



CS 498: Machine Learning System Spring 2026

Minjia Zhang

The Grainger College of Engineering

DL Inference

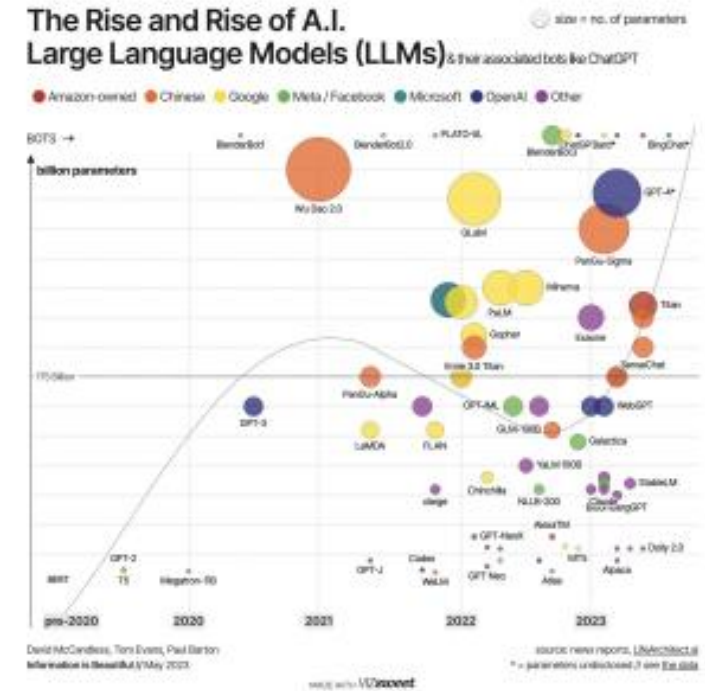
- Paged Attention

Objective: Learn why KV cache is important for LLM inference and how to manage KV cache efficiently

Serving LLMs is Expensive



- A ton of GPUs are required for production scale LLM services
- Nevertheless, each GPU only serve a handful of requests per second
 - For LLaMA-13B and moderate-size inputs, 1 A100 process < 1 request per second



Zuckerberg's Meta Is Spending Billions to Buy 350,000 Nvidia H100 GPUs

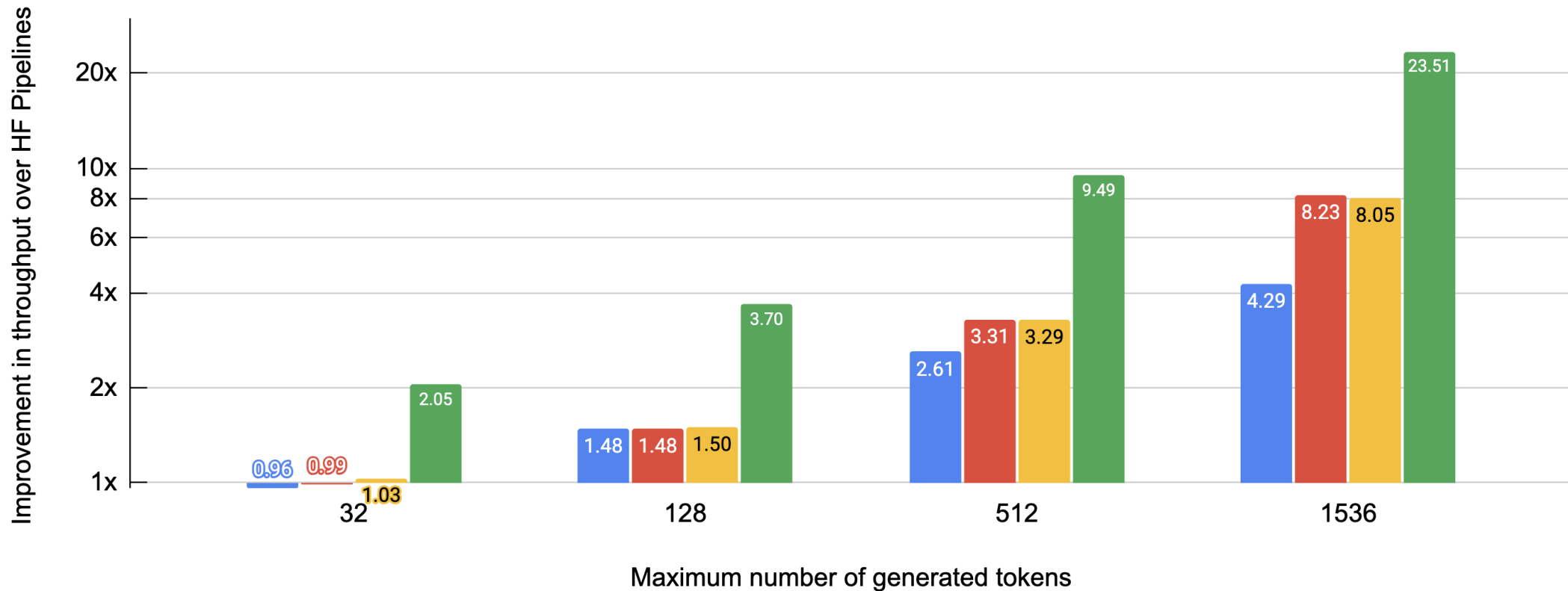
In total, Meta will have the compute power equivalent to 600,000 Nvidia H100 GPUs to help it develop next-generation AI, says CEO Mark Zuckerberg.

Throughput Improvement from Continuous Batching

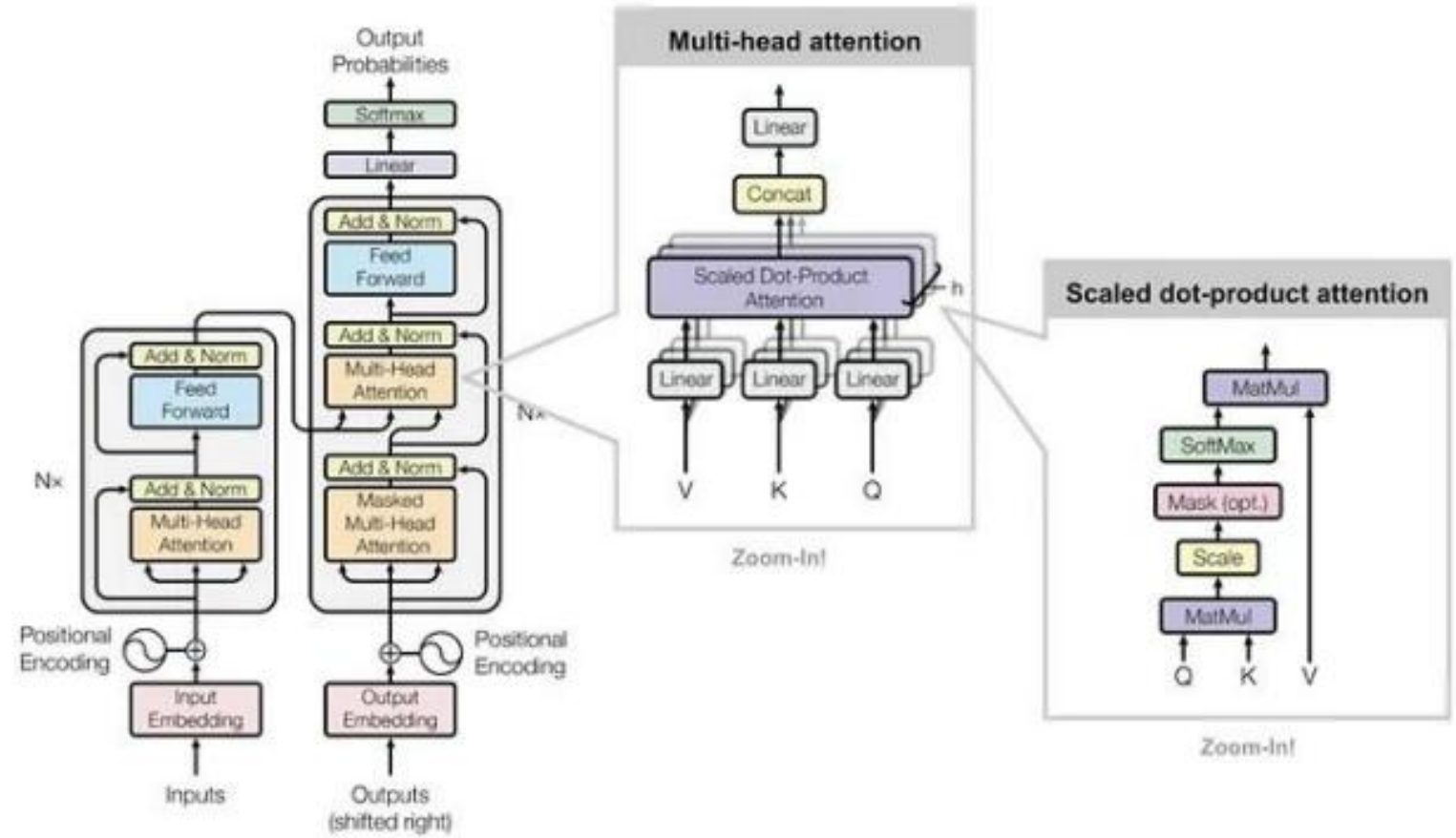


Throughput improvement over naive static batching vs. generated sequence length variance

