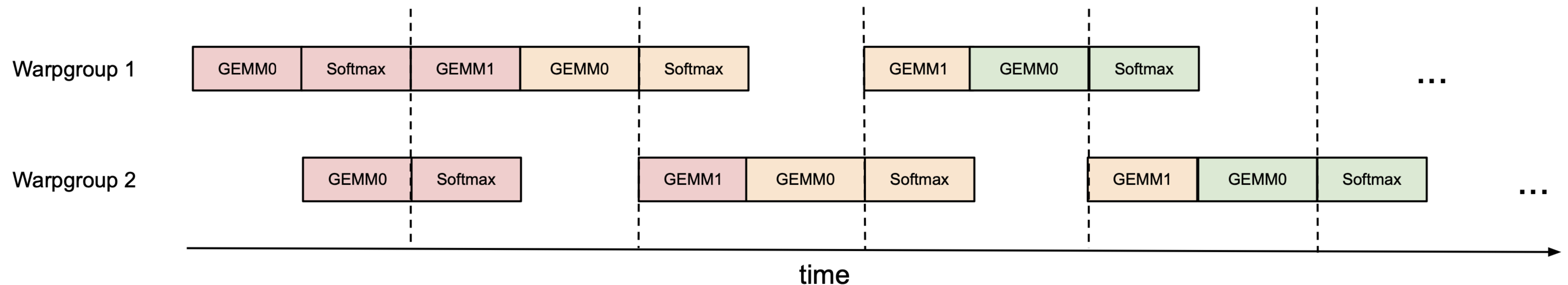
FP8 is even worse: WGMMMA and EX2 both take 1536 cycles.

We want to be doing EX2 while tensor cores are busy with WGMMMA.

# Inter-warpgroup Overlapping of GEMM and Softmax

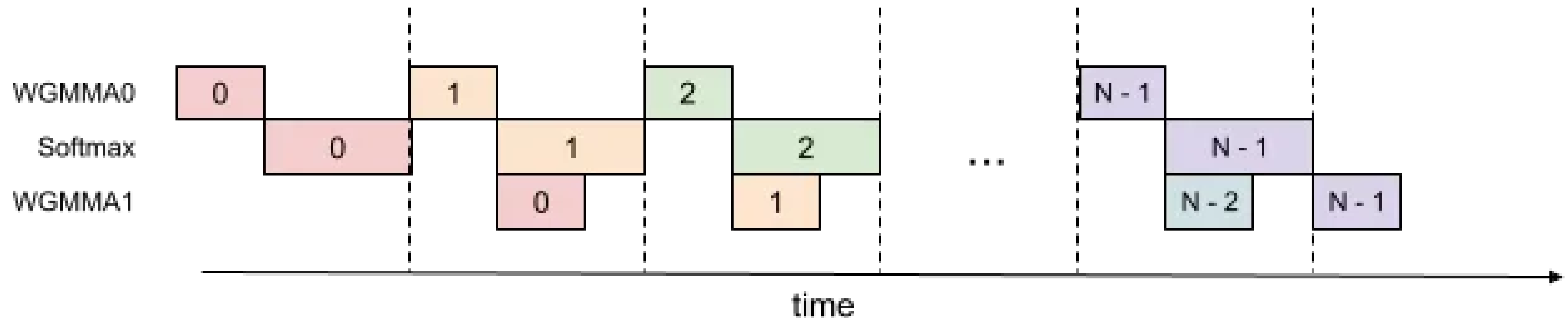
Easy solution: leave it to the warp schedulers!

This works reasonably well, but we can do better.



Pingpong scheduling using synchronization barriers (with `bar.sync`):  
580 TFLOPS -> 640 TFLOPS

# Intra-warpgroup Overlapping of GEMM and Softmax



2-stage pipelining: 640 TFLOPS -> 670 TFLOPS (but higher register pressure)

