

□ Agenda:

- more review of roofline model
- Basic Principles for Achieving High Perf. on Hardware (w/ focus on some computations in attention)
- Perf. Analysis of Transformers
- If time: KV caching, advanced techniques for inference efficiency

□ Sources:

- End of cs229s lecture on Hw: <https://docs.google.com/presentation/d/14hK7SmkUNfSEIRGyptFD2bGO7K9sJOTnwjAVg3vgg6g/edit?usp=sharing>
- cs229s Lecture on Analyzing Transformers:
<https://docs.google.com/presentation/d/1PV0cKnzcHRCAbJy3Oter1UdJME7T7WBEqLLYmKWJR4g/edit?usp=sharing>
- Blog post on roofline model: <https://dando18.github.io/posts/2020/04/02/roofline-model>

* Part 1: review of roofline model

- remember units of arithmetic intensity
- $AI = \frac{\text{Ops / time}}{\text{memory traffic}} \rightarrow \text{eg. } \frac{\text{Flops/second}}{\text{bytes/second}} = \frac{\text{flop}}{\text{byte}}$
- this metric is important because in general: memory ops are much slower than compute ops
 - on a CPU:
 - FLOP like add or multiply: n/cycle

→ I/O from register: ~1cycle

→ I/O to L1 cache: ~1ns
(2cycles)

→ I/O to L2 cache: ~4ns
(8cycles)

→ I/O to main memory: 100ns
(200cycles)

for 2GHz processor:
2 cycles = 1ns

- in general in our algorithm we should strive to do as much computation before loading

next data

- high AF: lots of computation per load

- low AF: low amt of computation per load

* consider AxPy operation:

$$y = ax + y \quad a \in \mathbb{R} \quad x, y \in \mathbb{R}^n$$

→ to impl in C:

```
void daxpy (size_t n, double a, const *double x,  
            const *double y);
```

```
for (size_t i=0 ; i<n ; i++) :
```

$$y[i] = a \cdot x[i] + y[i]$$

· number of flops: $2 \cdot n \rightarrow 1 \text{ add and 1 multiply per loop iteration}$

· memory traffic: $\left(2 \text{ loads} \cdot \frac{8 \text{ bytes}}{\text{load}} + 1 \text{ store} \cdot \frac{8 \text{ bytes}}{\text{load}} \right) \cdot n$
 $= 24n$

$$\cdot \text{Arithmetic intensity: } \frac{2n}{24n} = \frac{1}{12}$$

* consider matrix multiply for matrix $(M \times K) \cdot (K \times N)$

$$\rightarrow \text{from last time: } \frac{2MNK}{(MK + KN + MN)}$$

~ generally: MN will do more math per byte
for larger matrix sizes

* roofline model

• from HW specs we know:

- Peak Bandwidth: bytes/s

- Peak compute perf: flops/s

• the Actual perf. of an algorithm

i.e.

$$\min \begin{cases} \text{AI} \cdot \text{Peak bandwidth} \\ \text{Peak compute perf} \end{cases} \xrightarrow{\text{flop}} \frac{\text{flops}}{\text{byte}} \cdot \frac{\text{bytes}}{\text{s}} = \frac{\text{flop}}{\text{s}}$$

• intuition: no alg can have higher flop/s than hardware limit, but alg can be limited

by BW

* note that in reality, there might not be 1 BW slope or 1 compute ceiling.