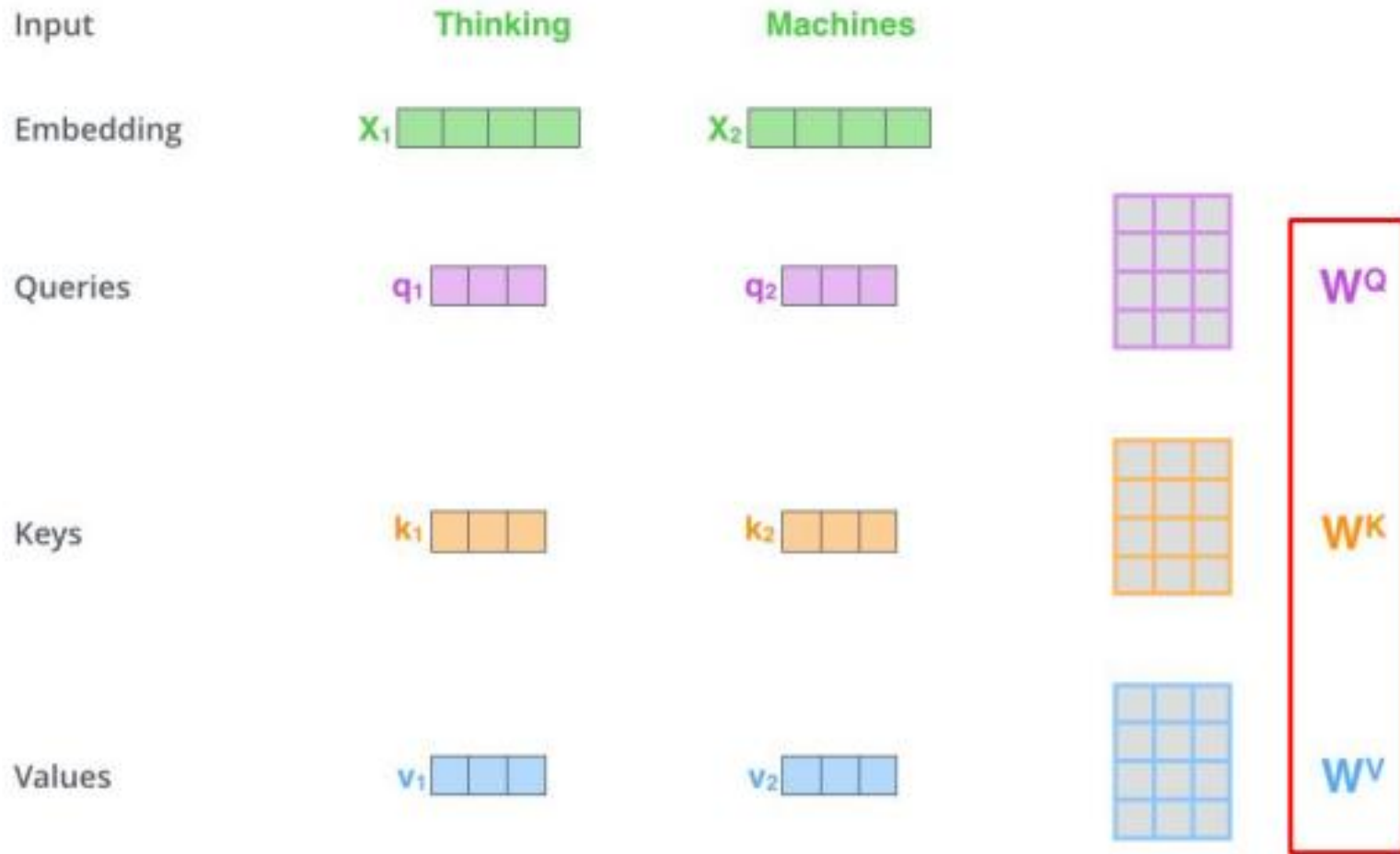
- Static batching (FasterTransformer)
- Continuous batching (text-generation-inference)
- Continuous batching (Ray Serve)
- Continuous batching (vLLM)