**Note:** Corrected image from talk.

# FlashAttention-3 Algorithm: Consumer View

---

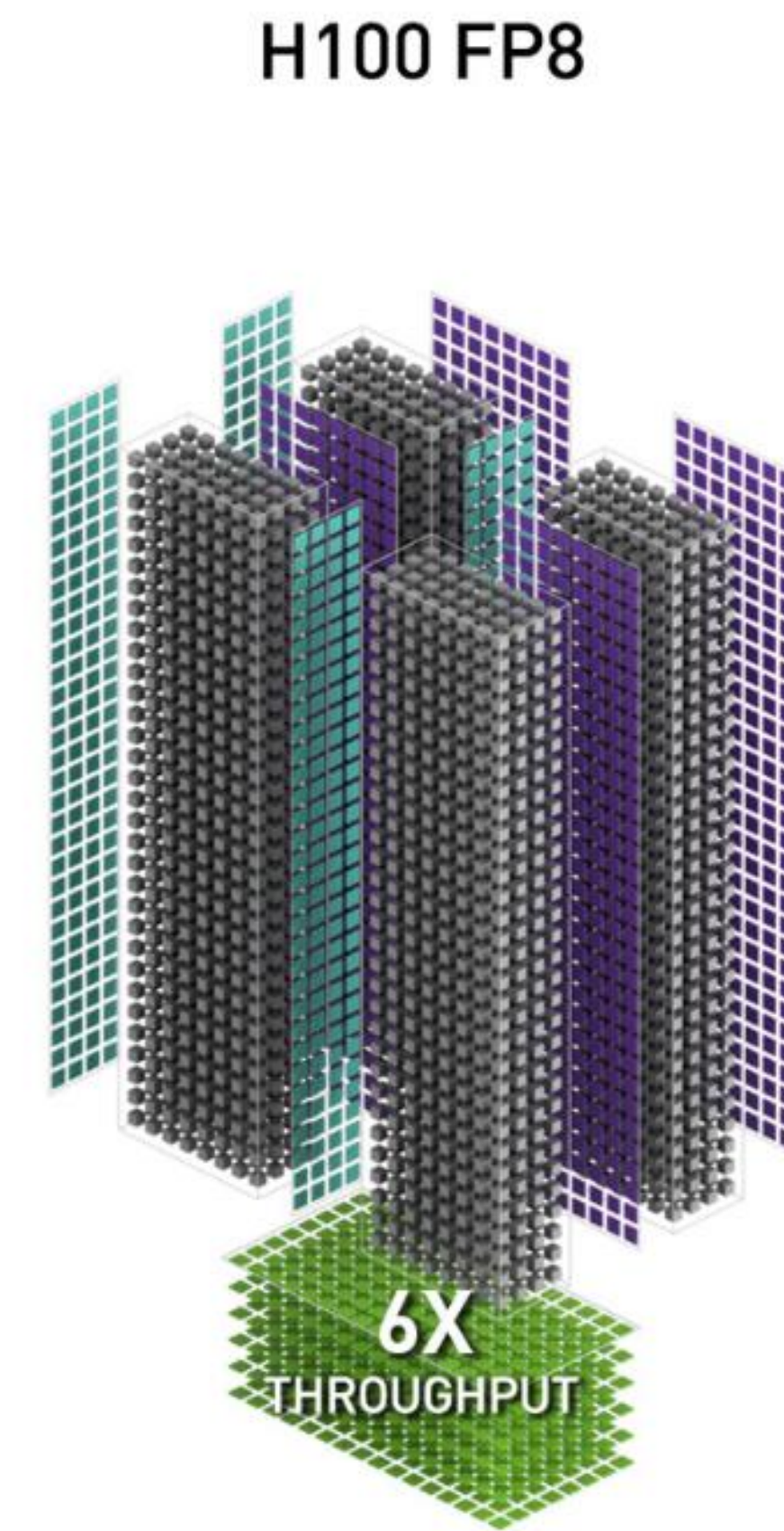
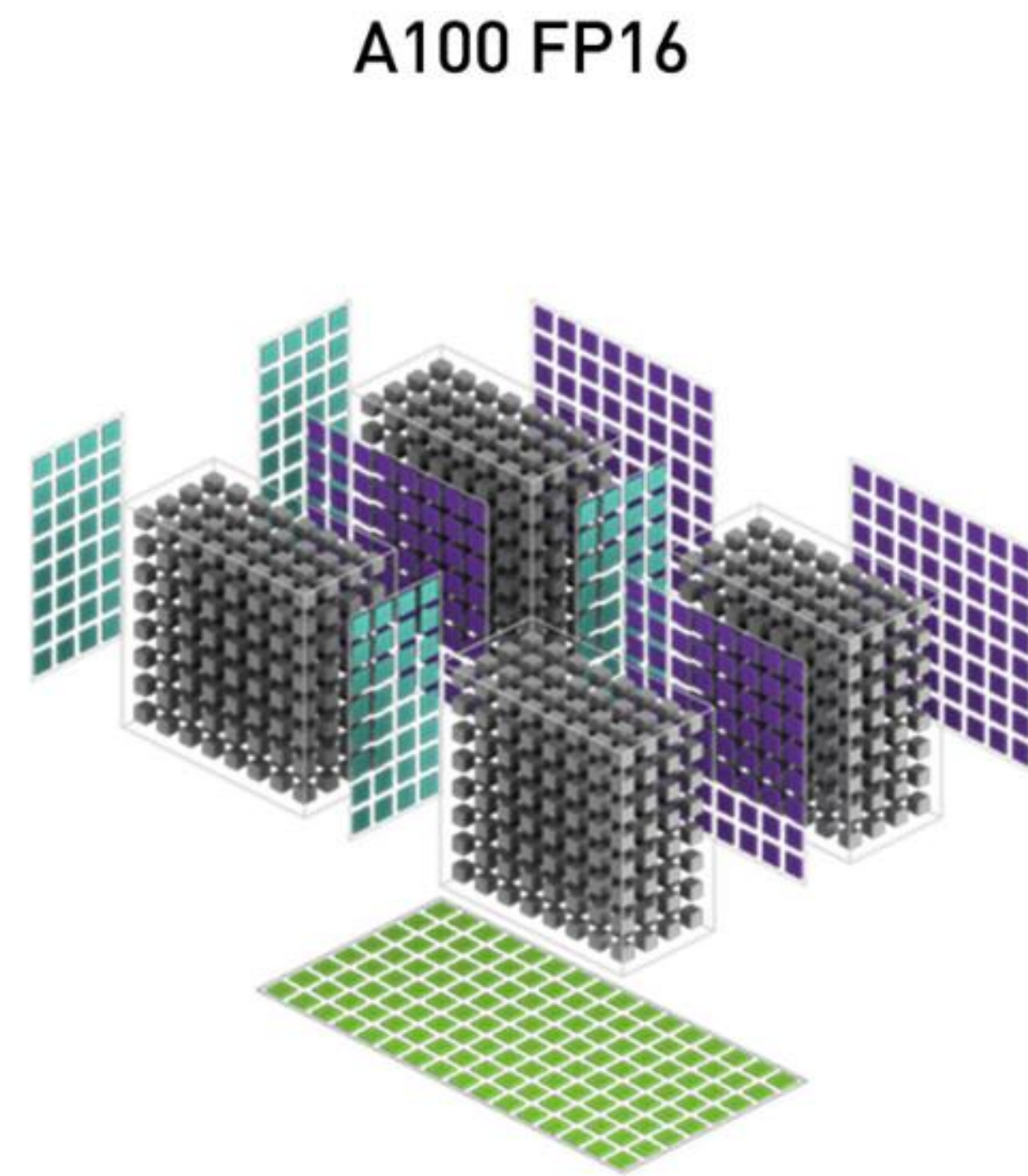
## Algorithm 2 FLASHATTENTION-3 consumer warpgroup forward pass

