

A agenda:

- (1) - as we didn't get to it last time: FLOPs analysis in transformer trainings
- (2) Revisit scaling in normal softmax attention
- (3) How is Attention Executed on GPUs today?
- (4) Revisit of Fusion, Tiling, Caching vs. recompilation, hardware-specific optimizations
- (5) Flash Attention: → we didn't get to this today!

→ credits:

- Lec 4 of CS 229s 23': https://docs.google.com/presentation/d/1PV0cKnzcHRCAbJy3OtER1UdJME7T7WBEqLLYmKWR4g/edit#slide=id.g24c4df0ad9e_0_52
- Lec 5 of CS 229s 23': https://docs.google.com/presentation/d/1kY5MOOJfWny7SvV4D7YeSAulOTMnnVJTKKG5DQYLid0/edit#slide=id.g24c4f25c1b8_0_6
- Lec 13 of CMU 15-442: <https://mlsyscourse.org>

(1) FLOPs Analysis in ML Training for Transformers

- reminder of storage vs recompilation and how it is present in backprop

→ consider:

$$L(w_1, w_2, w_3) = 2(w_1 \cdot w_2) \cdot w_3$$

$$\frac{\partial L}{\partial w_3} = w_3 \quad \frac{\partial L}{\partial w_2} = v \quad \frac{\partial v}{\partial w_1} = 2 \quad \frac{\partial v}{\partial w_2} = w_2$$

$\frac{\partial v}{\partial w_1}$... ideally

$$\frac{\partial L}{\partial w_2} = \omega_1$$

locally we reuse this!

so: $\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial v} \right) \left(\frac{\partial v}{\partial u} \right) \left(\frac{\partial u}{\partial w_1} \right)$

2) $\frac{\partial L}{\partial w_2} = \left(\frac{\partial L}{\partial v} \right) \left(\frac{\partial v}{\partial u} \right) \left(\frac{\partial u}{\partial w_2} \right)$

3) $\frac{\partial L}{\partial w_3} = v$

↳ reusing v ($v = 2v$)

required storing v from

fw pass

*tension b/w storing activations

to make backprop more efficient, and

minimizing memory use

* FLOPs analysis for MLP (multi layer perceptron)

$$x: B \times M$$

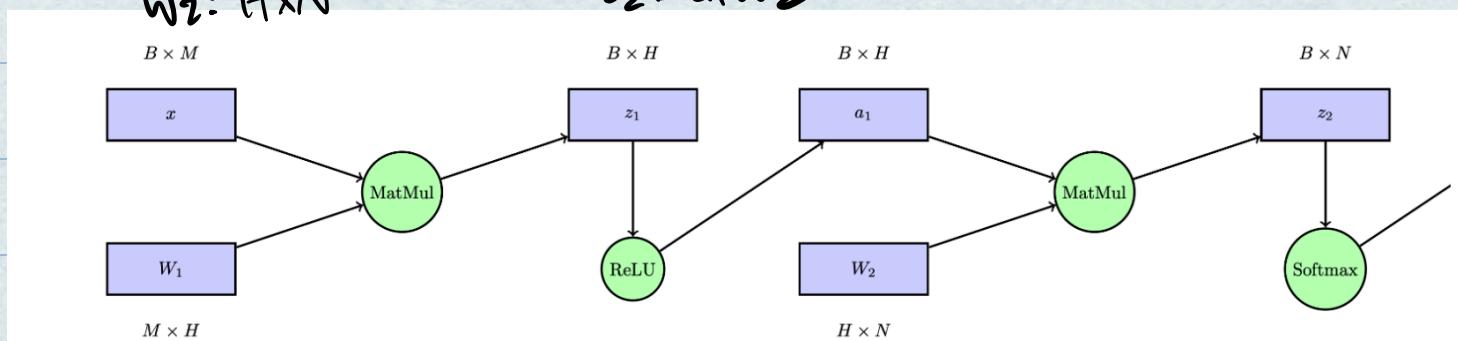
$$z_1 = xW_1$$

$$W_1: M \times H$$

$$a_1 = \text{ReLU}(z_1)$$

$$W_2: H \times N$$

$$z_2 = a_1 W_2$$



2BNH

- Let's consider $a_1 W_2$: 2BNH flop

→ what about backprop?