KV Cache



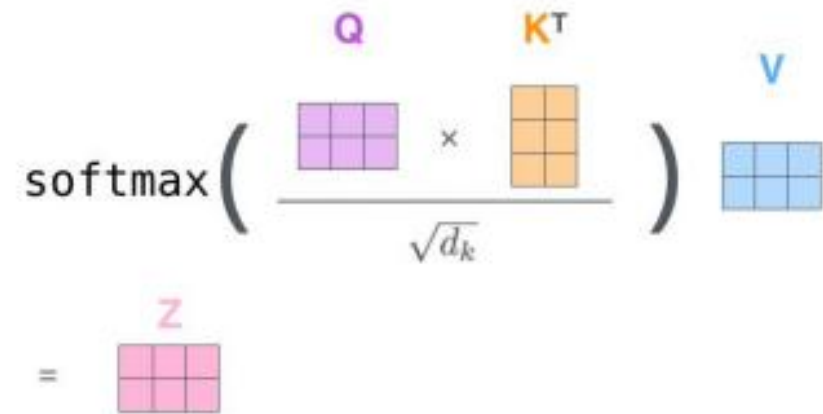
Self-Attention



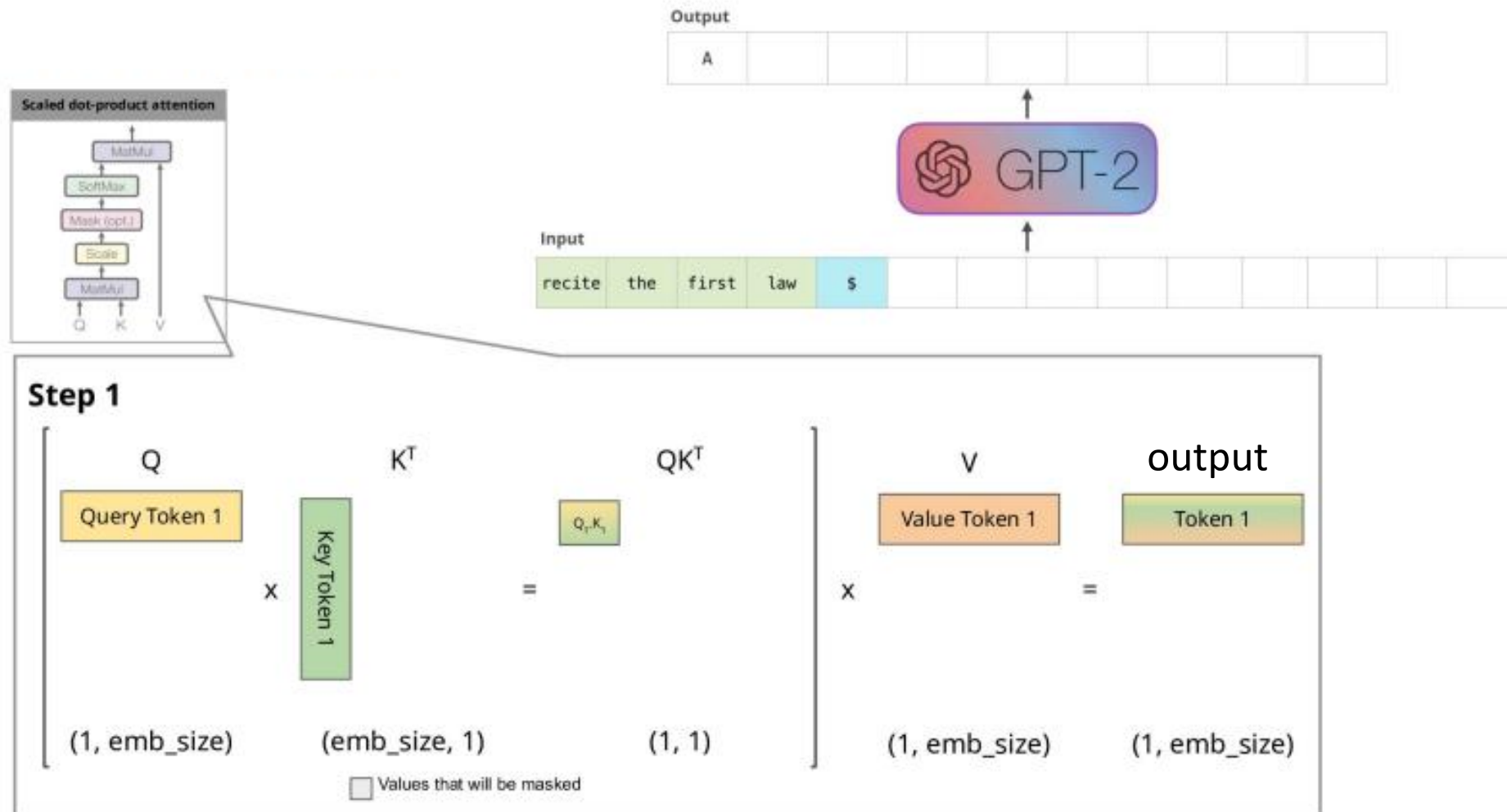
Self-Attention



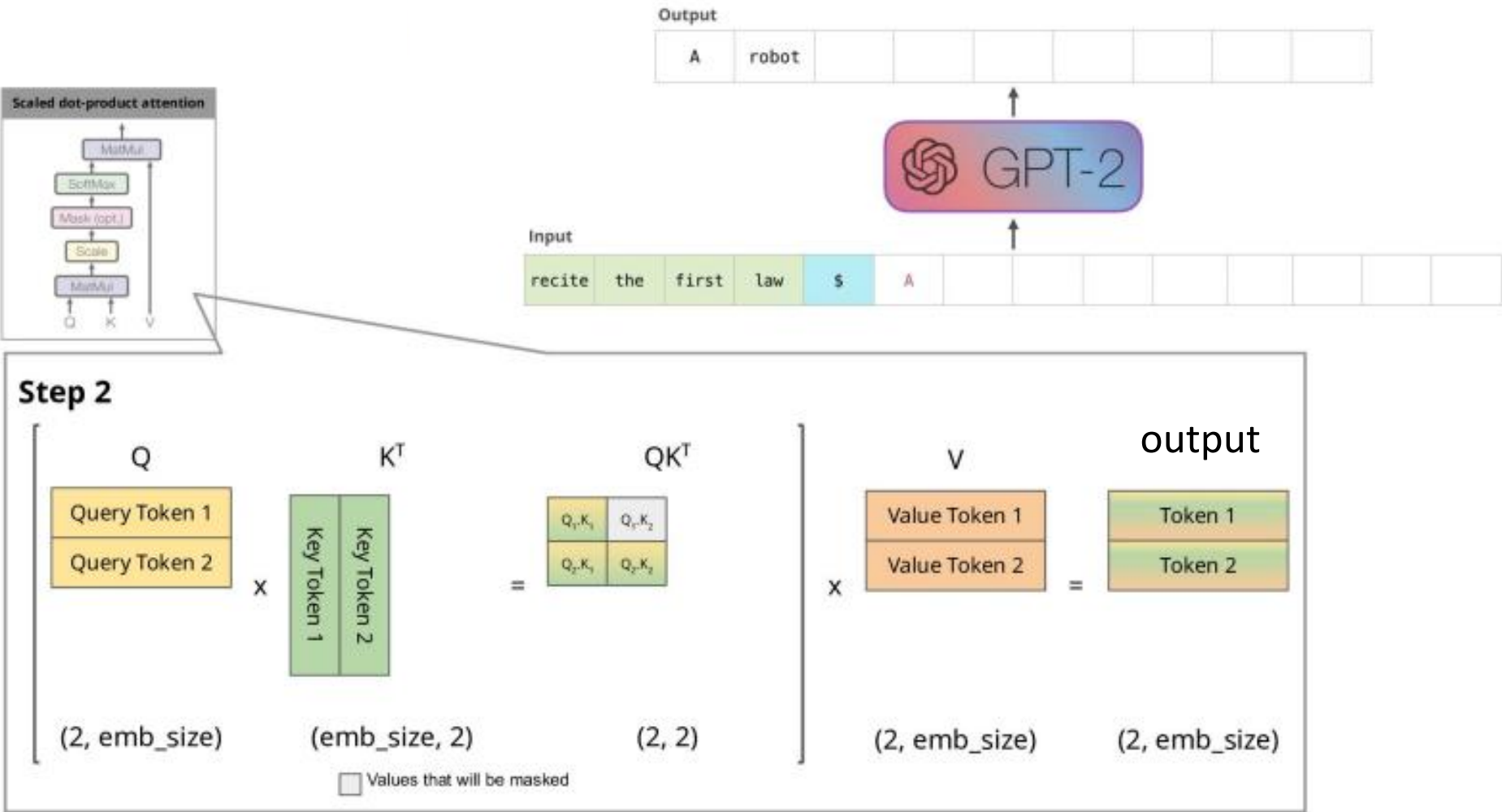
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V}$$



