

Agenda:

- Memory management in LLM serving & Paged Attention
- Batching, Throughput, Latency
- motivation for continuous batching
- "Continuous" Batching in LLM systems: if time

* First classification on division of Q,K,V for flash arm

+ consider a single attn head, where for one example,
Q,K,V are all $N \times d$

1) Divide Q into blocks of size $B_r \times d$

$$\hookrightarrow T_r = \frac{N}{B_r} \text{ blocks}$$

2) Divide K,V into blocks of size $B_c \times d$

$$\hookrightarrow T_c = \frac{N}{B_c} \text{ blocks}$$

↳ iterating over Q

→ For $1 \leq i \leq T_r$

• Load Q_i ↳ each block output
each loss sum exp

• Initialize $O_i^{(0)}$, $l_i^{(0)}$, $m_i^{(0)} \rightarrow \max$

→ For $1 \leq j \leq T_c$ ↳ iterating over K,V

• calculate local numerator, denominator

• using previous output, calculate this output (iterative)

↳ where warp scheduling comes into play.

→ splitting heads, Q over thread blocks to fill S.M.s

→ what about wops?

↳ within a block, FA1 splits K,V across wops, Q accessible to all wops

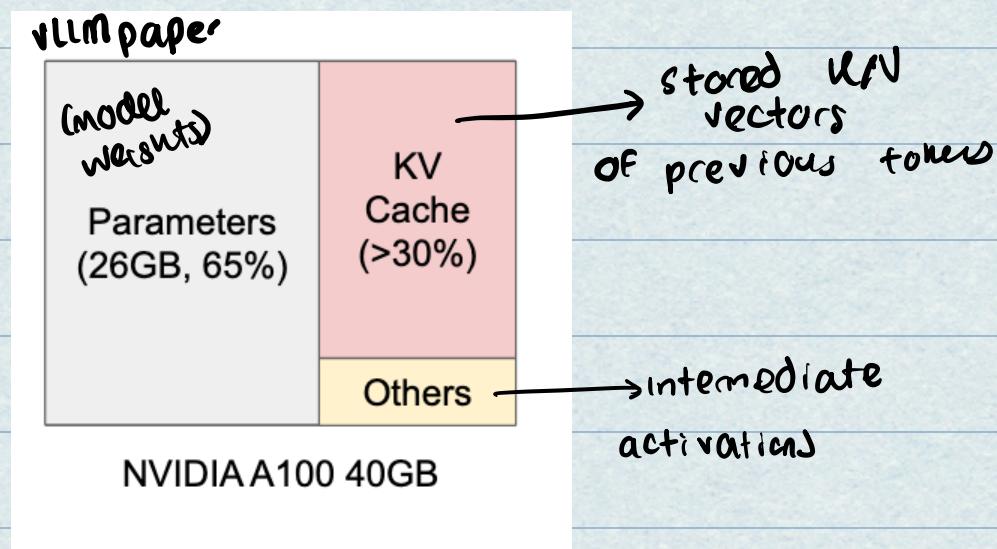
↳ instead, FA2: split Q across wops, K,V accessible by all

↳ think about inner inner part of for loop → we have specific K,V blocks

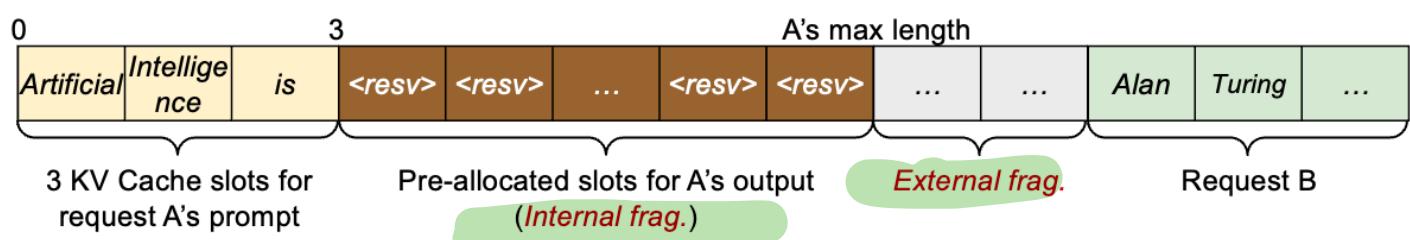
* see FA2 paper for more details! Q block → how do we use scrafe