---

- Require:** Matrices  $\mathbf{Q}_i \in \mathbb{R}^{B_r \times d}$  and  $\mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$  in HBM, key block size  $B_c$  with  $T_c = \lceil \frac{N}{B_c} \rceil$ .
- 1: Reallocate predetermined number of registers as function of number of consumer warps.
  - 2: On-chip, initialize  $\mathbf{O}_i = (0) \in \mathbb{R}^{B_r \times d}$  and  $\ell_i, m_i = (0), (-\infty) \in \mathbb{R}^{B_r}$ .
  - 3: Wait for  $\mathbf{Q}_i$  and  $\mathbf{K}_0$  to be loaded in shared memory.
  - 4: Compute  $\mathbf{S}_{\text{cur}} = \mathbf{Q}_i \mathbf{K}_0^T$  using WGMMMA. Commit and wait.
  - 5: Release the 0th stage of the buffer for  $\mathbf{K}$ .
  - 6: Compute  $m_i, \tilde{\mathbf{P}}_{\text{cur}}$  and  $\ell_i$  based on  $\mathbf{S}_{\text{cur}}$ , and rescale  $\mathbf{O}_i$ .
  - 7: **for**  $1 \leq j < T_c - 1$  **do**
  - 8:     Wait for  $\mathbf{K}_j$  to be loaded in shared memory.
  - 9:     Compute  $\mathbf{S}_{\text{next}} = \mathbf{Q}_i \mathbf{K}_j^T$  using WGMMMA. Commit but do not wait.
  - 10:    Wait for  $\mathbf{V}_{j-1}$  to be loaded in shared memory.
  - 11:    Compute  $\mathbf{O}_i = \mathbf{O}_i + \tilde{\mathbf{P}}_{\text{cur}} \mathbf{V}_{j-1}$  using WGMMMA. Commit but do not wait.
  - 12:    Wait for the WGMMMA  $\mathbf{Q}_i \mathbf{K}_j^T$ .
  - 13:    Compute  $m_i, \tilde{\mathbf{P}}_{\text{next}}$  and  $\ell_i$  based on  $\mathbf{S}_{\text{next}}$ .
  - 14:    Wait for the WGMMMA  $\tilde{\mathbf{P}}_{\text{cur}} \mathbf{V}_{j-1}$  and then rescale  $\mathbf{O}_i$
  - 15:    Release the  $(j \% s)$ th, resp.  $(j - 1 \% s)$ th stage of the buffer for  $\mathbf{K}$ , resp.  $\mathbf{V}$ .
  - 16:    Copy  $\mathbf{S}_{\text{next}}$  to  $\mathbf{S}_{\text{cur}}$ .
  - 17: **end for**
  - 18: Wait for  $\mathbf{V}_{T_c-1}$  to be loaded in shared memory.
  - 19: Compute  $\mathbf{O}_i = \mathbf{O}_i + \tilde{\mathbf{P}}_{\text{last}} \mathbf{V}_{T_c-1}$  using WGMMMA. Commit and wait.
  - 20: Rescale  $\mathbf{O}_i$  based on  $m_i$ . Compute  $L_i$  based on  $m_i$  and  $\ell_i$ .
  - 21: Epilogue: Write  $\mathbf{O}_i$  and  $L_i$  to HBM as the  $i$ -th block of  $\mathbf{O}$  and  $L$ .
-



# Low-precision: FP8



FP8 doubles WGMMMA throughput, but trade off accuracy

# Layout conformance challenges with FP8

FP8 WGMMMA: requires operand SMEM tensors to be memory-contiguous in the inner dimension (k-major)

It is standard for QKV to be memory-contiguous in the head dimension (BSHD).  
- Note: TMA can't change the contiguous mode.

For gemm0 (Q.K<sup>T</sup>), this is fine as is. For gemm1 (P.V), we need to transpose V.

**Solution:** In-kernel transpose of V in the producer warpgroup.  
Uses LDSM/STSM instructions with custom layouts and byte permute in-between.

## Layout conformance challenges with FP8, ctd.

We also need to reshape layout of scores accumulator for gemm1.

**Why?** FP32 accumulator layout differs from FP8 operand A layout.

T0 {d0, d1}	T1 {d0, d1}	T2 {d0, d1}	T3 {d0, d1}	T0 {d4, d5}	T1 {d4, d5}	T2 {d4, d5}	T3 {d4, d5}
T0 {d2, d3}	T1 {d2, d3}	T2 {d2, d3}	T3 {d2, d3}	T0 {d6, d7}	T1 {d6, d7}	T2 {d6, d7}	T3 {d6, d7}

Figure 3: FP32 accumulator register WGMMMA layout – rows 0 and 8, threads 0-3, entries 0-7.

T0 {a0, a1}	T0 {a2, a3}	T1 {a0, a1}	T1 {a2, a3}	T2 {a0, a1}	T2 {a2, a3}	T3 {a0, a1}	T3 {a2, a3}
T0 {a4, a5}	T0 {a6, a7}	T1 {a4, a5}	T1 {a6, a7}	T2 {a4, a5}	T2 {a6, a7}	T3 {a4, a5}	T3 {a6, a7}

Figure 4: FP8 operand A register WGMMMA layout – rows 0 and 8, threads 0-3, entries 0-7.

**Note:** Can use in-kernel “transpose” to write out row permutation of V such that we can avoid shuffle instructions for the reshape (but not byte permute).

# Persistent Kernels in FlashAttention

**Idea:** Decouple *physical* CTAs from *logical* work tiles, launching fixed number of CTAs. Can then overlap epilogue of current work tile with prologue loads of next work tile.