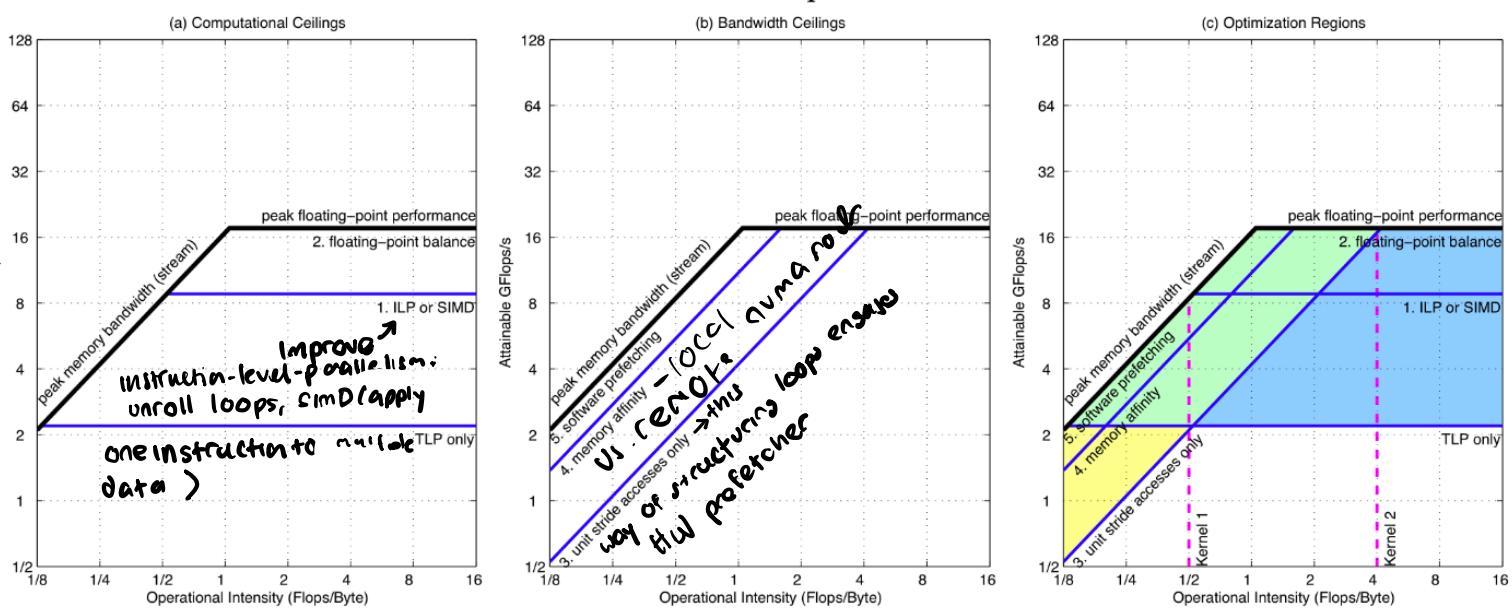
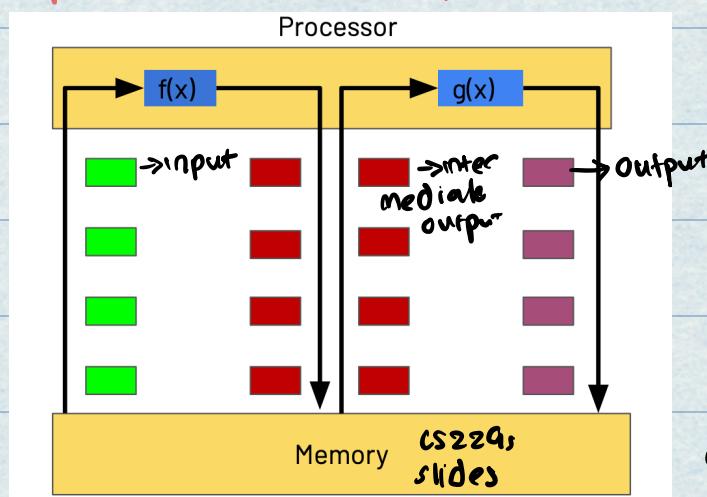


Figure 2. Roofline Model with Ceilings for Opteron X2.

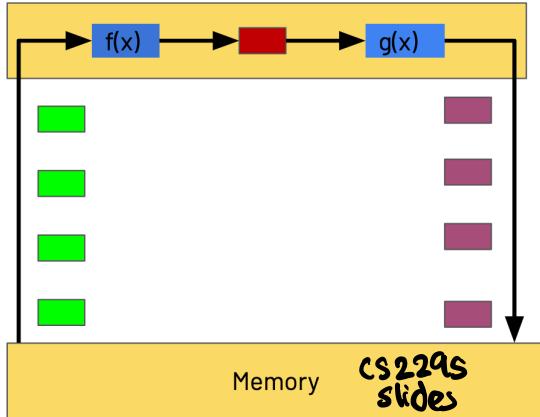
- Depending on where exactly your algorithm falls - roofline also suggests a series of compute and memory optimizations you can apply to reach peak of both memory and flops/second

*Part 2 of lecture: Basic Principles of Achieving High Perf. on Hardware

1) Fusion: save trips to /from memory by performing composite data on processing units



→ we write outputs of $f(x)$ back to memory, load them back, before then computing $g(x)$