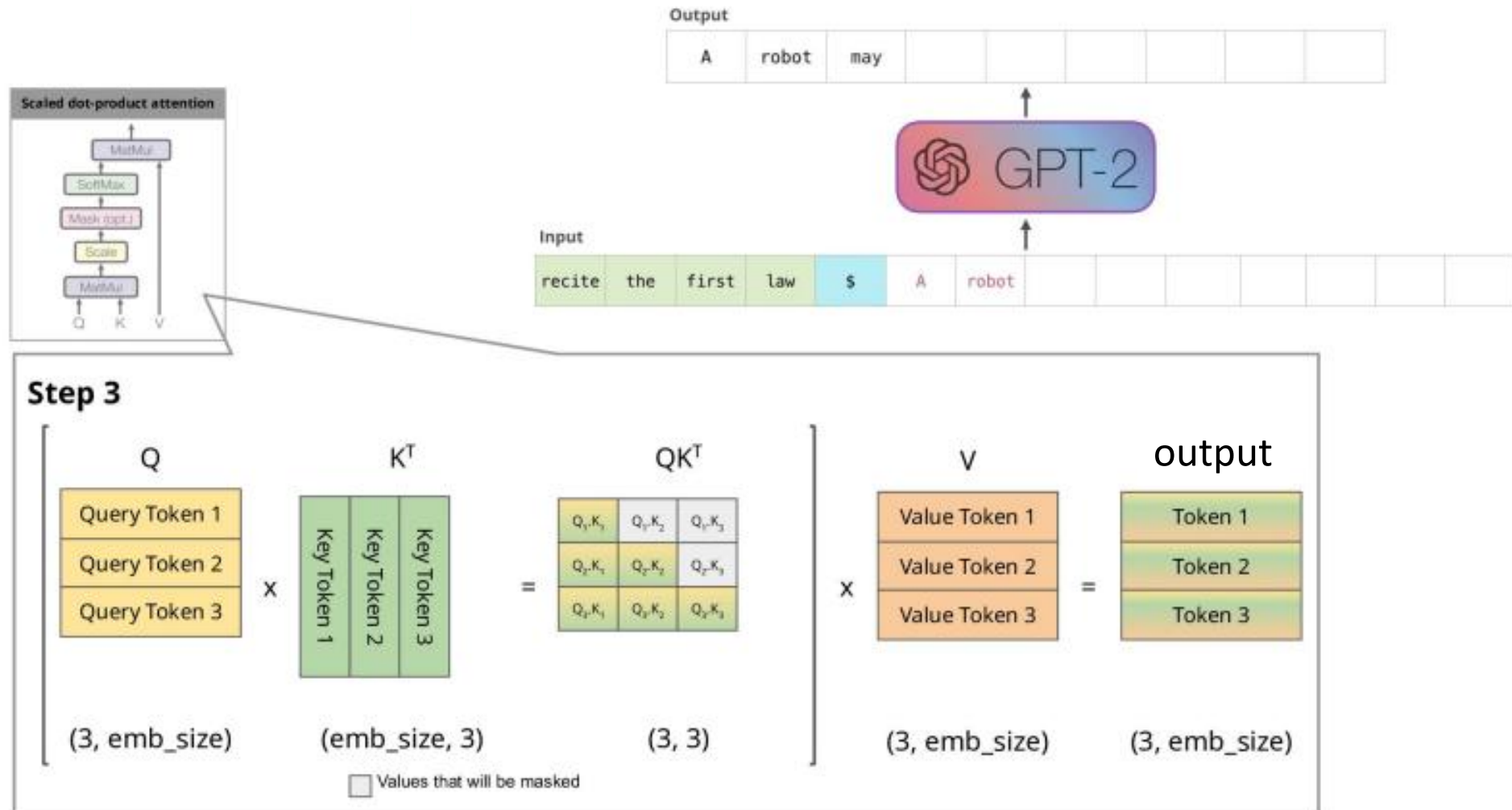
Naïve Self-Attention



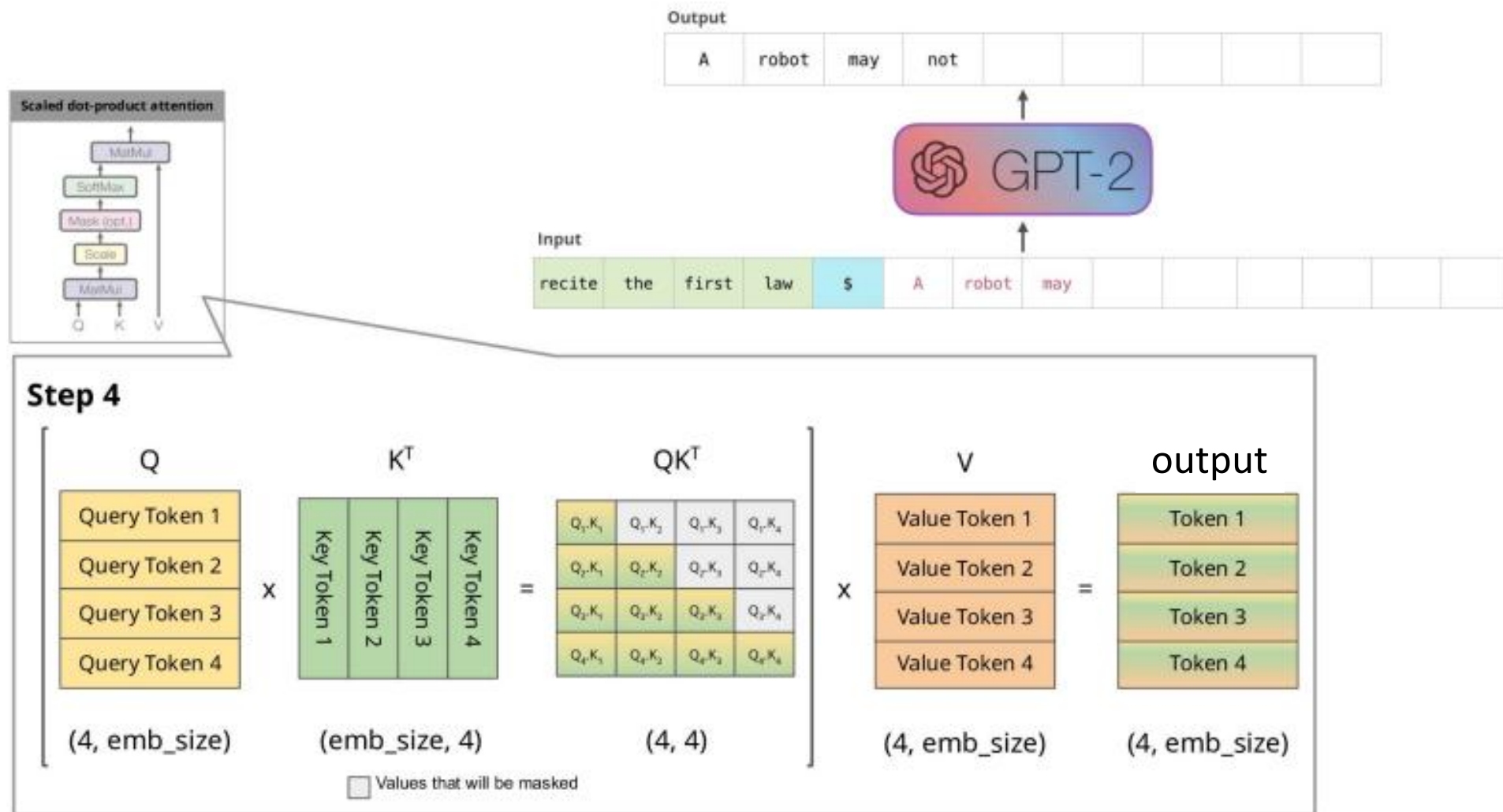
Naïve Self-Attention



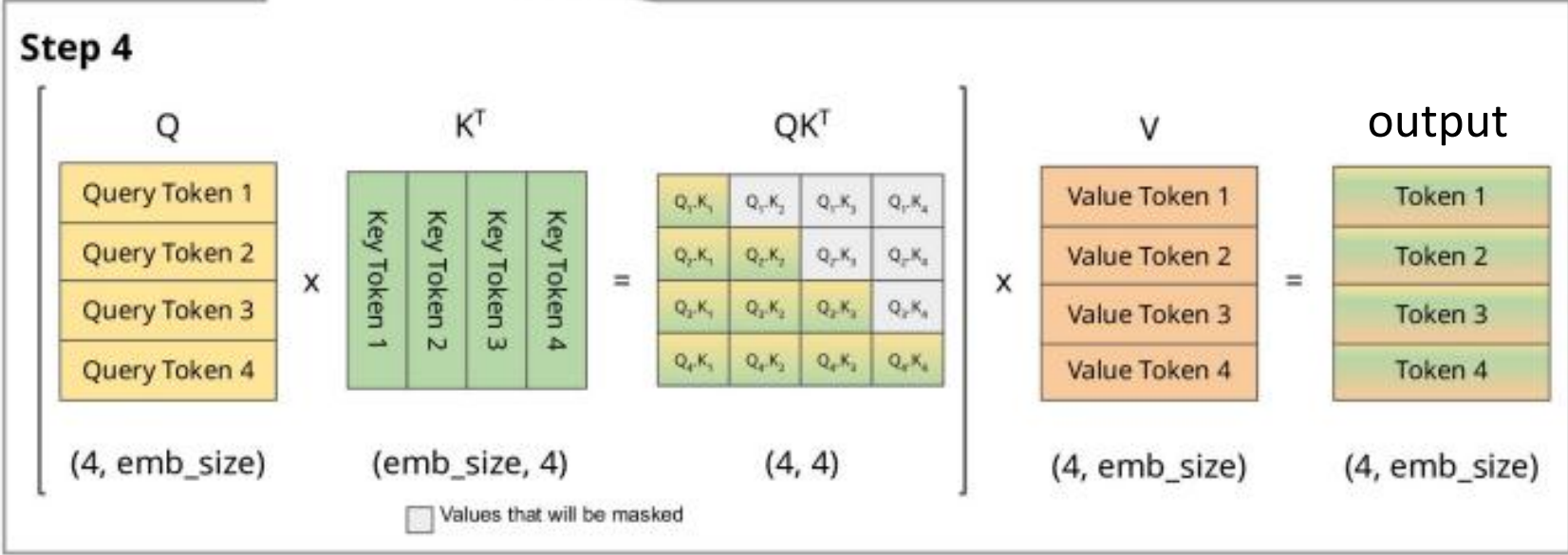
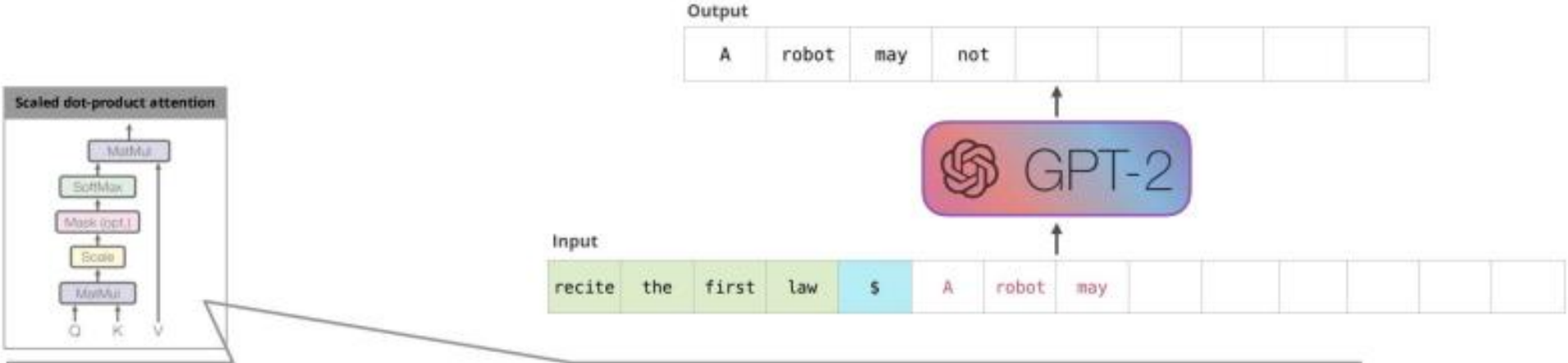
Naïve Self-Attention