$$\frac{\partial L}{\partial w_2} = q_1 \cdot \bar{T} * \frac{\partial L}{\partial z_2}$$



$$H \times B \cdot B \times N =$$

$\frac{1}{2} HBN$ flop

$$\frac{\partial L}{\partial q_1} = \frac{\partial L}{\partial z_2} \cdot w_2^T$$



$$B \times N \times N \times H$$

= $\frac{1}{2} BNH$ flop

→ and we would need to cache q_1

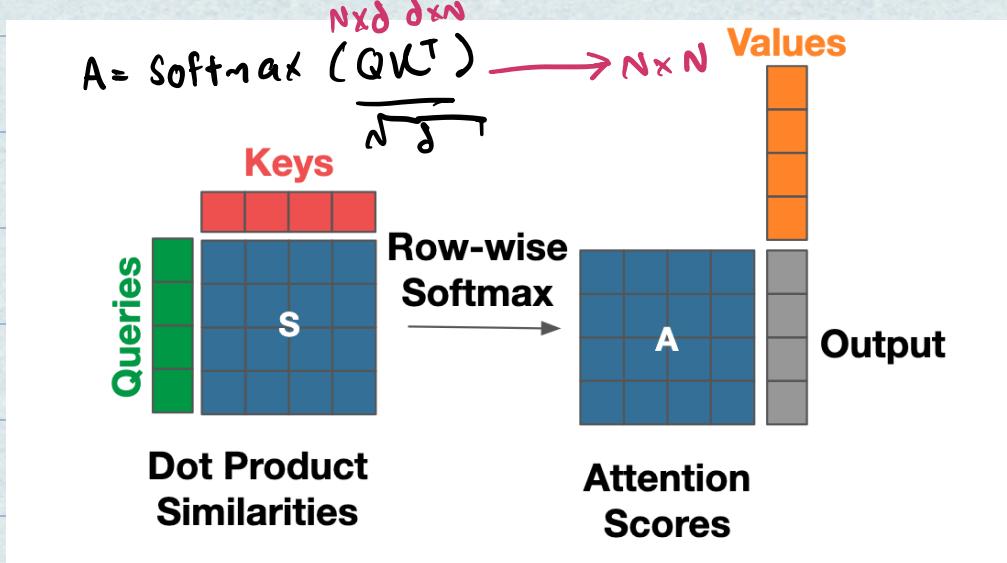
* interesting pt:

→ backprop is 2x longer than fw
($4BNH$ vs. $2BNH$)

(2) Revisit scaling in normal, softmax, attend.

→ recall: Q, K, V matrices are $N \times D$ dimension

→ therefore: final attention scores are $N \times N$



→ high level takeaway (though we went through math)

In the last class: Attention MEMORY and COMPUTE scales quadratically $O(N^2)$ in the INPUT SEQUENCE LENGTH

→ but it would be nice if our models could scale to long sequences:

- thousands of words in books or docs
- long-range dependencies between definitions and concepts at multiple scales (within a problem, chapter, across entire book)
- thousands of timesteps in single second of raw audio ; long-range dependencies across audio sequences . . .

→ next time: we will look at SW methods to improve this complexity. but today:

* IS THERE A FAST EXACT ATTENTION ALGORITHM? *

(3) HOW IS ATTENTION EXECUTED ON GPUs TODAY?

• kernel = some "unit" of operation

a GPU involves:

1. Load data from HBM to

SRAM / registers (local to

each streaming multiprocessor

2. compute (kernel can also

store intermediate data in SRAM)

3. write data (result) back to HBM

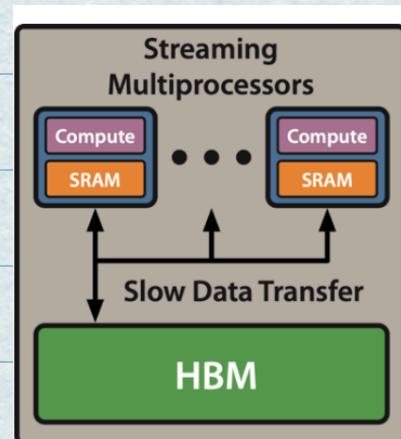


Photo: Dao et al., 2022

- GPUs can execute d.p. recalls in parallel across streaming multiprocessors!

