Optimizing Attention for Modern Hardware



Tri Dao https://tridao.me

Motivation: Modeling Long Sequences

Enable New Capabilities

NLP: Large context required to understand books, plays, codebases.

Computer vision: higher resolution can lead to better, more robust insight.





Close Reality Gap

Open New Areas

Time series, audio, video, medical imaging data naturally modeled as sequences of millions of steps.





Efficiency is the Bottleneck for Modeling Long Sequences with Attention

Context length: how many other elements in the sequence does the current element interact with. Increasing context length slows down (or stops) training



How to efficiently scale models to longer sequences?

Background: Attention is the Heart of Transformers



Transformer

Encoder

Background: Attention Mechanism



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

Attention scales quadratically in sequence length N



Value

Output



Background: Approximate Attention



Survey: Tay et al. Long Range Arena : A Benchmark for Efficient Transformers. ICLR.2020

Is there a fast, memory-efficient, and exact attention algorithm?

Approximate attention: tradeoff quality for speed fewer FLOPs

Our Observation: Attention is Bottlenecked by Memory Reads/Writes



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

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





Attention prob = row-wise normalized similarity score



Output

Background: GPU Compute Model & Memory Hierarchy



<u>Blogpost</u>: 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: Compute by Blocks

Challenges:

(1) Compute softmax normalization without access to full input.

(2) Backward without the large attention matrix from forward.

Approaches:

(1) Tiling: Restructure algorithm to load block by block from HBM to SRAM to compute attention.

(2) Recomputation: Don't store attn. matrix from forward, recompute it in the backward.

Attention Computation Overview



normalization constant

Compute by blocks: easy to split Q, but how do we split K & V? 10



Tiling – 1st Attempt: Computing Attention by Blocks



11

Tiling – 2nd Attempt: Computing Attention by Blocks, with Softmax Rescaling

Tiling

Decomposing large softmax into smaller ones by scaling.

Animation credit: Francisco Massa

1. Load inputs by blocks from HBM to SRAM.

2. On chip, compute attention output with respect to that block.

3. Update output in HBM by scaling.

Recomputation (Backward Pass)

By storing softmax normalization from forward (size N), quickly recompute attention in the backward from inputs in SRAM.

Attention	Standard	FlashAttentic
GFLOPs	66.6	75.2 <mark>(↑13%</mark>
HBM reads/writes (GB)	40.3	4.4 (↓9x)
Runtime (ms)	41.7	7.3 (↓6x)

FlashAttention speeds up backward pass even with increased FLOPs.

FlashAttention: 2-4x speedup, 10-20x memory reduction

2-4x speedup — with no approximation!

10-20x memory reduction — memory linear in sequence length

Challenge: Optimizing FlashAttention for Modern Hardware

FlashAttention-2 is highly optimized on A100, reaching 70% utilization

But, FA-2 only gets to 35-40% utilization on H100!

FlashAttention-3: Optimizing FlashAttention for Hopper Architecture

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

- **1. New instructions** on H100:
 - WGMMA: higher throughput MMA primitive, async - **TMA**: faster loading from gmem <-> smem, async, saves registers

2. Asynchrony

- Inter-warpgroup overlapping: warp-specialization, pingpong scheduling - Intra-warpgroup overlapping: softmax and async matmul

3. Low-precision – FP8

- Solve for layout conformance, in-kernel V transpose

- **Plus:** Persistent kernels, LPT scheduling for causal attention, GQA packing, MLA
 - **Upshot**: 1.6-3x speedup on **Hopper**, algorithmic ideas apply for **Blackwell**

New Instructions on Hopper: WGMMA & TMA

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

TMA: accelerate gmem -> smem copy, saves registers

Both WGMMA and TMA are **asynchronous** instructions: threads issue them and can then do other work while they execute.

Asynchrony: Overlapping GEMM and Softmax

Why overlap?

Special Function Units (SFU) have very low throughput relative to Tensor Cores.