FIRST PART OF LECTURE: PAGED ATTENTION

- vLLM's original main technique
- core question: how do we **MANAGE** GPU memory (for KV cache, specifically)



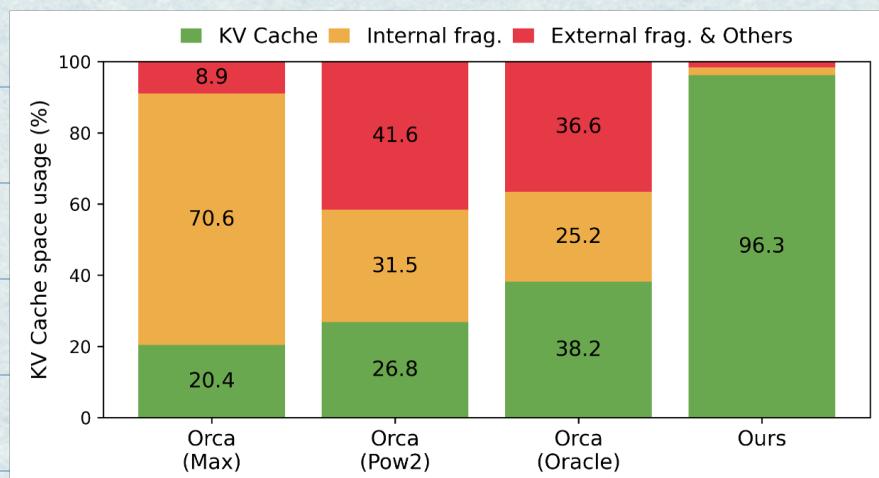
- a few things make memory management hard:
 1. For each request doing autoregressive decoding, we DON'T KNOW how long the generation will be
 2. Across requests, users can provide dif. max request lengths
- * previous engines allocated contiguous spaces of memory in KV cache for each request's max memory:



* some observations:

- A and B have DIF max request lengths
 - engine pre-allocates slots up to max length
 - this "static memory management":
 - leads to 2 types of fragmentation
 - 1) Internal frag: LLM does not predict up to max request length. wasted space
 - 2) External frag: when allocating memory requests of dif. sizes, there will always be some unused space
- + go study memory allocation schemes for more info

FRAAGMENTATION IS SIGNIFICANT:



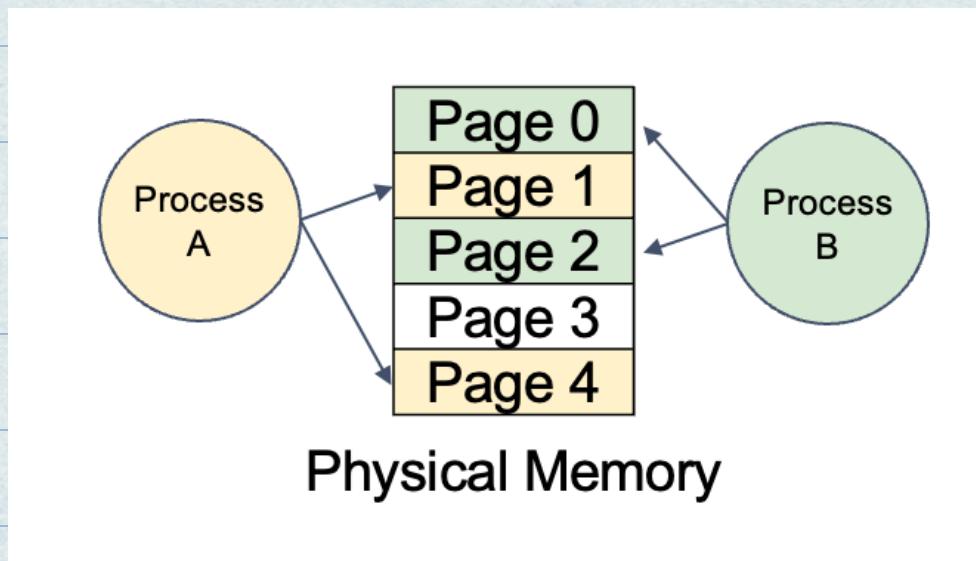
(from VLLM paper):
- in previous baselines -
only 20 - 40 % of
available memory for
KV cache is used
for useful data