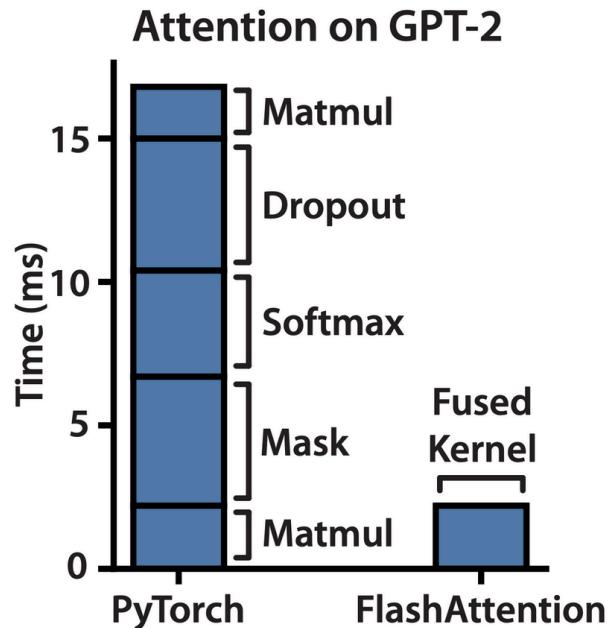
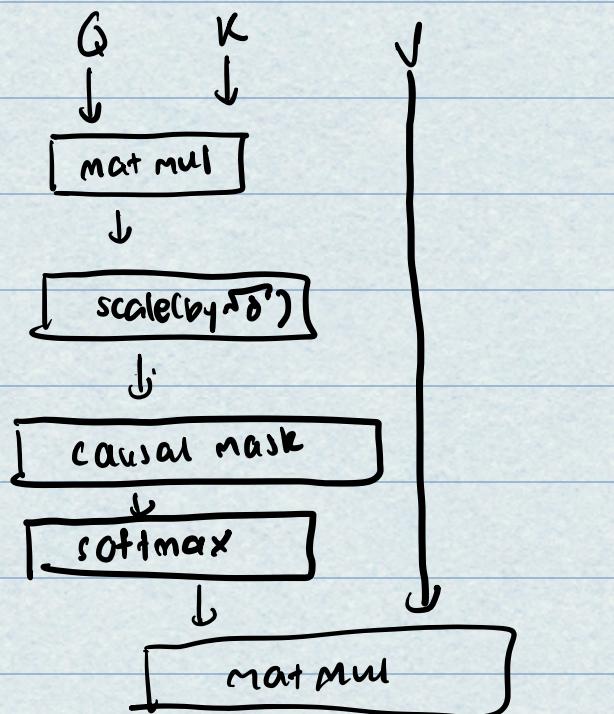


-Fusion prevents unnecessary trips to HBM

-critical for I/O bound ops!

* Fusion for attention (flash attention!)

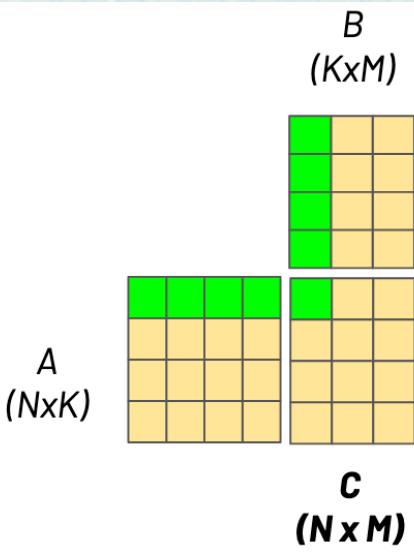
recall:



[FlashAttention: Fast and Memory-Efficient Exact Attention with IO-Awareness](#)

2) Parallelization (in context of matrix multiply)

- recall : A 100 GPU has 108 streaming multiprocessors
- how do we ensure each sm has as much work as possible
- idea: perhaps parallelize across OUTPUT cells
 - ↳ ideally at least 108 output items



→ we can create one thread PER output cell value and put each thread in a sp. thread block - the dif. thread blocks will go to dif. SMs

→ what would arithmetic intensity be?

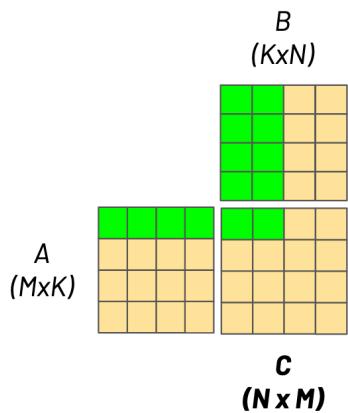
→ compute: Dot product of row of A and col. of B (K adds, K multiplies)
 $= 2K$ FLOPs

→ mem. (K) (K)
 1) read row of A, col. of B
 2) write one val. for $C^{(1)}$] = $2K+1$ memory accesses

$$\rightarrow \text{arithmetic intensity} = \frac{2K}{2(2K+1)} + \text{multiply by 2 for fp16}$$

2) Idea: Blocking / Tiling w/ 1D tiling

- each thread computes 2 tile items in output matrix



→ compute: dot product of 2 cols of B, one row of A = $4K$

→ memory:
 1) read row of A = K
 2) read 2 cols of B = $2K$
 3) write 2 outputs of C = 2] = $3K+2$ memory accesses

$$\rightarrow \text{new AI} : \frac{4K}{2(3K+2)}$$

+ depending on problem size - might set (turn off)
caching

2) 2D tiling

\rightarrow each thread computes 2D tile of
output C

\rightarrow compute = $8K$ flops (4 output cells)