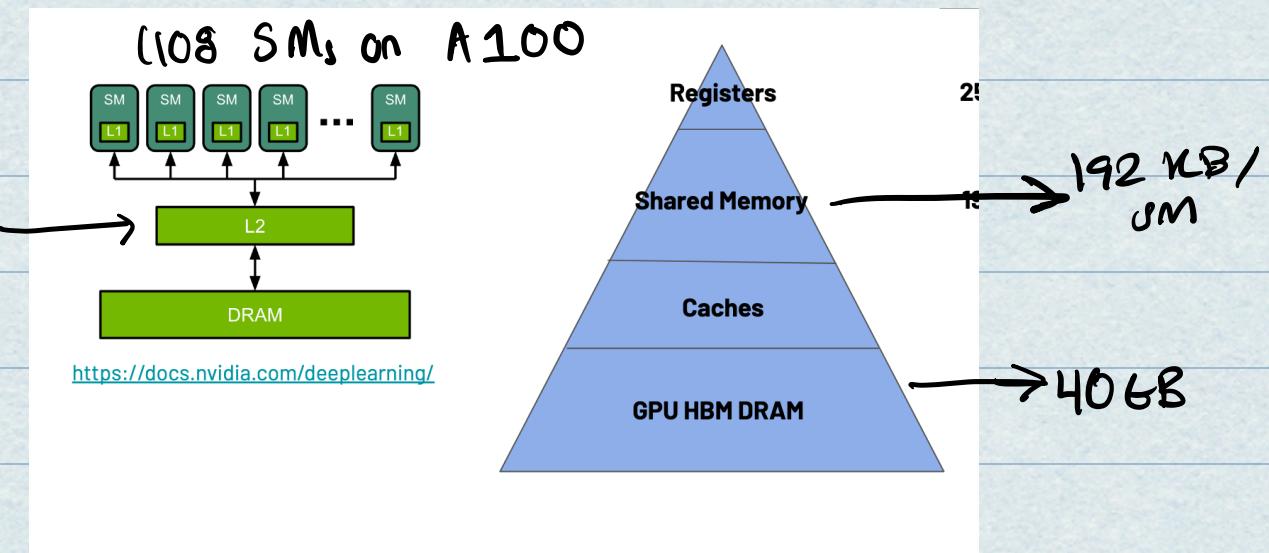
- as we explored in a previous lecture, SRAM/HBM form a memory hierarchy:

- SRAM = fast but small

- HBM = slow but larger

on A100:

~3x smaller,
but order of
magnitude faster



- how does attention currently run on a GPU?

- in most regimes - memory

is the bottleneck in

attention computation (only
compute bound when generating
208 tokens, as we saw last

class)

- in GPU - compute speed outpaces

memory speed!

- 2 further issues:

1. Attention implementations do not carefully account for

	A100 80GB PCIe	A100 80GB SXM
FP64		9.7 TFLOPS
FP64 Tensor Core		19.5 TFLOPS
FP32		19.5 TFLOPS
Tensor Float 32 (TF32)		156 TFLOPS 312 TFLOPS*
BFLOAT16 Tensor Core		312 TFLOPS 624 TFLOPS*
FPI6 Tensor Core		312 TFLOPS 624 TFLOPS*
INT8 Tensor Core		624 TOPS 1248 TOPS*
GPU Memory	80GB HBM2e	80GB HBM2e
GPU Memory Bandwidth	1,935 GB/s	2,039 GB/s

> compute BW is ~ 300 FLOPS

ops/s

HBM memory
BW is ~ 2 TB/s

reads and writes to diff. levels of GPU memory

(e.g. SRAM vs. HBM)

2. Popular DL frameworks (Pytorch) do not expose ways for finegrained memory management

* you can't fit attention computation into L1 cache - so you end up doing many reads and writes between S.M. and HBM.