- GPU device memory is the bottleneck - we really shouldn't be wasting it!

KEY INSIGHT OF VLLM WORK:

- we've already solved a similar problem w/ OS-level paging and virtual memory:

- OS divides up memory into "Pages"
 - will map program LOGICAL pages to Physical pages
- ↳ so logical pages need not be physically contiguous:



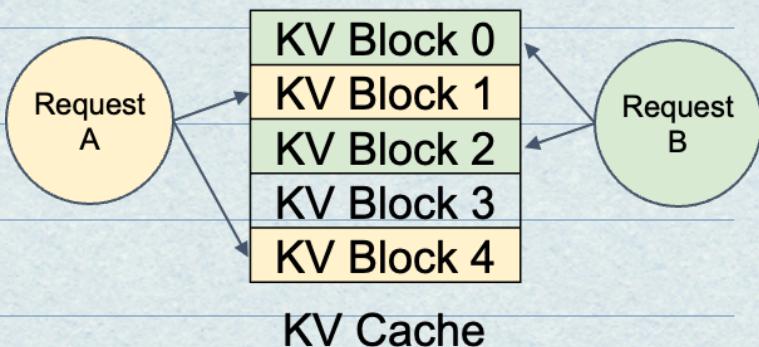
But
each
process
THINKS
it has logically
contiguous
memory

Let's do something similar for the KV cache!

- partition KV cache into "blocks"

- to minimize fragmentation, blocks for

dif. requests do NOT have to
live next to each other:

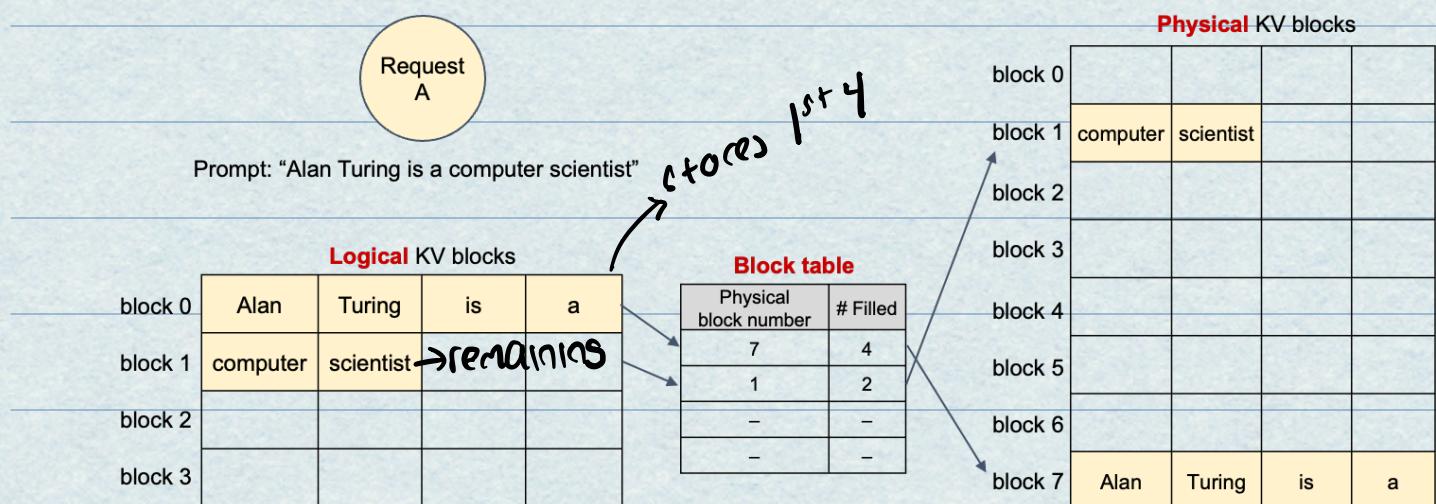


KV-block = fixed-size contiguous memory

chunk that stores KV states from left to right.