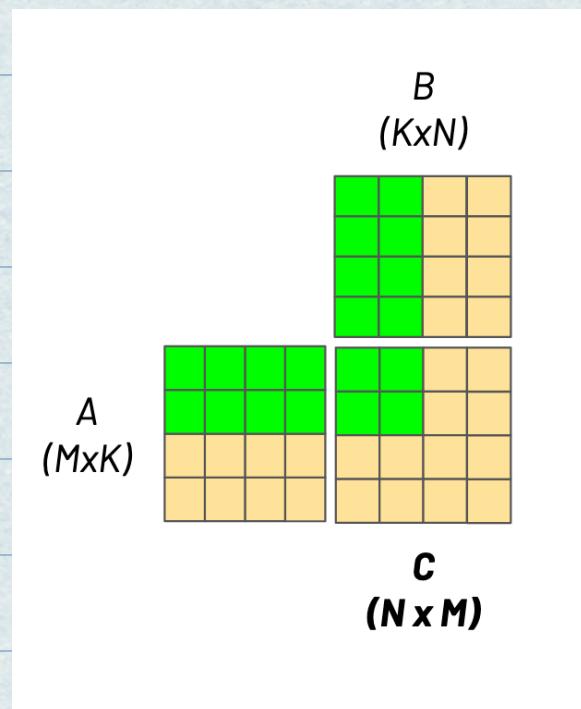
\rightarrow mem. access = $4K + 4 \rightarrow$ but n FPIB

$$\text{AI: } \frac{8K}{2(4K+4)}$$

\rightarrow general idea of blocking/
tiling: assign TILES of work

to each parallel processing unit -
exploits some mem. locality (reuse of memory
being loaded)

- note that AI does depend on K, though



3) Idea 3: caching and recomputation

- for compute bound workloads - it makes sense to CACHE
- for memory bound workloads -

Instead of caching - lets recompute
- recall auto diff. for backprop

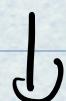
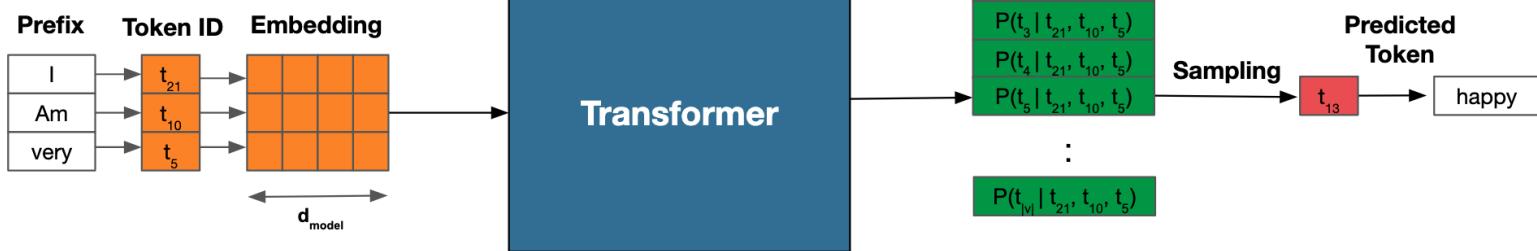
- we need to store activations to compute backward pass-gradients

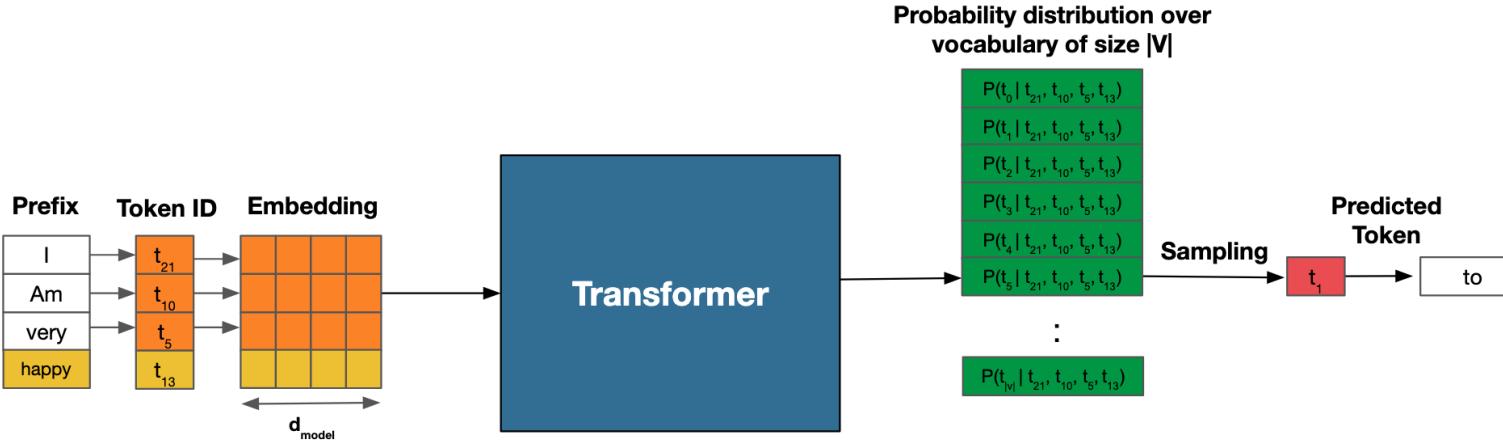
* In some cases - maybe don't store all intermediate activations; maybe better to actually recompute forward passes!