Naïve Self-Attention

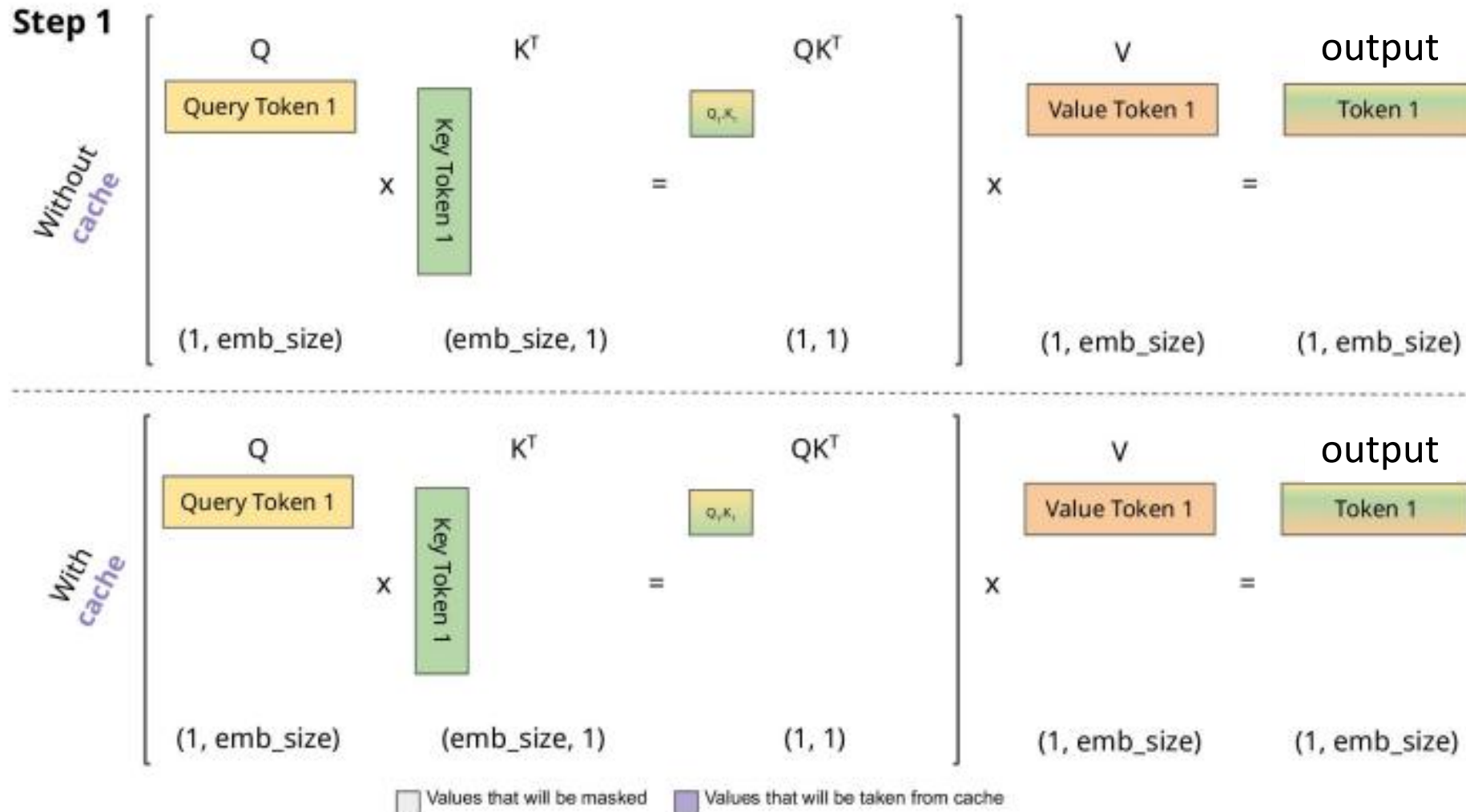


Naïve Self-Attention

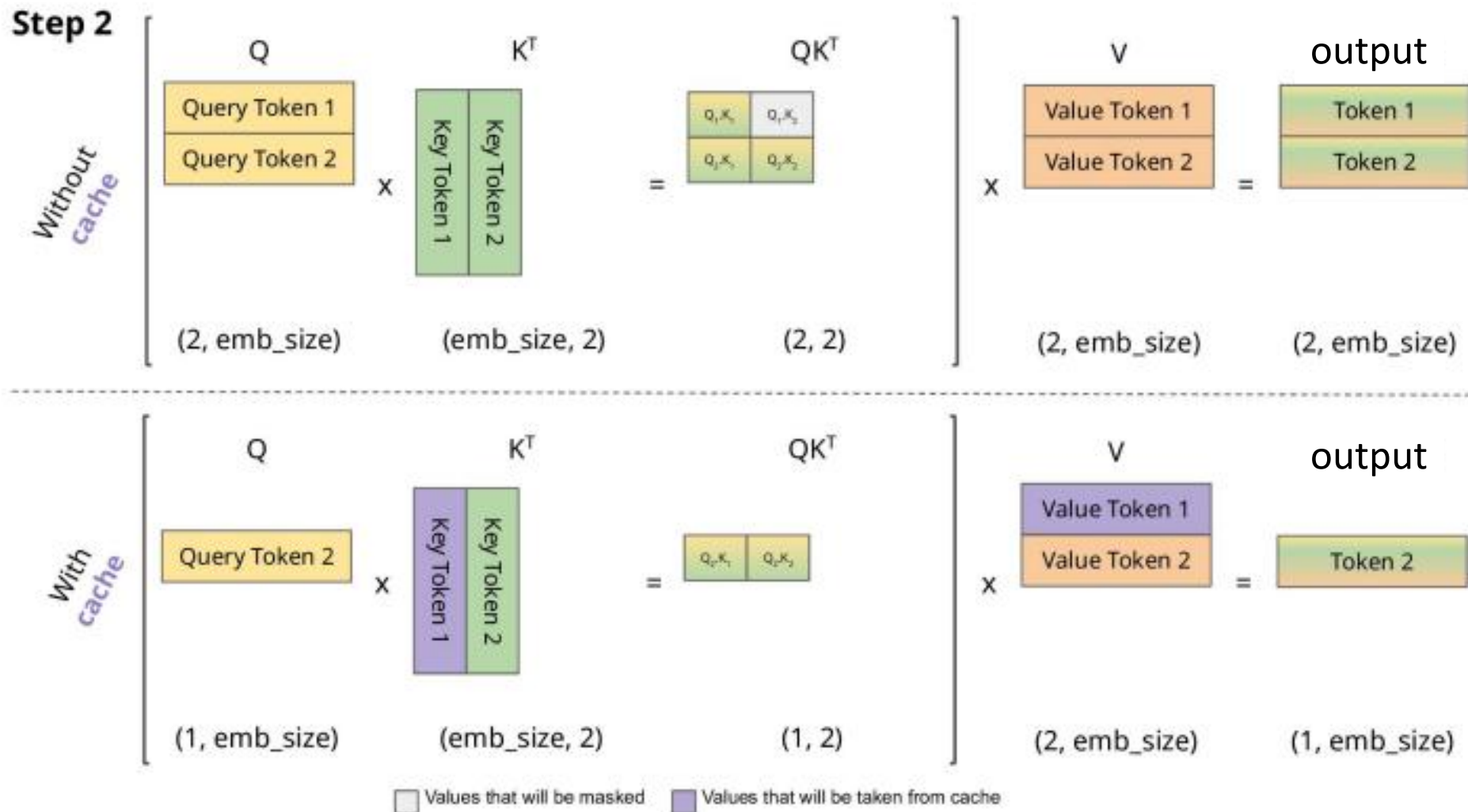


Question: Issue with naive attention?