- as a reminder: what is happening after Q, K, V are

calculated: $N = \text{seq. length}$ $B = \# \text{atten heads}$
 $D = \text{dimension}$

Memory Accesses:

- Read Q, K
- Write QK^T
- Read/write masking
- read/write Dropout
- read/write Softmax
- read V
- write multiply by V

Sizes	read and write?	flop	Total
$2NDB$	1	2	$4NDB$
NNB	1	2	$2N^2B$
* you will fill in table during PLQ:			

Compute:

Op:

Flops:

- Compute QK^T
masking, dropout,
softmax

$N \times D \times N$

Matmul

$B \cdot N \cdot D \cdot N \cdot 2$

notes made
amir take in lecture!

elementwise/
or reduction

$B \cdot N \cdot N \cdot 2$ → N^2 matrix

- multiply by V

$N \times N \times D$

Matmul

$B \cdot N \cdot N \cdot D \cdot 2$

softmax = element-wise exponent, sum, divide

- note that we are just talking about **attention** not feed-forward/MLP layers
- we are also ignoring Q, K, V projections:

or creating $Q = W_Q \cdot \text{inp}$, $K = W_K \cdot \text{inp}$, $V = W_V \cdot \text{inp}$

* observation: attention op has memory-boundedness because it is a series of VERY low-compute op (masking, dropout, softmax) w/ reads/writes between these ops!

* Let's Profile Attention!

recall : $AI = \frac{\# \text{ops}}{\# \text{bytes accessed}}$

if $AI > \frac{\text{HW compute BW} \cdot \text{compute bound}}{\text{HW mem BW}}$

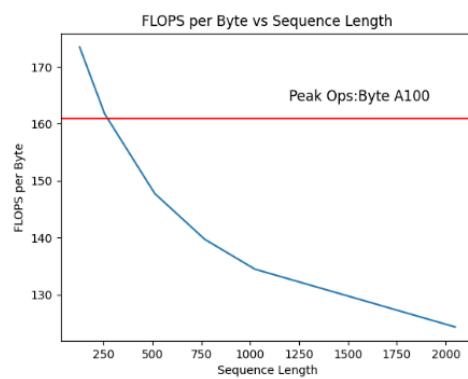
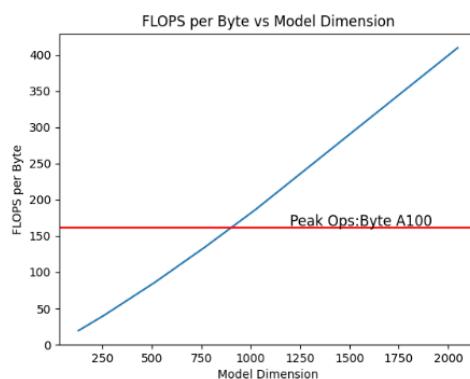
if $AI < \frac{\text{HW compute BW} \cdot \text{mem bound}}{\text{HW mem BW}}$

for A100: peak ratio is 161 ($\frac{312 \text{ Flops}}{19356 \text{ Bytes}}$)

\hookrightarrow PCIe H BW

mem BW

- we can (ignoring B) plug in values for d and D above:



* observation: as sequence length increases, AI \downarrow

less than peak ratio for A(LOD whenever $N > 250$)

↳ but in full attn block- $Q \cdot K^T$ could be compute bound, mat mul w/ V could be compute bound

- feedforward networks could be mem. bound

• recall AI of matmul:

what happens as $N \uparrow$?

$2MNK$

$2(KM + NK + MN)$

$\cdot K = M = 8192, N = 128: AI = 124.1$

↳ mem. bound

$\cdot K = M = N = 8192: AI = 2730.6$

↳ compute bound!

* takeaway: as $N \uparrow$ for attn block, becomes more mem bound, but for mat mul, it becomes more compute bound (imagine mat mul for QK^T or w/ V)

3) RECAP OF PRINCIPLES OF ADDRESSING ATTENTION

BOTTLENECK

• **Fusion** → save trips to/from memory by performing composite data ops on processing units