* Let's look at caching vs. recompilation in context of ↗ AUTOREGRESSIVE GENERATION ↗

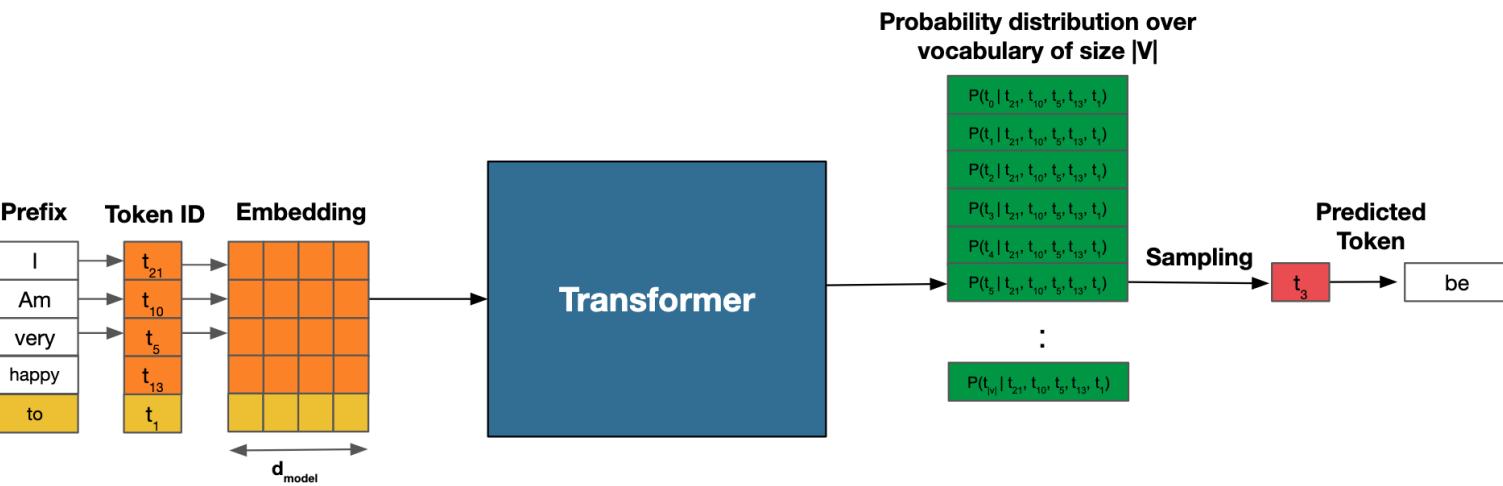
- In Inference - we take current prompt sequence - and predict next word

autoregressive = "predict one token at a time" - can't predict next token until entire current sequence has been produced

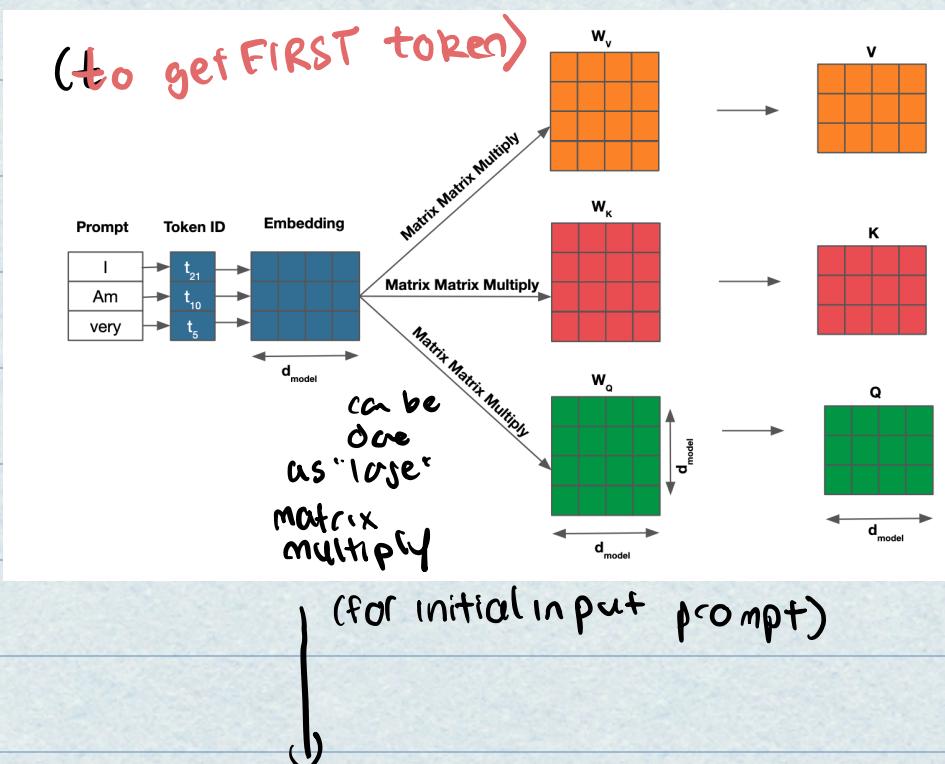




↓



How does this manifest in terms of math being done?