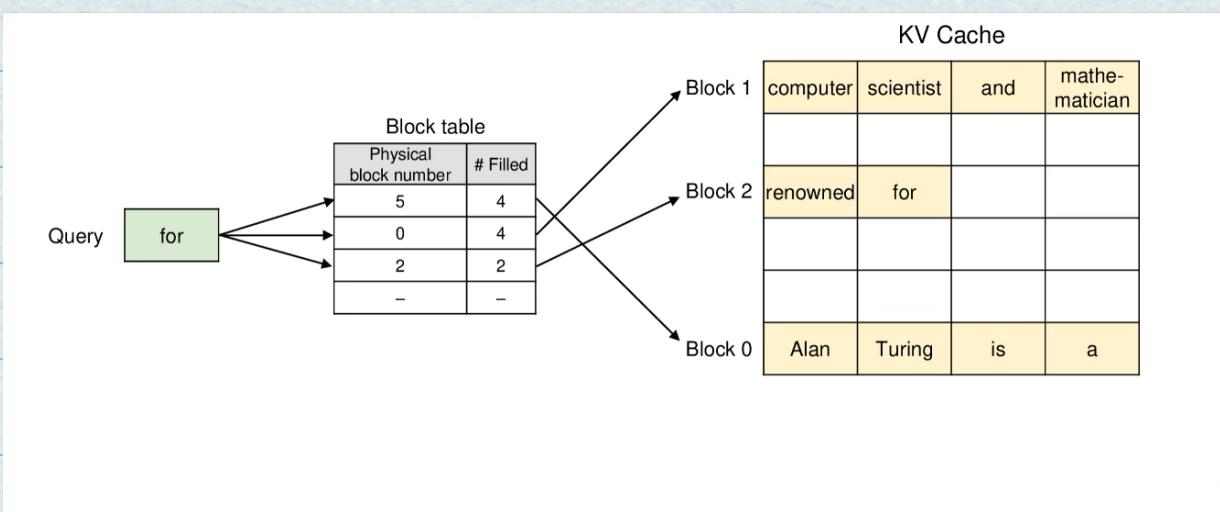
↳ block-size of 4 would be able to store k/v projections for 4 tokens

* now we can VIRTUALIZE THE IN CACHE:



↳ request "thinks" B0 and B1 are contiguous, but they are physically not: (w/ block/
page table in between)

* How do we compute attention?



- we need to compute attention for query "for"- w.r.t respect to each token before; but KV blocks cannot be fetched all at once!

- iterate over blocks containing KV cache for this request:

- for each logical block, set row of block table

- look up physical block location

- calculate attn of query w.r.t those keys/values:

- we are trying to calculate an attn

score for each token so far:

- consider block j and we are trying to compute attn score of token i (i is query)

$$A_{ij} = e^{\frac{q_i \cdot K_j^T}{\sqrt{D}}}$$

$$o_i = \sum_{j=1}^{l/B} V_j A_{ij}^T$$

→ iteration for $t=1..l/B \rightarrow$ all

blocks so far! → gets w proper softmax sum

→ in o_i : again iterate from $j=1$ to i/B to get all values

contains keys/ value
for token index:

$$B \cdot (j-1) + 1 \dots B \cdot j$$

↳ for $B=4$:

$K_1 \rightarrow$ keys for tokens 1..4

$K_2 \rightarrow$ keys for tokens 5..8

$V_1 \rightarrow$ vals for token

$V_2 \rightarrow$ vals for tokens 1..4 token 5..8

$$\begin{aligned} l &= 4(1-1) + 1 \\ s &= 4(2-1) + 1 \end{aligned}$$

→ One more note: how do we make space for KV of newly predicted tokens?