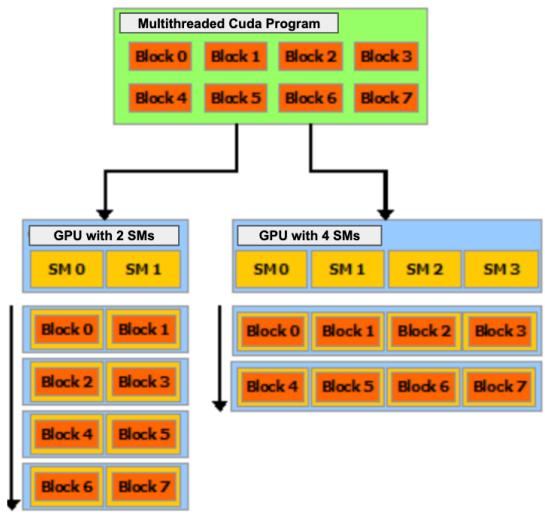
• **Tiling**: assign subsets of computation to parallel workers - ends up reducing memory traffic

↳ remember: GPU has multiple SMs (streaming multiprocessors)

- thread blocks = given to independent SMs to execute

↳ blocks usually organized into 1D, 2D, 3D

Grids



- blocks contain fine ads

- SMs schedule w/cps, or

groups of 32 threads,

to execute same intr.

Max # active threads =

$$(\# \text{SMs}) \cdot \left(\frac{\text{Max blocks}}{\text{SMs}} \right) \cdot \left(\frac{\text{Max threads}}{\text{block}} \right)$$

→ Tiling:

- recall: memory access within a thread block:

- global memory: seen by all threads (in HBM)

- shared memory: seen by threads in SAME block

↳ one thread can pull in data multiple
threads need to see

- registers: private per thread

* what is an efficient mem. loading rule of
thumb?

1) Coalesced mem. access - threads

In same warp should access adjacent cells
of memory

2) Tiling: partition data into subsets

so subset fits into shared mem.

computation on these tiles should be performed INDEPENDENTLY

each elt. accessed twice)

Access order

thread _{0,0}	$M_{0,0} * N_{0,0}$	$M_{0,1} * N_{1,0}$	$M_{0,2} * N_{2,0}$	$M_{0,3} * N_{3,0}$
thread _{0,1}	$M_{0,0} * N_{0,1}$	$M_{0,1} * N_{1,1}$	$M_{0,2} * N_{2,1}$	$M_{0,3} * N_{3,1}$
thread _{1,0}	$M_{1,0} * N_{0,0}$	$M_{1,1} * N_{1,0}$	$M_{1,2} * N_{2,0}$	$M_{1,3} * N_{3,0}$
thread _{1,1}	$M_{1,0} * N_{0,1}$	$M_{1,1} * N_{1,1}$	$M_{1,2} * N_{2,1}$	$M_{1,3} * N_{3,1}$

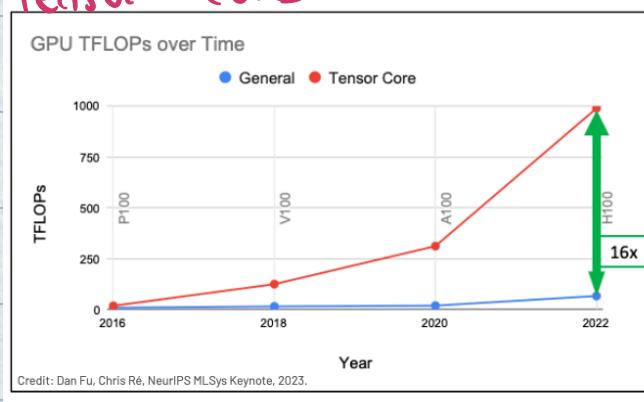
- observe that the threads overlap in some of the data they access
- if they can collaborate (e.g. share data) in some

SM's shared mem) -
save 1/2 of mem access

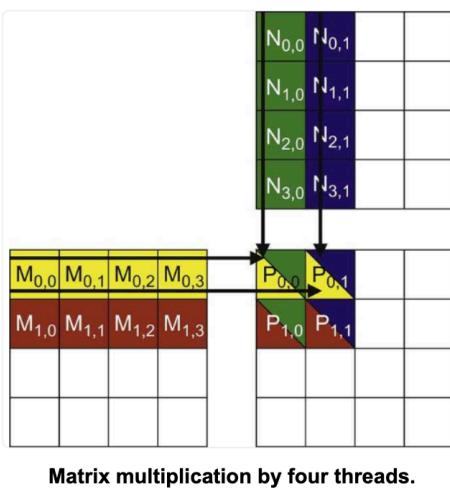
→ use Faster Specialized HW units (e.g.) **tensor cores**