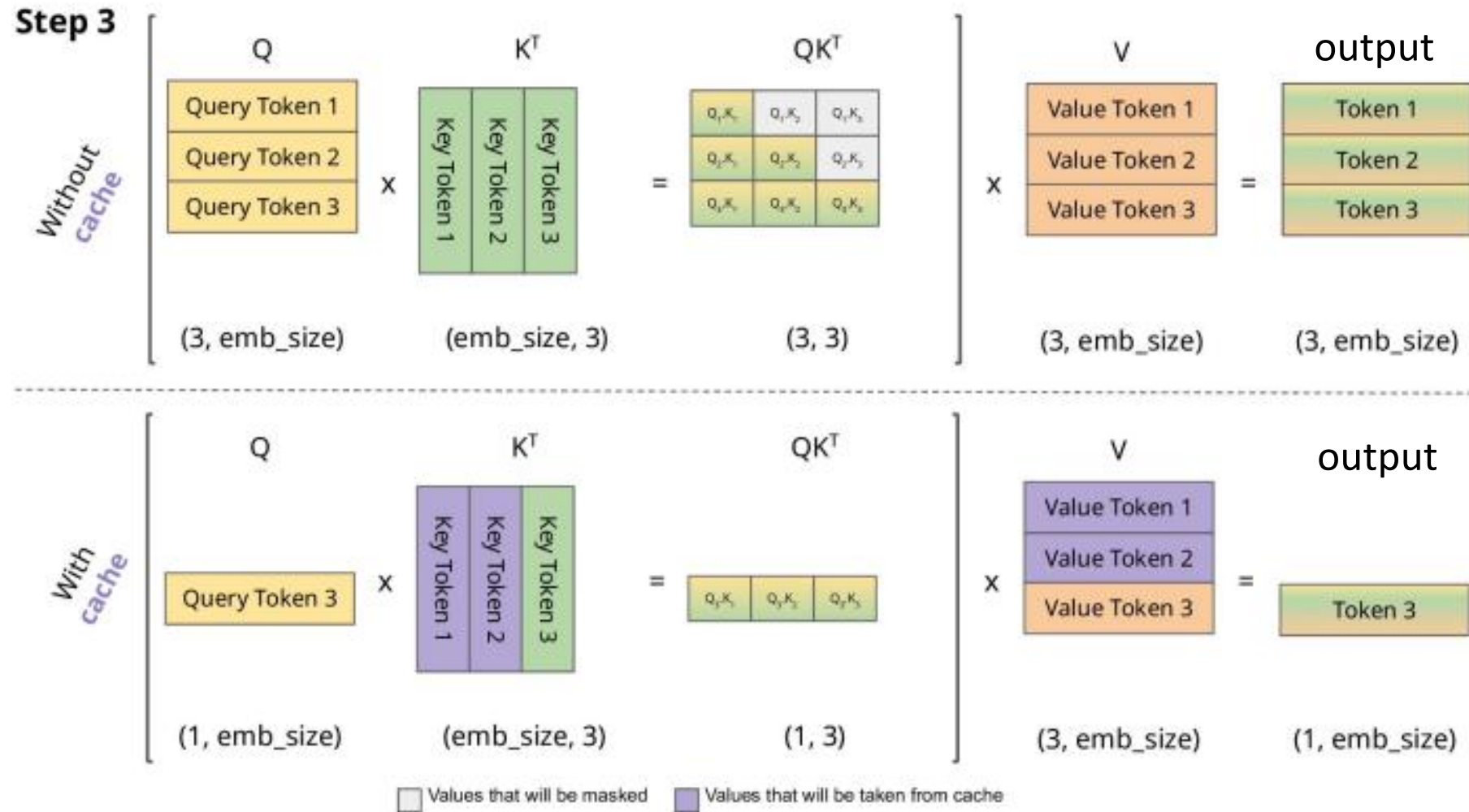
Self-Attention with KV Cache



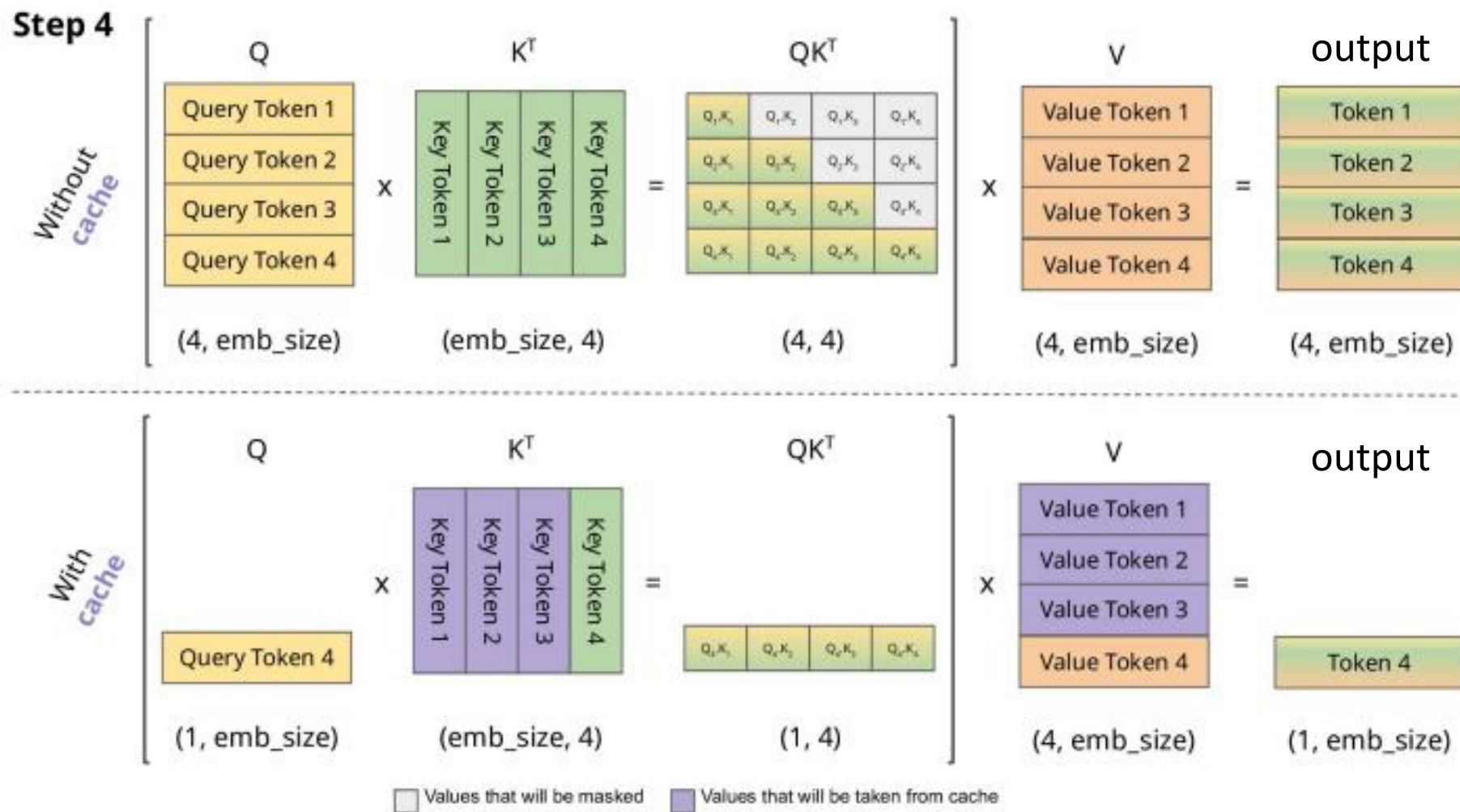
Self-Attention with KV Cache



Self-Attention with KV Cache



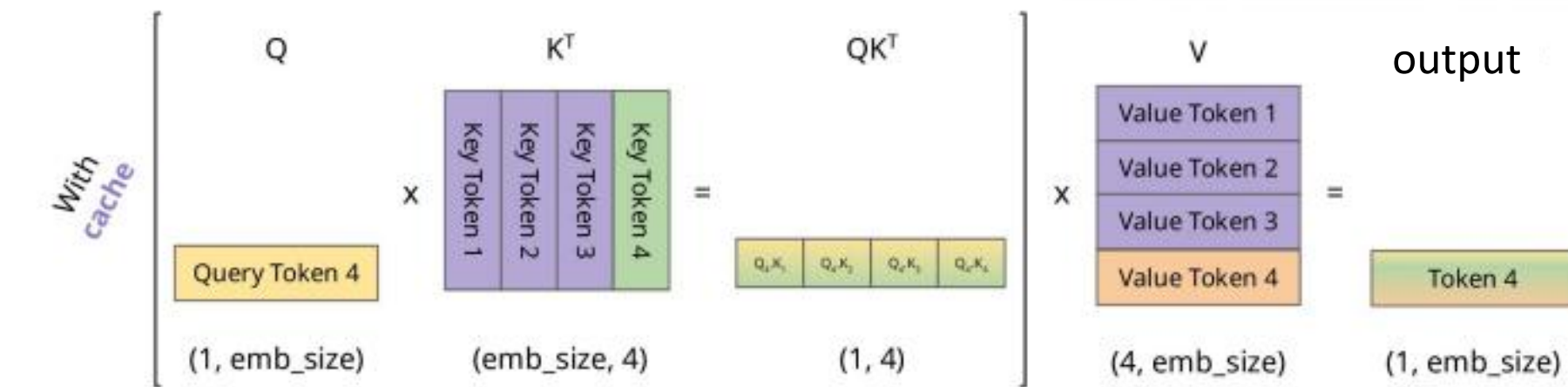
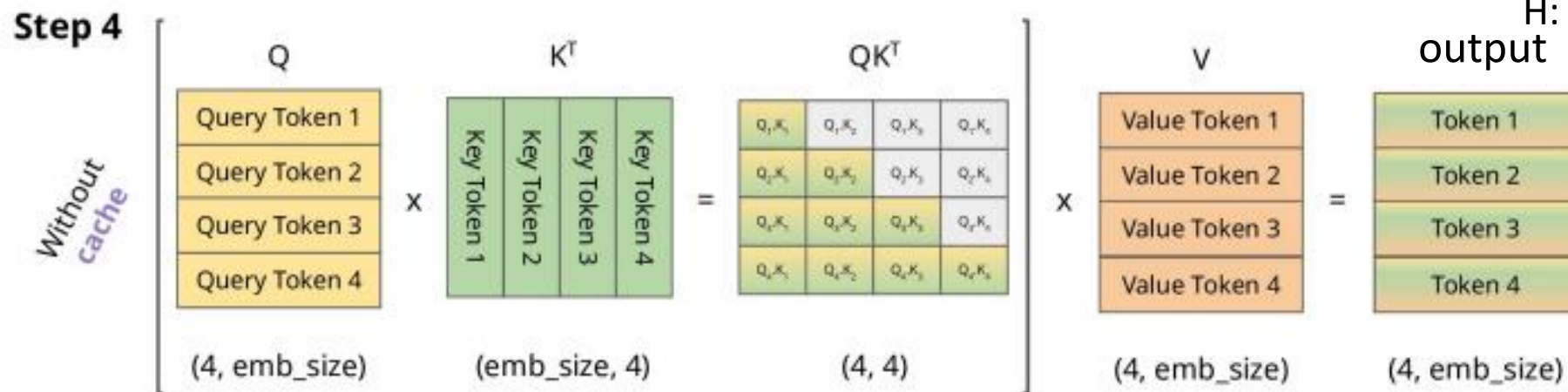
Self-Attention with KV Cache



Self-Attention with KV Cache



S: sequence length
H: model embedding size



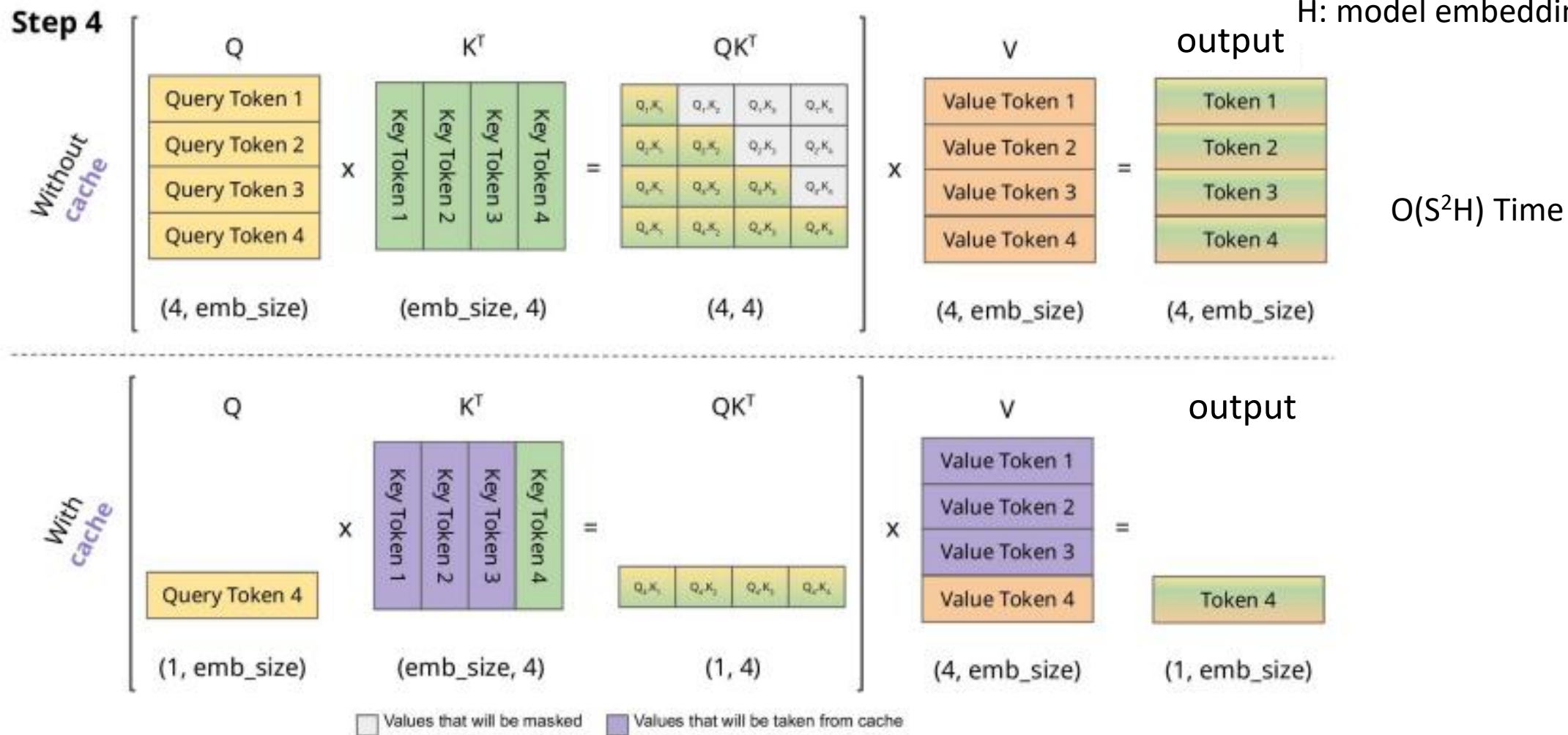
□ Values that will be masked ■ Values that will be taken from cache

Question: Time and memory complexity without and with cache?