↳ first: try to fill end of existing block

↳ then: allocate new blocks

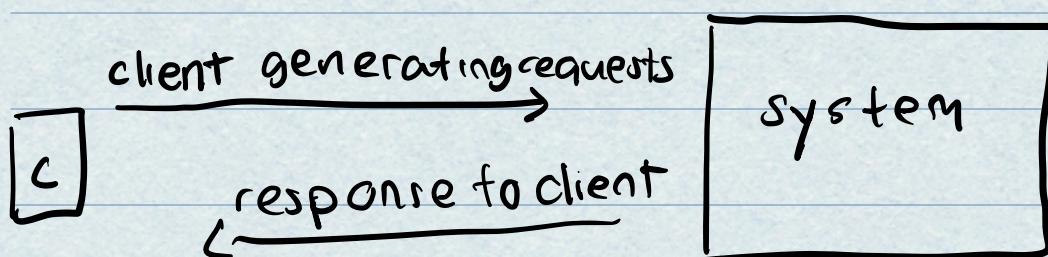
*these kind of techniques open way for:

- prompt KV caching **ACROSS** requests

- prompt "tree" caching (e.g., when requests share some lowest common prefix)

... and more:

PART 2 OF LECTURE: BATCHING, THROUGHPUT, LATENCY



* **latency**: how long does it take for the system to process a request? measurement: time (e.g., sec or ms)

↳ generally - it's easier to measure 2-way latency (there, processing, and back)

- in some cases - the time b/w C

and S is negligible - so we are effectively measuring processing latency

* **throughput**: how many bytes are being transferred over the network in unit time?

↳ measured in Data /s e.g. bytes/s

- throughput could include "not-useful" bytes (e.g., overheads like packet headers)

- therefore, another useful measure is:

* **goodput** - how much USEFUL DATA is being sent over network

↳ could be measured in something like "requests /s"

* **utilization**: how much of the available processing capacity are we using?

↳ we saw this w/ memory vs. compute utilization

* **Discuss for 5 minutes**: why can achieving

BOTH low-latency (across all users) and high

goodput be conflicting goals?

- we can achieve low-latency at expense of goodput or throughput:

- always admit 1 request to be processed at a time
- drop all requests that arrive while the 1 request is being processed
- the 1 request gets all of system's resources

• we can achieve high throughput at expense of latency:

- wait for B requests to arrive
- execute B at same time! usually beneficial if some amortization
 - ↳ $B \times N \times d \rightarrow$ dimension (e.g. load some model)
 - batch size
 - seq. length
 - weights for B items

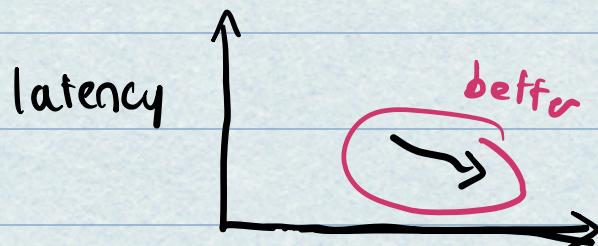
- but for requests at beginning of B :

we've waited to execute them

* How do we think about / measure throughput and latency?

* plug for CS167S (being offered in Fall)

- let's first ignore batching → assume we get no amortization from executing requests at some time:

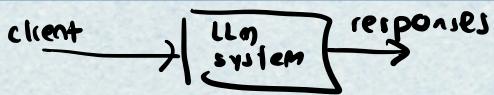


Discussion

- why is lower right of graph better?

throughput / throughput (measured)

- how would we measure this?