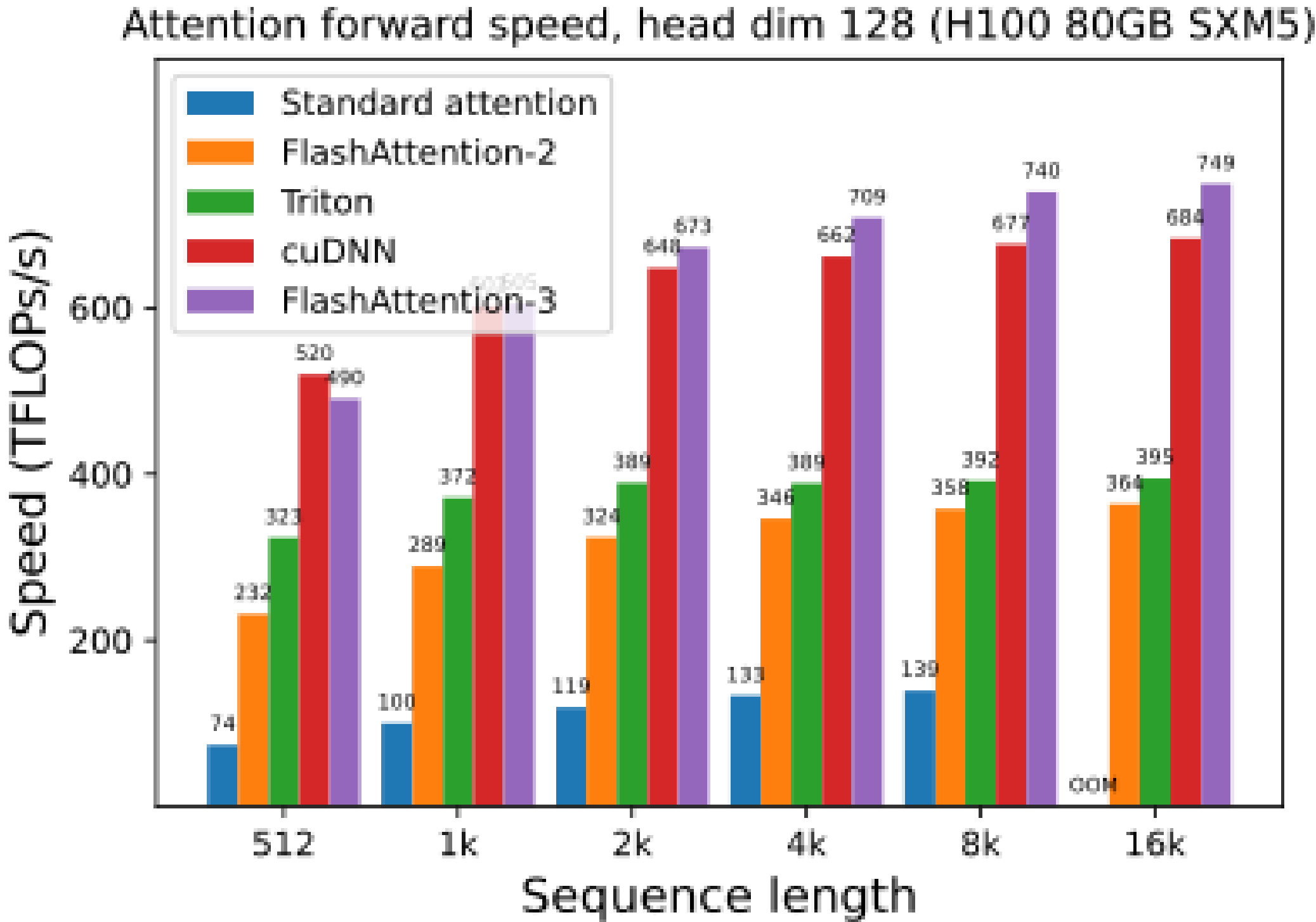
**Example:** SeqLen = 4096, Heads = 8, Batch = 4. Fix BlockM = 128, so mblocks = 32. Have  $32 * 8 * 4 = 1024$  work tiles to process in the kernel.

Without persistent kernel, launch CUDA grid with dims = (32, 8, 4), so 1024 CTAs. 1-to-1 mapping of CTAs with work tiles.