Self-Attention with KV Cache



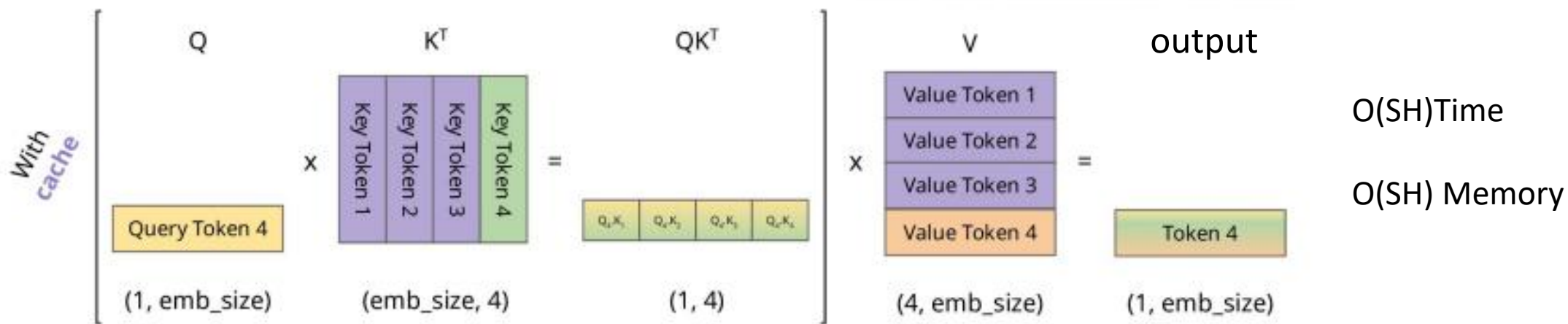
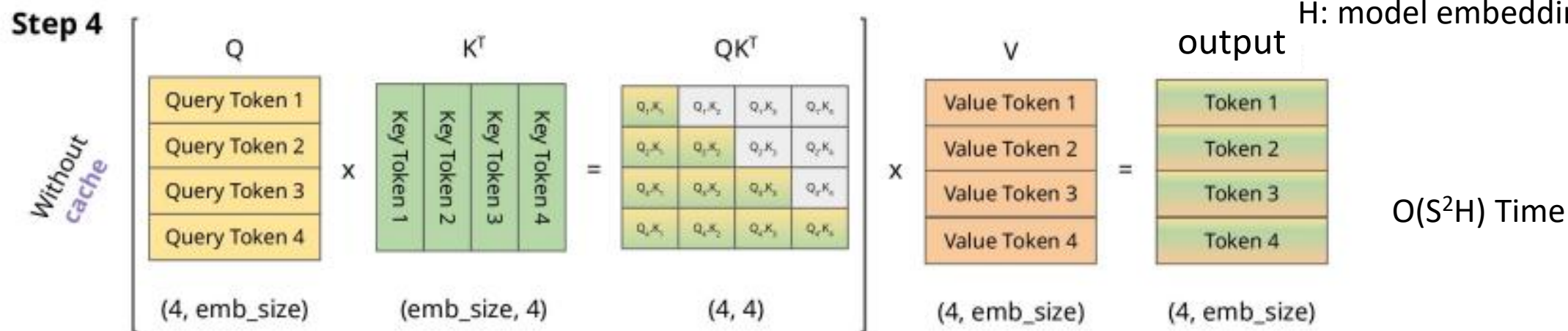
S: sequence length
H: model embedding size



Self-Attention with KV Cache

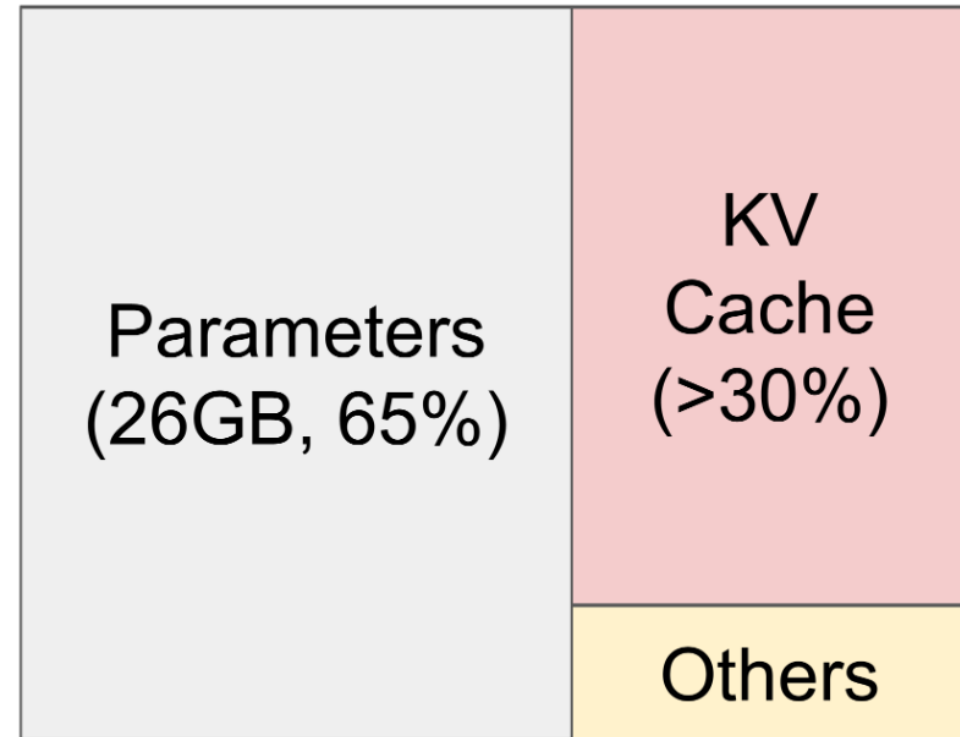


S: sequence length
H: model embedding size



Values that will be masked
 Values that will be taken from cache

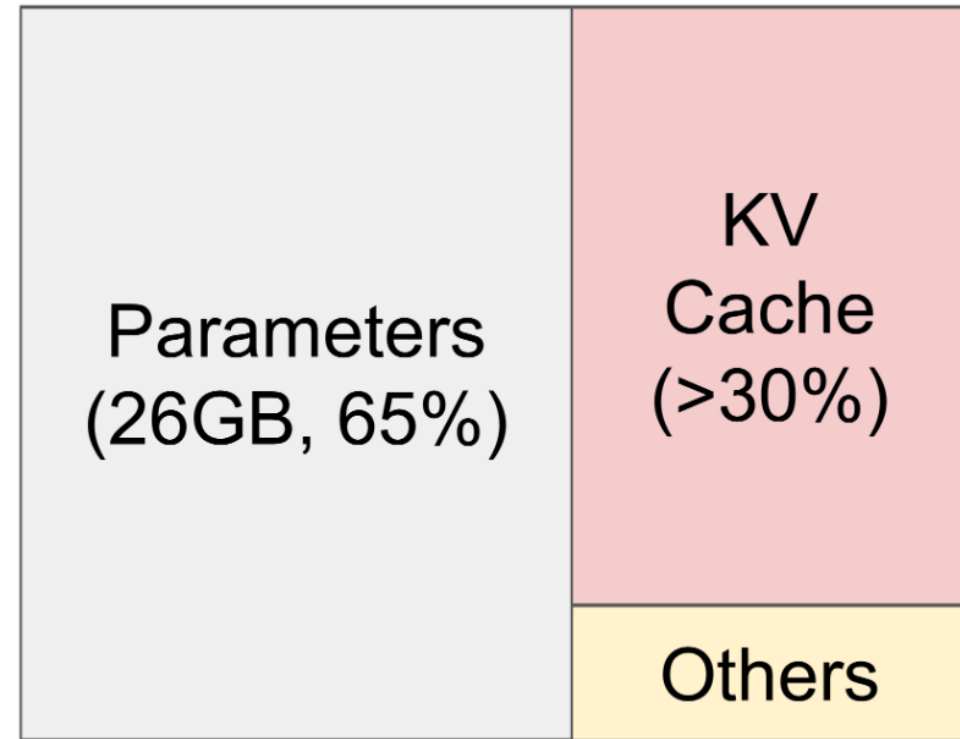
- KV cache
 - Sequence length
 - # Hidden
 - # Layers
 - Batch size



LLaMA-13B on A100 40 GB

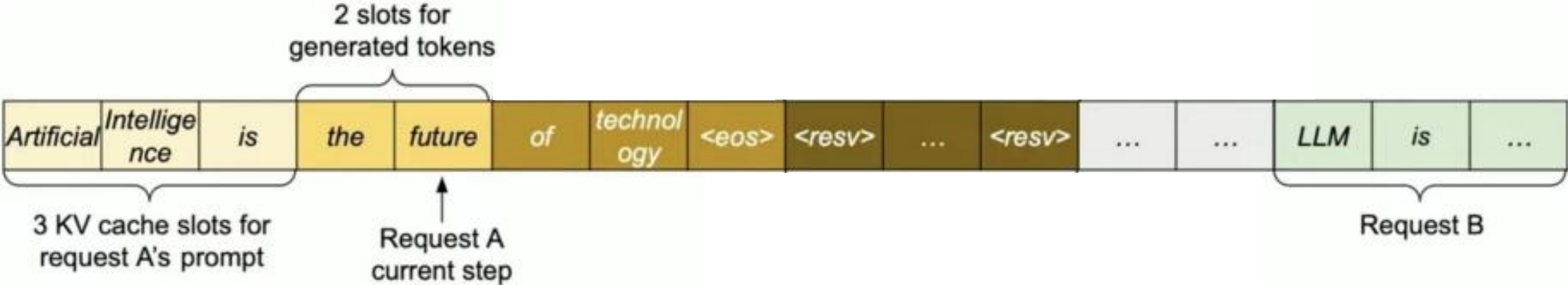
- KV cache
 - Sequence length
 - # Hidden
 - # Layers
 - Batch size

How about we allocate some contiguous memory per request:
 $\text{max-seq-len} * \text{hidden-dim} * \text{num-layers}$?