Example: headdim 128, block size 128 x 128 FP16 WGMMA: 2 x 2 x 128 x 128 x 128 = 8.4 MFLOPS, 4096 FLOPS/cycle -> **2048 cycles** MUFU.EX2: 128 x 128 = 16k OPS, 16 OPS/cycle -> **1024 cycles**

MUFU.EX2 takes 50% the cycles of WGMMA! FP8, or Blackwell is even worse: WGMMA and EX2 both take 1024 cycles. We want to be doing EX2 while tensor cores are busy with WGMMA.

19

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

time

20

Intra-warpgroup Overlapping of GEMM and Softmax

Per warpgroup, can finally exploit **asynchrony** of WGMMA. Overlap GEMM1 for kth iteration with softmax for (k+1)th iteration. \bullet

- GEMM1 are now live concurrently.

2-stage intra-warpgroup overlapping: 640 TFLOPS -> 670 TFLOPS

Uses more registers since accumulator for next GEMM0 and operand for current

time

2	1
Ζ	Т

Low-precision: FP8

A100 FP16

H100 FP8

FP8 Tensor Cores double WGMMA throughput, but trade off accuracy

FP8 Attention with Incoherent Processing

the outliers.

Reduces quantization error by 2.6x on normally distributed QKV data with 0.1% entries given large magnitude (to simulate outliers).

	M	ethod	Baseline	e FP16	FLASHATTENTION	-2 FP16	FLASH	ATTENTION-3 FP16	
	R	MSE	3.20	e-4	1.9e-4			1.9e-4	
Meth	od	Baseli	ine FP8	FLASH	Attention-3 FP8	No block	k quant	No incoherent proc	essing
RMS	E	2.	4e-2		9.1e-3	9.3	e-3	2.4e-2	

- Can multiply Q and K with a random orthogonal matrix (Hadamard) to "spread out"
- Note: $S = Q.K^T = (QJ)(KJ)^T$ for orthogonal matrix J since J.J^T = I by definition.

Persistent Kernels: Hiding Prologue & Epilogue

Motivation: Tensor Cores are so fast, Prologue/Epilogue latency become non-trivial.

Idea: Fixed number of CTAs (= num SMs), persistent across work tiles. • Asynchronous TMA store to overlap epilogue, prologue, and mainloop.

Persistent Kernels: Hiding Prologue & Epilogue

Persistent kernels: 670 TFLOPS -> 700 TFLOPS

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

26

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

1		
2		
3		

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

T = 2

1		
2	2	
3	3	
1		

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

1			
2	2		
3	3	3	
1	1		

2		

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

1			
2	2		
3	3	3	
1	1	1	

2		
2		
3		

30

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

1			
2	2		
3	3	3	
1	1	1	1

2		
2	2	
3	3	

2	1
J	Т

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

T = 6

1			
2	2		
3	3	3	
1	1	1	1

2			
2	2		
3	3	3	
1			

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

T = 7

1			
2	2		
3	3	3	
1	1	1	1

2			
2	2		
3	3	3	
1	1		

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

1			
2	2		
3	3	3	
1	1	1	1

2			
2	2		
3	3	3	
1	1	1	

Challenge: Unequal work per tile

Example: 2 batches, 3 workers (SMs)

1			
2	2		
3	3	3	
1	1	1	1

T = 9

Work distribution: [9, 5, 6] blocks. Longest tile is always scheduled last!

2			
2	2		
3	3	3	
1	1	1	1

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

SIAM J. APPL. MATH. Vol. 17, No. 2, March 1969

BOUNDS ON MULTIPROCESSING TIMING ANOMALIES*

1. Introduction. It is well known (cf. [5], [6], [8]) to workers in the field of parallel computation that a multiprocessing system consisting of many identical processors acting in parallel may exhibit certain somewhat unexpected "anomalies," even though the system operates under a rather natural set of rules; e.g., it can happen that increasing the number of processors can *increase* the length of time required to execute a given set of tasks. In this paper we study a typical model of such a multiprocessing system, and we determine the precise extent by which the execution time for a set of tasks can be influenced because of these timing anomalies. A special case of this model will be shown to generate an interesting numbertheoretic question, partial answers to which are given in the latter half of the paper.