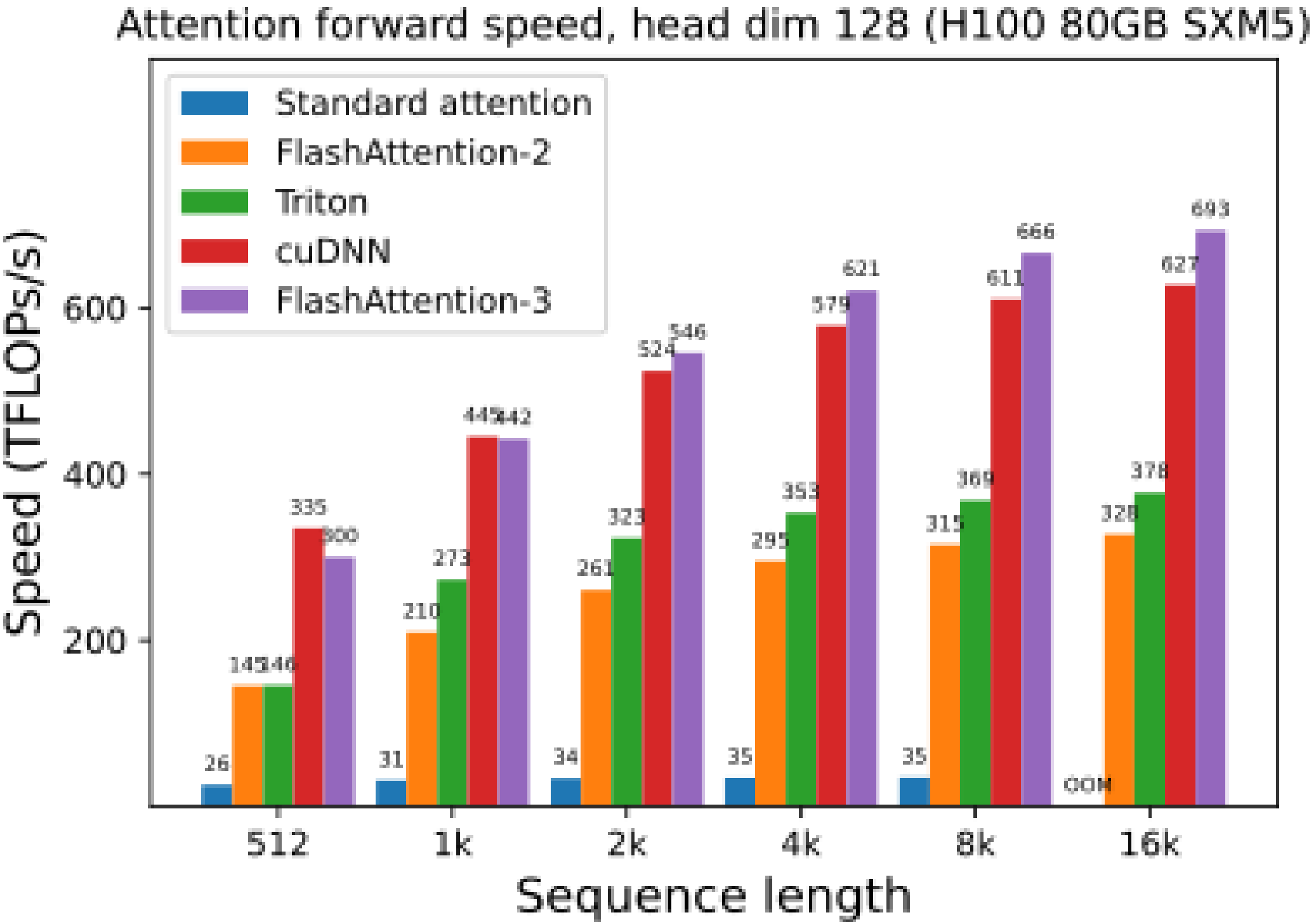
With persistent kernel on H100 SXM5 GPU, launch 132 CTAs = num SMs. Each CTA runs over a fraction of the 1024 work tiles.

Can dynamically allocate work tiles to next available CTA in persistent kernel. Helps with load balancing when doing causal masking.