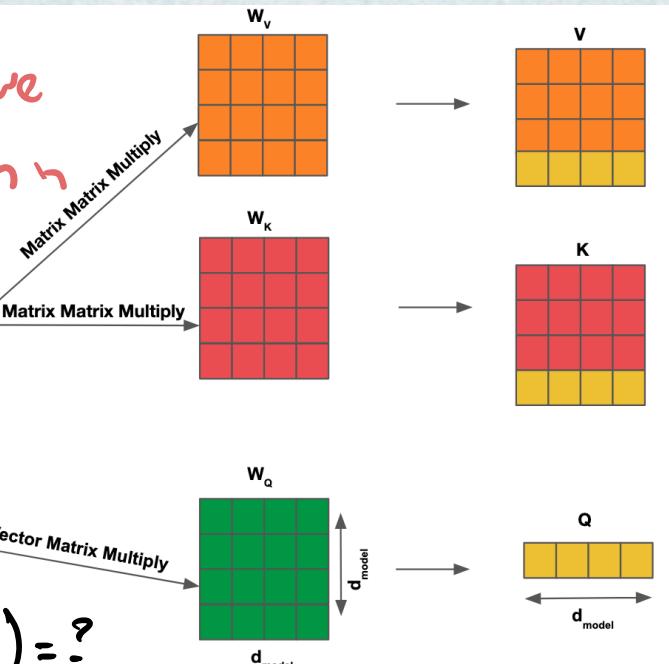
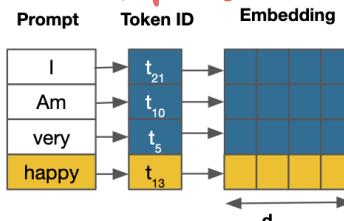


→ for each token in sequence we calculate:

- key and query vectors
- for each token t_i
 - find attention scores of all tokens in sequence
 - use attention scores to get value matrices

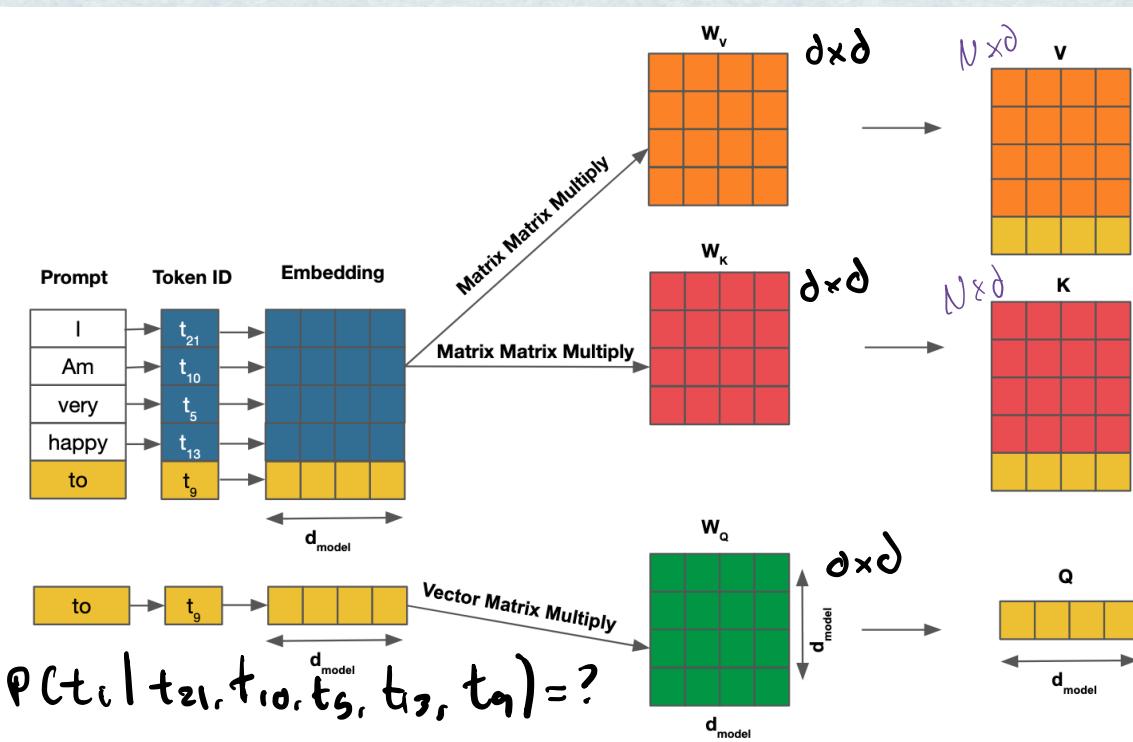
- But what happens when we predict next token in sequence?

* know that we have "happy" how do we add NEXT token in sequence



→ for next token:
- we only need to calculate one new query vector, one additional key vector; and one resulting value vector

Similarly:



- again:
- we only need to generate one query vector
- and add one row each to value and key matrices

* what are the FLOPs in the attention calculation?