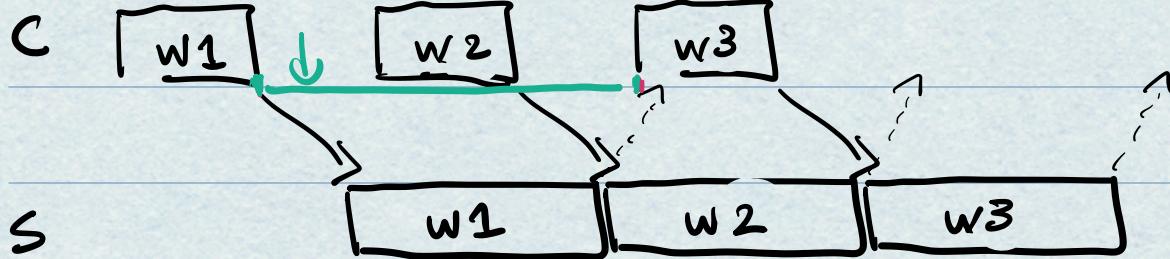


* input: offered load

↳ client sends at some rate r

* output: achieved load

latency distribution
latency of w1



• $w_2 - w_1 = 1/(\text{client rate}) \rightarrow$ even:
(or $w_3 - w_2$)

• blue line: latency of $w_1 \rightarrow$ if we measure all latencies, we can get a distribution

- here: we are assuming client sending at even rate, all requests take same amount of time, system can't operate on more than 1 request per time

• let's break those assumptions!

1) in real world, clients may not be

sending evenly at some # of requests

per second (requests evenly spaced out over second)

↳ usually systems will have some kind of load balancer in front and enough replicas will be provisioned to handle load

→ take cs1675 to find how to model this in system measurement!

2) requests may NOT take same amt of time

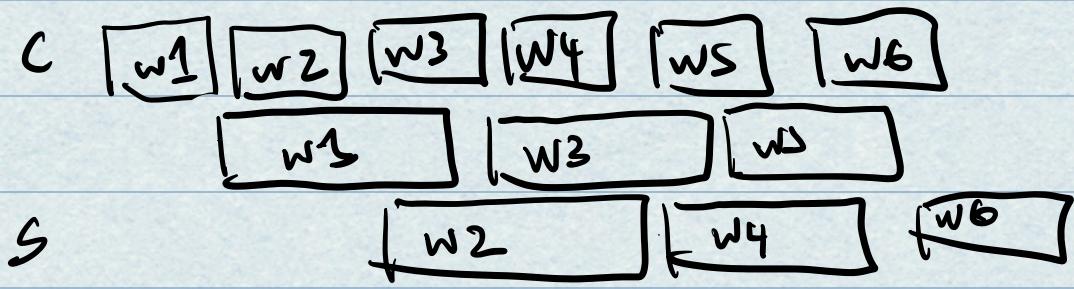
→ consider LL ms.

→ req 1 may generate 50 tokens,
req 2 may generate 2 tokens

* this is one reason why generative AI systems are different:

· for CNN, all images take same amt of time to process!

3) our system may have capacity to execute multiple reqs at same time



] can execute 2 simultaneously.

- how might the size of inference batches be controlled on a GPU?
 - consider available memory to store KV cache!

- we have 7B param model, 40 GB GPU
 - 24 layers
 - $D = 2048$
 - seq. length (for all examples)
 - = 1024

1) memory to load model:

$$7 \times 10^{12} \cdot d = 14 \text{ GB}$$

2) remaining memory for KV cache: ~26GB

per token overhead for KV cache:

$$\begin{aligned} & \text{(K,V)} \quad d \times d \overset{\text{(fp16)}}{\times} n_{\text{layers}} \times d_{\text{model}} \\ &= 4 \times 24 \times 2048 \\ &= 200 \text{ k bytes} = .0002 \text{ GB} \end{aligned}$$

3) How many tokens can we fit in KV cache?

$$\frac{26 \text{ GB}}{.0002 \text{ GB}} \approx 130 \text{ K Tokens}$$

4) For fixed seq. length, how many batch

items can we fit? $\frac{130 \text{ K}}{1024} = 128$

• naive algorithm: (at LLM GPU)

- 1) keep model weights loaded.
- 2) Wait for 128 inference requests.
- 3) Compute KV cache for all 128 requests and their prompts (tokens so far).

- But why might this algorithm be bad?

1: prioritizing throughput over latency

2: not all inference requests length 1024

*next time: continuous batching / scheduling