# BF16 Benchmark: 1.6-2.0x speedup



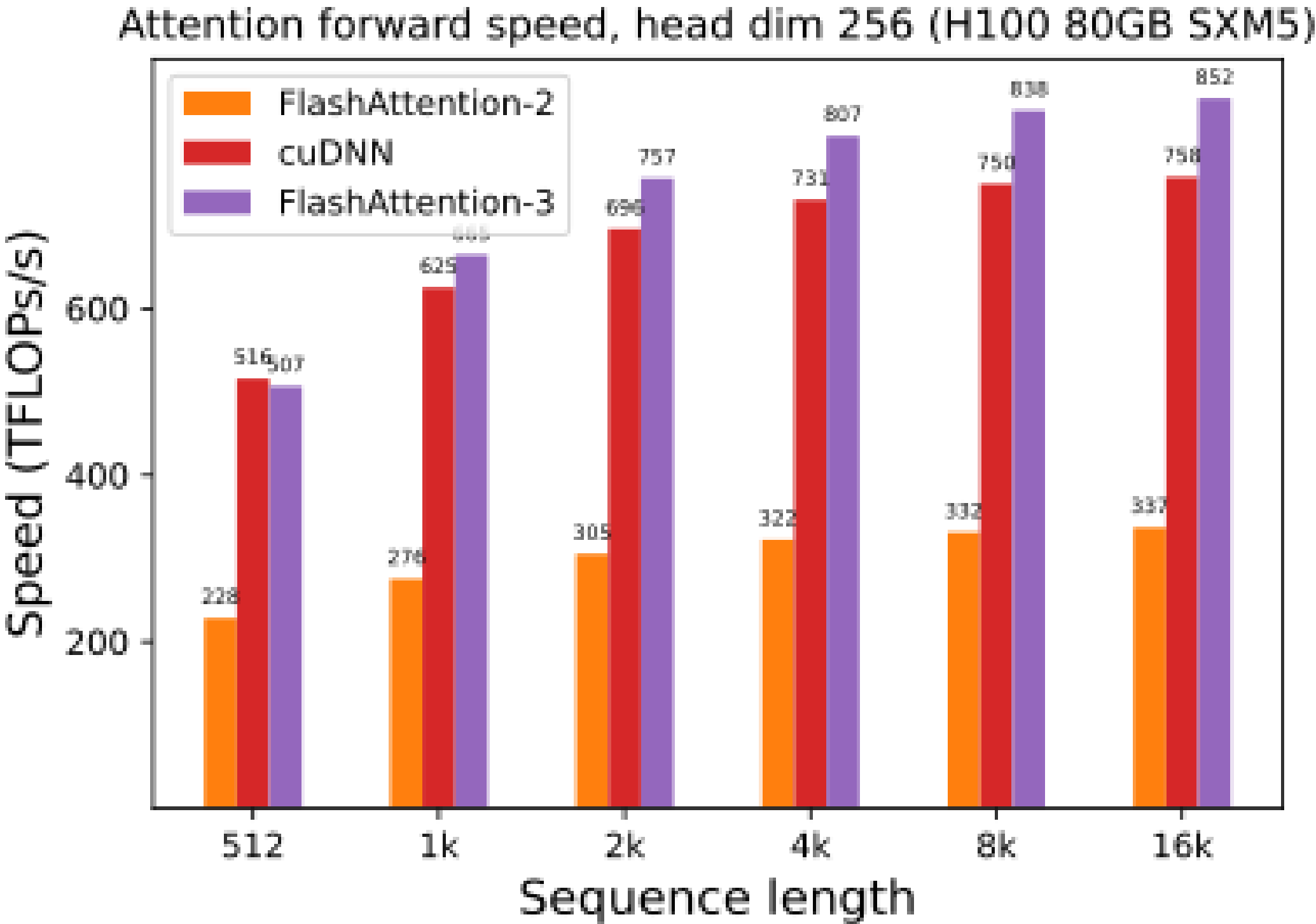
Without causal mask



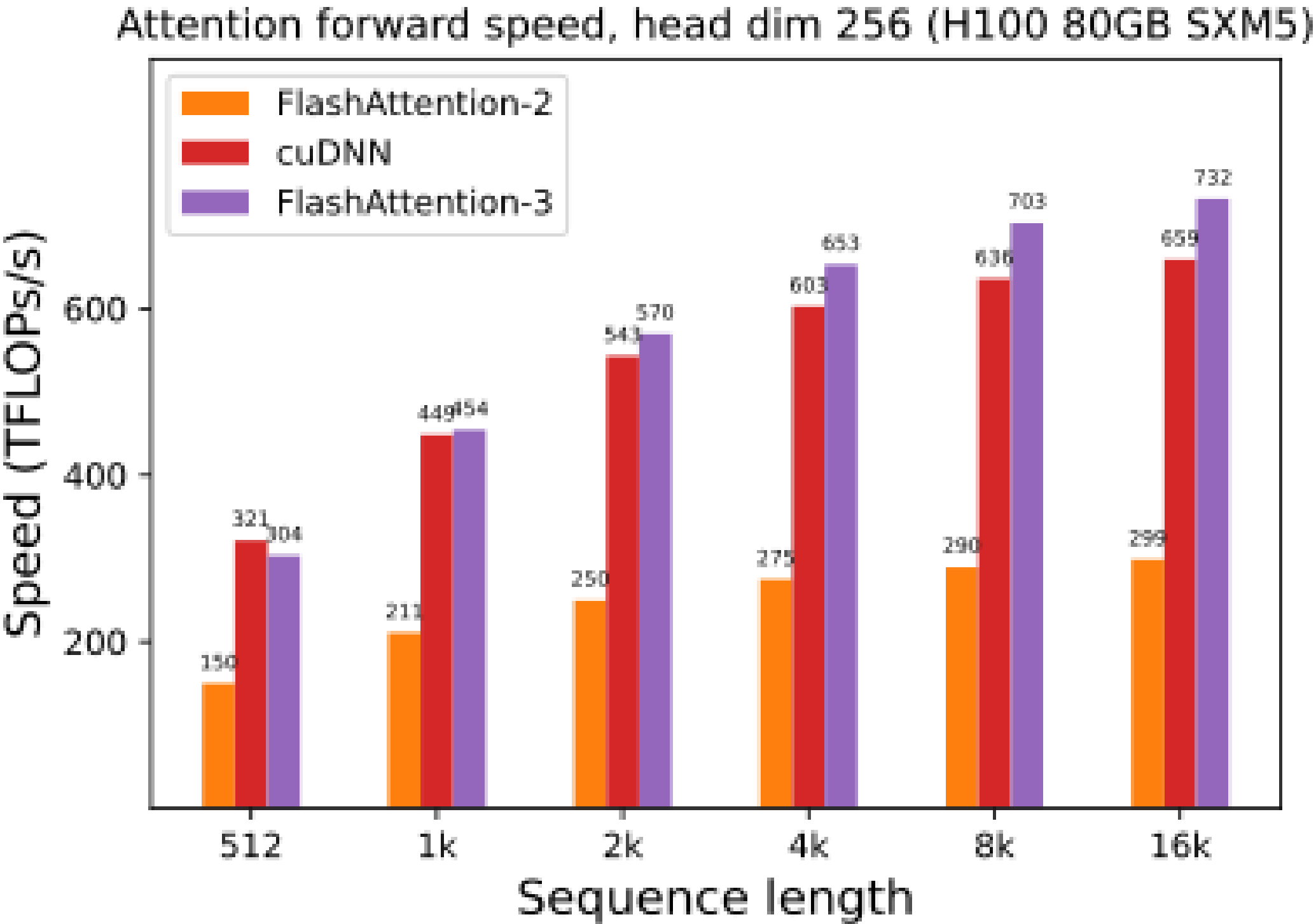
With causal mask



# BF16 Benchmark: reach up to 850 TFLOPS

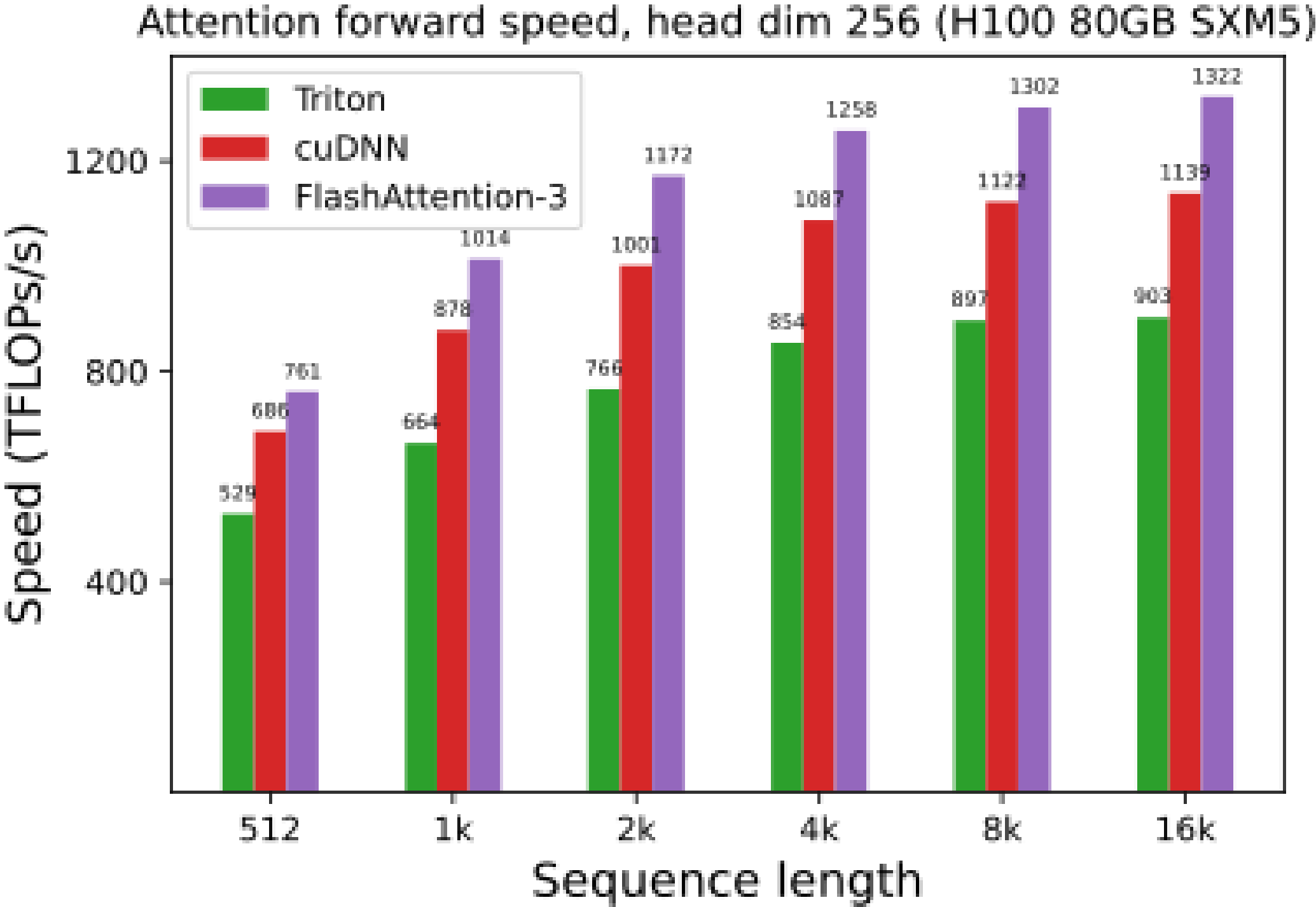


Without causal mask

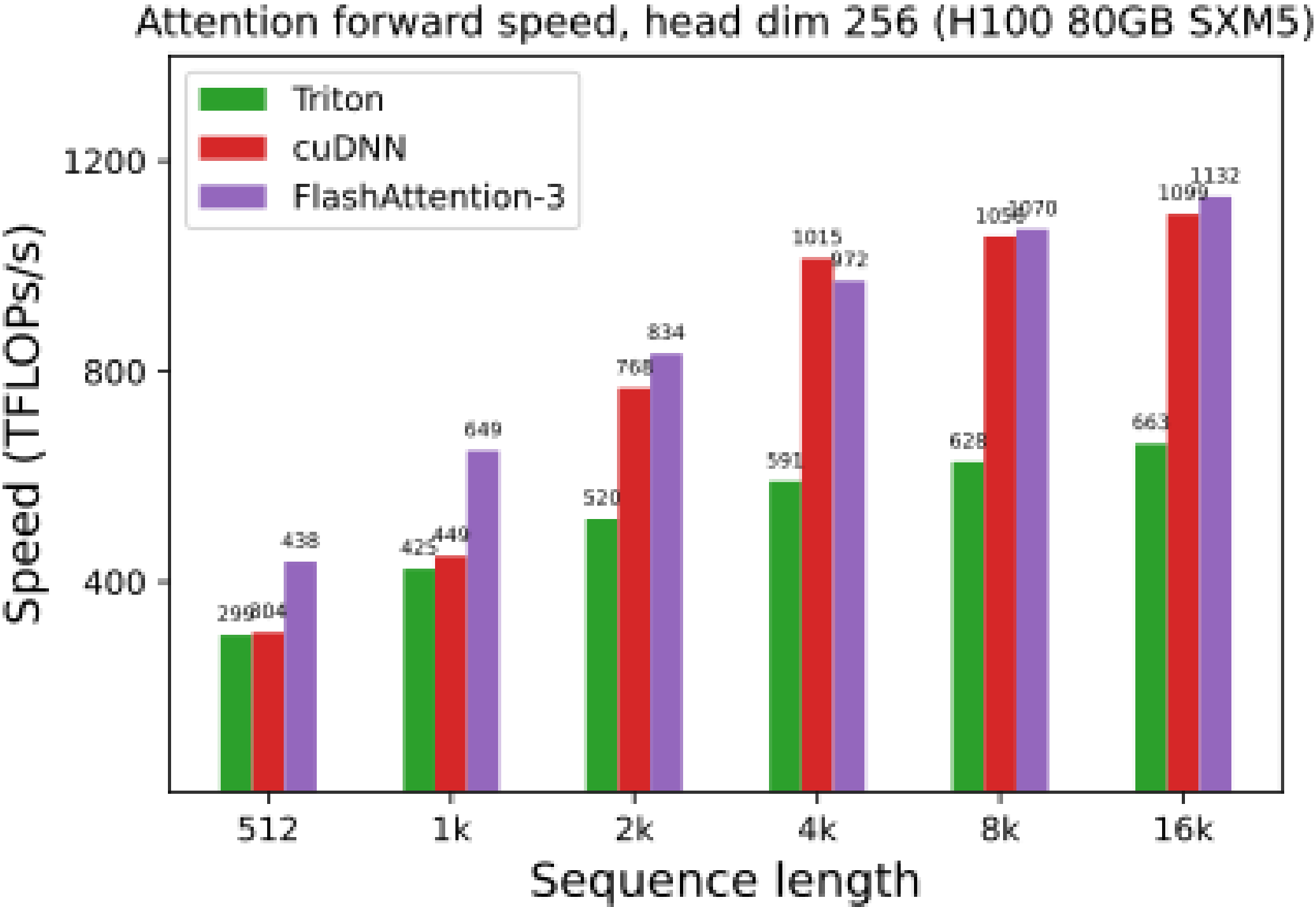


With causal mask

# FP8 Benchmark: up to 1.3 PFLOPS



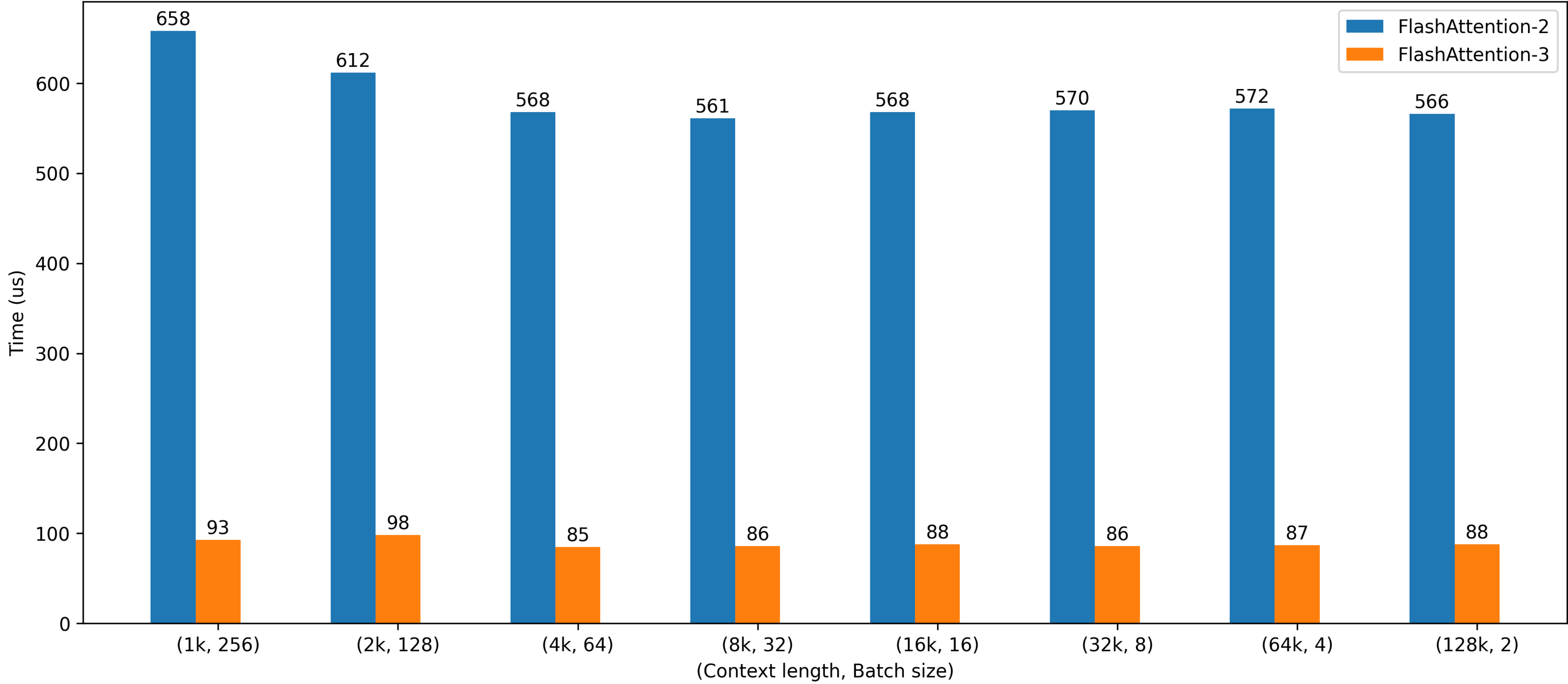
Without causal mask



With causal mask

# BF16 Decode Benchmark for MQA. Lower is better!

BF16 Attention, head dim 128 (H100 80GB PCIe).  
MQA 16, query sequence length = 4.





# FlashAttention-3 for Decoding Inference

For decoding, query length is short (on the order of a few tokens), while context length is long (for example, 128k).