→ let's say we are pre-processing input sequence to produce FIRST token prediction?

• computing K: $O(N \cdot d \cdot d) = O(2Nd^2)$

• computing Q: $O(N \cdot d \cdot d) = O(2Nd^2)$

• multiplying Q and K: $S = QK^T = O(2N^2d)$

• letting $A = \text{softmax}(S)$ = had to account for exactly, but not bottleneck

• computing V: $O(2Nd^2)$

→ now we've predicted happy. how do we predict next word?

1. compute K for full sequence: $O(2Nd^2)$

2. compute new query vector: $O(2d^2)$

vector matrix multiply
 $1 \times d$ $d \times d$ \rightarrow

3. compute S: $QK^T = O(2Nd)$
 $1 \times d$ $d \times N = 1 \times N$

4. A = softmax(S): had to account for exactly

5. compute V for entire sequence so far: $O(2Nd^2)$

6. Compute $A \times V = \text{attention weights for next pred.} : O(2Nd)$
 $1 \times N$ $N \times d = 1 \times d$

* Key idea: there's a TON of repeated computation

If we do this for EACH NEW TOKEN

- Introduce concept of ** KV CACHING **

→ in producing any subsequent token

after the FIRST one, we DON'T need to

calculate all of K and V again

(as most of it is unchanged):

→ we ^{should} reuse K and V as much as possible:

- new protocol for each NEW token:

1. compute K/Q/V VECTORS for new token:

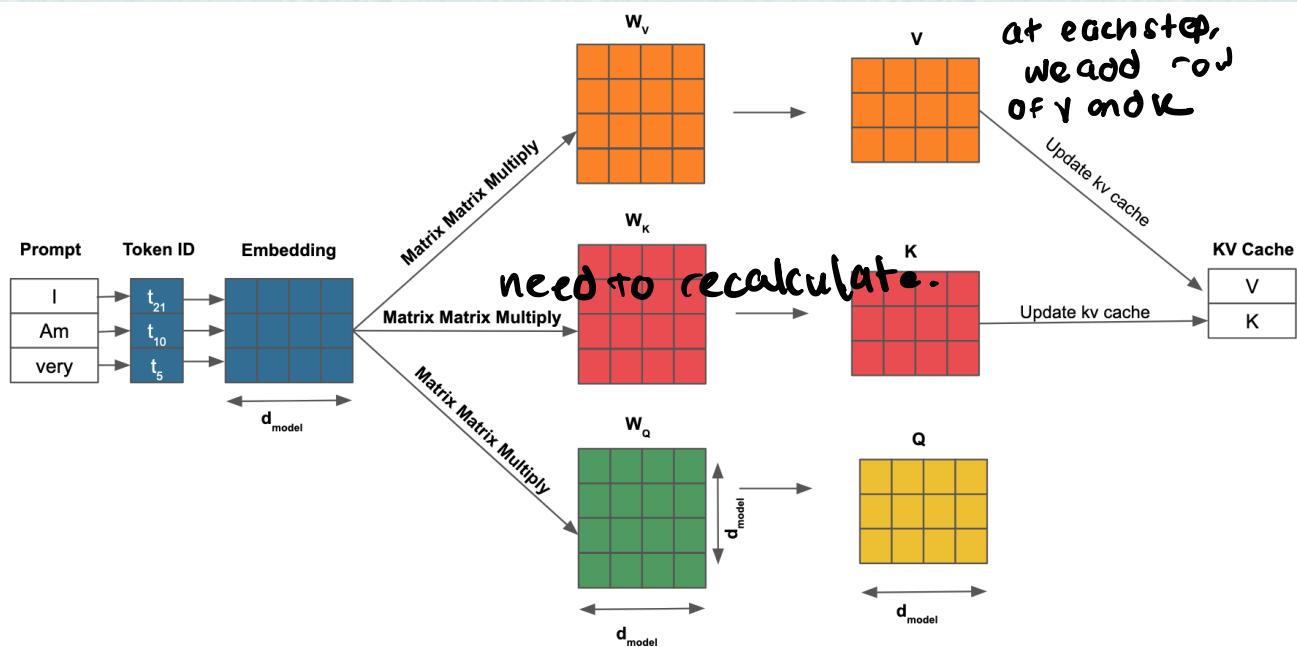
$$\rightarrow \text{all } l \cdot d \times d \times d \text{ MMs} = O(3 \cdot 2d^2)$$