R. L. GRAHAM[†]

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

3 1

T = 1

2		

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

2	2	

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

40

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

41

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

2			
3	3		
2	2	2	2

42

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

2	2		
3	3	3	
2	2	2	2

43

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

T = 7

Work distribution: [7, 7, 6] vs [9, 5, 6] previously

2			
2	2		
3	3	3	
2	2	2	2

Challenge: Unequal work per tile

Idea: Assign work to workers using longest-processing-time-first

Caveat: Take care to not spill L2 cache (using L2 swizzling)

Causal attention speed 670 TFLOPS -> 730 TFLOPS

45

BF16 Benchmark: 1.8-2.2x speedup

Without causal mask

CUDA tool kit 12.8 Triton 3.1 cuDNN 9.6

Attention forward speed, causal, head dim 128 (H100 80GB SXM5)

Causal Attention 730-750 TFLOPS \approx Matmul speed!

46

BF16 Benchmark: Reach up to 840 TFLOPS!

Attention forward speed, head dim 256 (H100 80GB SXM5)

Without causal mask

CUDA tool kit 12.8 Triton 3.1 cuDNN 9.6

Attention forward speed, causal, head dim 256 (H100 80GB SXM5)

With causal mask

47

FP8 Benchmark: Up to 1.3 PFLOPS!

Without causal mask

CUDA tool kit 12.8 Triton 3.1 cuDNN 9.6

Attention forward speed, causal, head dim 256 (H100 80GB SXM5)

48

Optimizations for Decoding Inference: Old and New

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

From FA-2, have **Flash Decoding**: split along the KV sequence length to occupy the GPU with enough work.

GQA Packing: compute for multiple query heads per KV head

WGMMA tile is 64 wide in the M dimension. This is wasted for short query length! However, we can pack multiple query heads to fill out WGMMA tile for MQA/GQA.

Unpacked 64x128 Query Tile for a single head (4 query tokens)

FA-2 already did this for the case of **one** query token, which is a simple reshape.In FA-3, we extend to the more involved case of arbitrary query length.Also benefits some compute-bound situations with tile quantization effects

Packed 64x128 Query Tile (16 query heads, 4 query tokens)

MQA 16, query sequence length = 4.

Multi-head Latent Attention (MLA): Warp specialization for large head dim

52

WG1 does both QK matmul and PV matmul WG2 only does PV matmul

Multi-head Latent Attention (MLA): Warp specialization for large head dim

Even seqlen_q = 1 (decoding 1 token) already hits compute-bound regime!

Blackwell's new SASS instructions: Reducing issuing pressure

Challenge: each SM can issue 4 instructions per cycle, but FMA throughput is 128 ops/cycle, so every instruction needs to be FMA to achieve peak throughput.

Approach: New instructions:

- FADD2, FMUL2, FFMA2: 2 ops per instruction Accessible through PTX (add.f32x2, mul.f32x2, fma.f32x2) and through intrinsics (___fadd2_rn, ___fmul2_rn, ___ffma2_rn)

- FMNMX3: 3-input floating-point maximum, reduce no. of instructions Not directly accessible but compiler will generate if you use fmax(a, fmax(b, c))

CUTLASS team. Example 77. <u>https://github.com/NVIDIA/cutlass/tree/main/examples/77_blackwell_fmha</u>

CUTLASS team. Example 77. <u>https://github.com/NVIDIA/cutlass/tree/main/examples/77_blackwell_fmha</u>

Summary – FlashAttention

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

Key algorithmic ideas: asynchrony, low-precision - Persistent kernels with LPT scheduling for causal attention - For inference: **Split KV** (Flash-Decoding) and **GQA packing**

Upshot: faster training, better models with longer sequences

Code: https://github.com/Dao-AILab/flash-attention https://github.com/NVIDIA/cutlass/tree/main/examples/77 blackwell fmha