- modern GPUs contain tensor cores - or specialized compute for performing MM, ($D = AB + C$)
- tensor cores can generally multiply 16×16 tiles

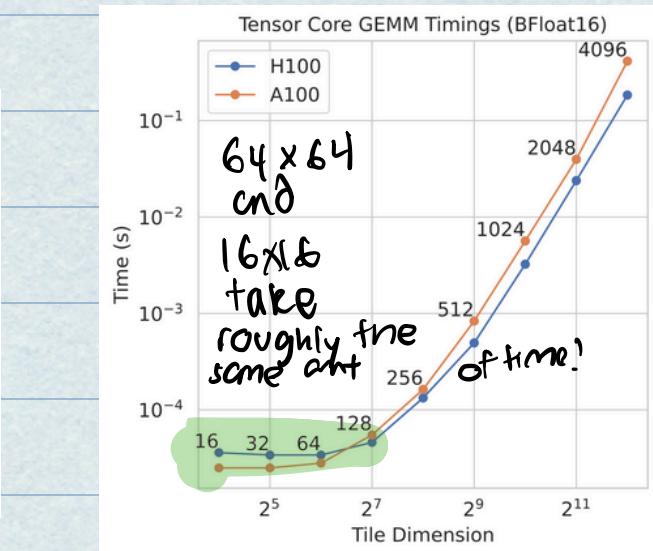
* key idea: algorithms should leverage tensor cores. On later and later GPUs, increasing speed diff. w/o tensor cores



<https://nichijou.co/cuda7-tiling/>



* key idea: select tile sizes to keep tensor cores busy (launching take times)



→ Idea 4: Caching and Recompute

- Backprop requires storing intermediate values from FW pass, adding to reads/writes

- when memory bound (reads/writes expensive relative to compute) - recompute activation

- when compute bound (compute expensive relative to reads/writes) - store activation

• **FLASH ATTENTION **

- overview: new implementation of standard attention^(exactly), that carefully accounts for reads/writes by using above ideas:

1) recompute for memory bound workloads, cache for compute bound workloads

2) Tiling

3) Fusion

- for simplicity, we'll only discuss: $\text{Softmax}(QK^T) \cdot V$
↳ you can sometimes add masking / dropout along w/ softmax, can be fused together

- overview of results:

→ HBM accesses, given $N = \text{sequence length}$
 $d = \text{dimension}$
 $M = \text{SRAM size}$

$$\text{flash attn: } O\left(\frac{N^2 d^2}{M}\right)$$

standard attn: $O(Nd + N^2)$ → about 9x
fewer HBM writes

Attention	Standard	FlashAttention
GFLOPs	66.6	75.2 ($\uparrow 13\%$)
HBM reads/writes (GB)	40.3	4.4 ($\downarrow 9x$)
Runtime (ms)	41.7	7.3 ($\downarrow 6x$)

Model implementations	OpenWebText (ppl)	Training time (speedup)
GPT-2 small - Huggingface [87]	18.2	9.5 days (1.0x)
GPT-2 small - Megatron-LM [77]	18.2	4.7 days (2.0x)
GPT-2 small - FLASHATTENTION	18.2	2.7 days (3.5x)
GPT-2 medium - Huggingface [87]	14.2	21.0 days (1.0x)
GPT-2 medium - Megatron-LM [77]	14.3	11.5 days (1.8x)
GPT-2 medium - FLASHATTENTION	14.3	6.9 days (3.0x)

• 3x faster training compared to huggingface!

• 15% faster pre-training than BERT pre-training

record in Mperf (the olympics of efficient

model training)

• also in some cases results in higher quality models - can train on longer sequences!

- Let's incrementally build flash attn:

• approach: lets avoid materializing the full $N \times N$ mat \times during the attn. computation

→ Try 1: what if we loaded all of Q and all of K into stream-multiprocessor's SRAM to compute attn output?

↳ an $N \times N$ matrix needs $2N^2$ space

↳ but a S.M. has 20-40 MB of SRAM

Context length (N) | Size of $O(2N^2)$

500 | 0.5

1000 | 2MB

4000 | 32 MB

But we could be trying to fit multiple BATCHES of inputs