2. multiplying $Q \cdot K^T = O(2Nd)$
 $l \cdot d \quad d \cdot N$

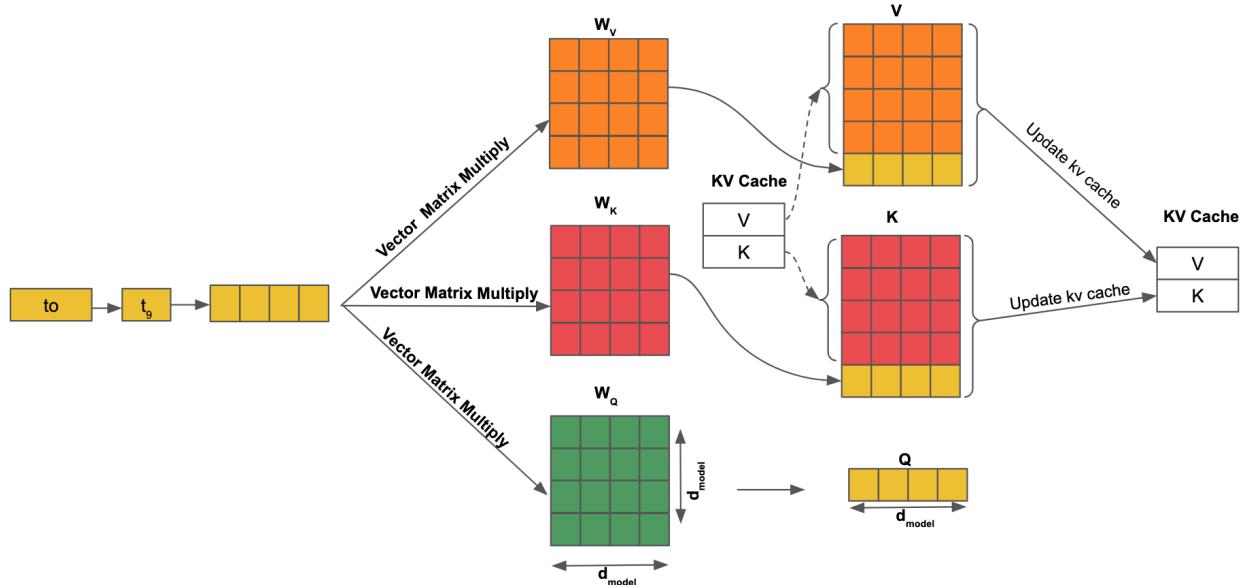
3. A = softmax(S): hard to account for exactly

4. multiplying A by V : $O(2Nd)$
 $l \times N \quad N \times d$

new picture:



keep V/K cached in
GPU DRAM so we don't recalculate!



* this was all for SINGLE HEAD attention, would be diff. for multi-head

* In summary, for each new token our compute is:
(w/ KV-caching)

$$O(3 \cdot 2d^2) + O(2Nd) + O(2Nd)$$

→ let's try to understand WHEN KV caching is worth it:

→ total compute for calculating KV per token

$$O(2 \cdot n_{\text{layers}} \cdot 2d^2)$$

↗ k and v ↘ # attention layers ↑ matrix multiply for
 layers (1x d) - (d x 0)

→ total storage for k and v vector

per token

$$O(2 \cdot 2 \cdot n_{\text{layers}} \cdot d)$$

↗ # attention layers ↗ dimension of
 k and v vectors

k and Fp16

* so this needs to be stored PER TOKEN

→ recall on A100:

$$\text{compute} = 312^{12} \text{ FLOPs/s}$$

$$\text{mem BW} = 1.5^{12} \text{ Bytes/s}$$

→ How do we calculate arithmetic intensity
of calculating $W_k \cdot x$ and $W_v \cdot y$

1) Compute → compute term above:

$$O(2 \cdot n_{\text{layers}} \cdot d^2)$$

↳ represents $W_k \cdot x$ and $W_v \cdot x$

2) Memory:

1) Loading W_k and W_v weights is:

$$O(2 \cdot 2 \cdot n_{\text{layers}} \cdot d_{\text{model}}^2)$$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
 $K, V \text{ fp16} \quad \# \text{ attn} \quad W_k \text{ and } W_v$
 n_{layers} one $d \times d$ matrices

2) Writing output of $W_k \cdot x$ and $W_v \cdot x$

(activation) is from pluggins term

above:

$$O(2 \cdot 2 \cdot n_{\text{layers}} \cdot d)$$

- since loading weights is more,
we will just consider that for simplicity:

$$T_{\text{compute}} = \frac{4n_{\text{layers}} \cdot d^2}{312 \cdot 10^{12} \text{ FLOPs}} \quad \begin{matrix} \text{compute} \\ \text{to} \\ \text{nem.} \\ \text{ratio} \end{matrix} = \frac{T_{\text{mem}}}{T_{\text{math}}} = \frac{312 \cdot 10^{12}}{1.5 \cdot 10^{12}}$$

$$T_{\text{mem}} = \frac{4n_{\text{layers}} d^2}{1.5 \cdot 10^{12} \text{ Bytes/s}} \quad \boxed{= 208}$$

* what does 208 mean here?

- once we've loaded weights into memory,
it is worth it to compute 208 tokens

(run computation 208 times)

- fewer than 208: memory bound
 - greater than compute bound:
- we will explore in
a later class -
another way to
exploit this is
via batching across
requests!

* so when does KV caching make sense?

- if we are compute bound - KV cache makes sense!
(trade off computation for free extra storage)