## Optimizations:

- 1) Split along the KV sequence length to occupy the GPU with enough work  
- already in FlashAttention-2 as **Flash-Decoding**. Reuse the same algorithm.
- 2) **GQA packing**: Pack multiple query heads into a single query tile.

# Summary – FlashAttention-3

**Fast** and **accurate** attention optimized for modern hardware

Key algorithmic ideas: **asynchrony**, **low-precision**

- for inference: **Split KV** (Flash-Decoding) and **GQA packing**.

Upshot: **faster** training, **better** models with **longer** sequences

Code: <https://github.com/Dao-AI-Lab/flash-attention>

# Building FlashAttention-3 with CUTLASS

Overall structure has three classes:

1. CollectiveMainloop for load and mma.

2. CollectiveEpilogue for store.

3. TileScheduler to manage work loop for persistent kernel.

Each class has its own set of kernel parameters  
(to\_underlying\_arguments).

```
1 CollectiveMainloop collective_mainloop
2 CollectiveEpilogue collective_epilogue
3 if producer:
4     reg_dealloc(LoadRegisterRequirement)
5     work_tile = get_initial_work()
6     while(work_tile.is_valid()):
7         auto block_coord = work_tile_info.get_block_coord()
8         collective_mainloop.load(block_coord)
9         work_tile = get_next_work()
10 else:
11     reg_alloc(MmaRegisterRequirement)
12     work_tile = get_initial_work()
13     while(work_tile.is_valid()):
14         auto block_coord = work_tile_info.get_block_coord()
15         Tensor t0r0 = collective_mainloop.mma(block_coord)
16         collective_epilogue.store(t0r0, block_coord)
17         work_tile = get_next_work()
```