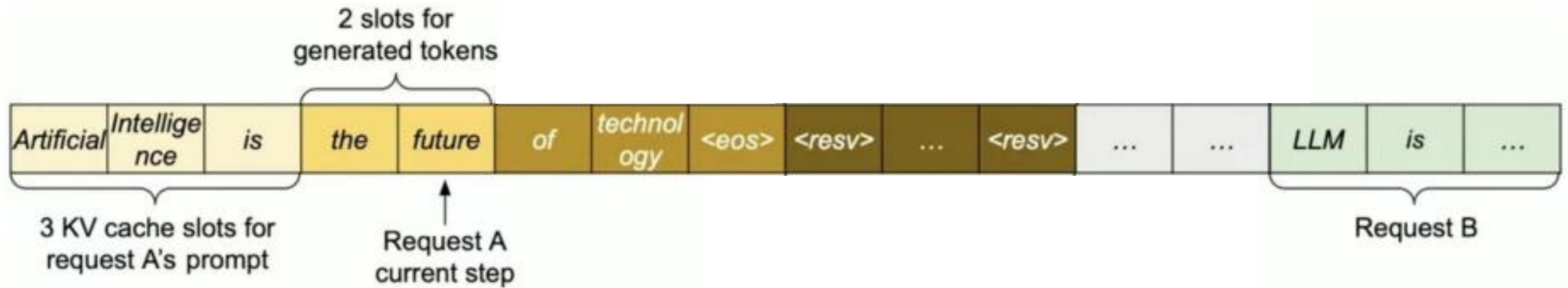


LLaMA-13B on A100 40 GB

Motivation: Memory Waste in KV Cache

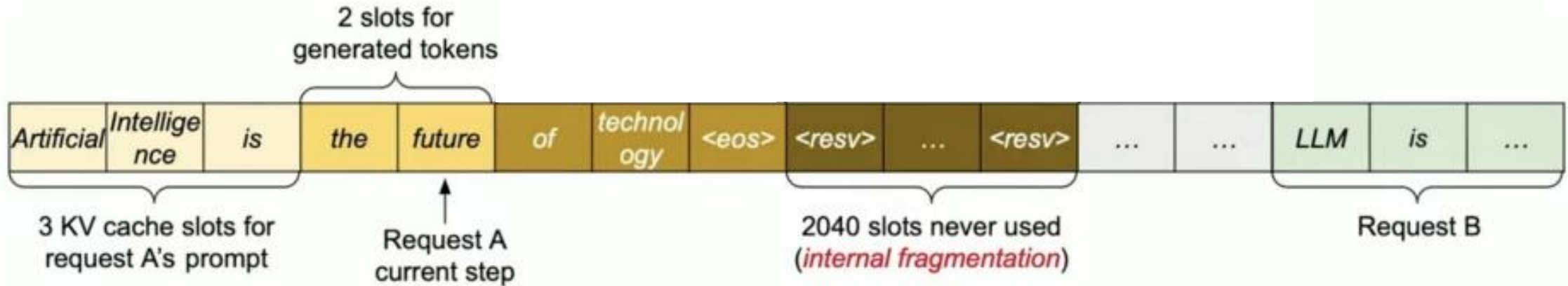


Motivation: Memory Waste in KV Cache



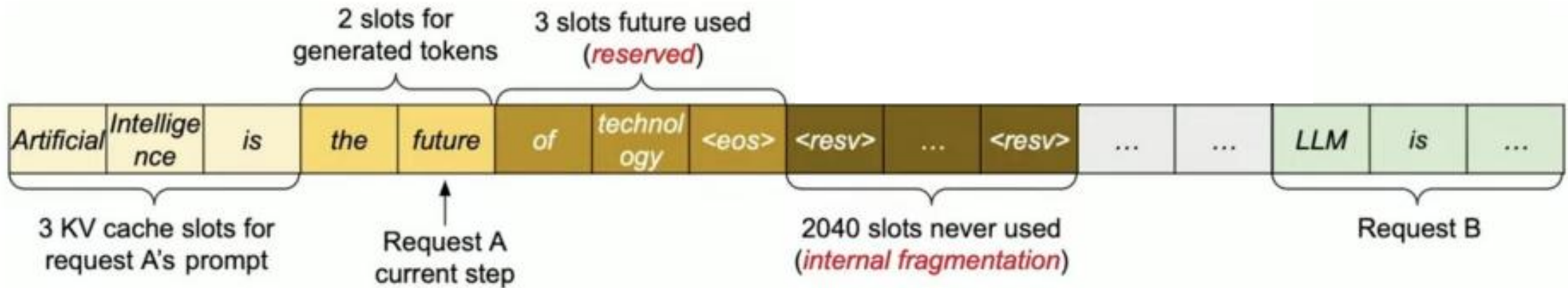
Question: Where is memory being wasted?

Motivation: Memory Waste in KV Cache



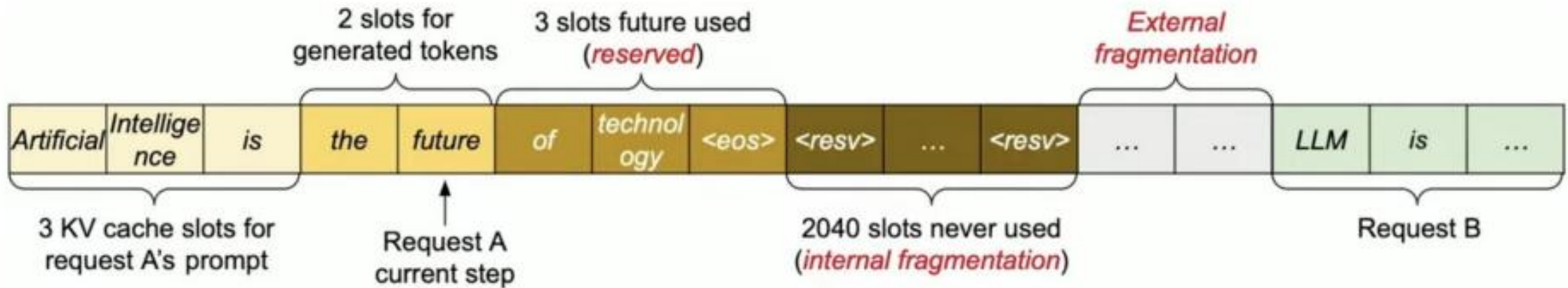
- **Internal fragmentation:** memory reserved for a request but left unused because the final output length is shorter than expected

Motivation: Memory Waste in KV Cache



- **Internal fragmentation:** memory reserved for a request but left unused because the final output length is shorter than expected
- **Reservation:** memory currently unused but reserved for future use by the same request

Motivation: Memory Waste in KV Cache

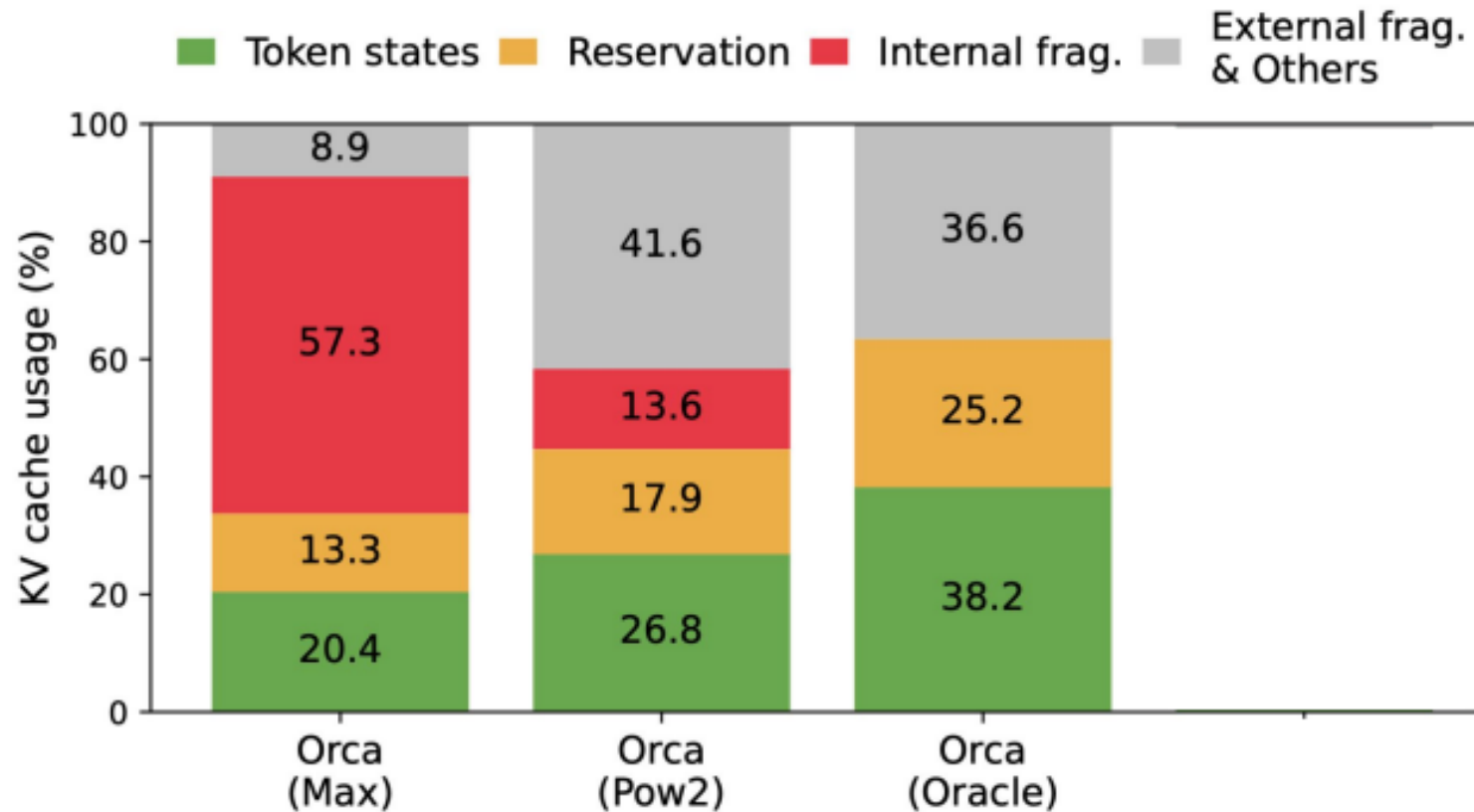


- **Internal fragmentation:** memory reserved for a request but left unused because the final output length is shorter than expected
- **Reservation:** memory currently unused but reserved for future use by the same request
- **External fragmentation:** small unused memory gaps between allocations caused by varying sequence lengths across different requests

Motivation: Memory Waste in KV Cache



- Only **20-40%** of KV cache is used to store the actual token states



- PagedAttention algorithm allows for storing attention key-value pairs in **non-contiguous blocks** in the memory



Key and value vectors

Block 1	years	ago	our	fathers
Block 2	brought	forth		
Block 0	Four	score	and	seven

block size = 4

vLLM: Efficient Memory Management using PagedAttention



- PagedAttention algorithm allows for storing attention key-value pairs in **non-contiguous blocks** in the memory
- Each KV block is a fixed-size contiguous chunk of memory that stores KV cache from left to right

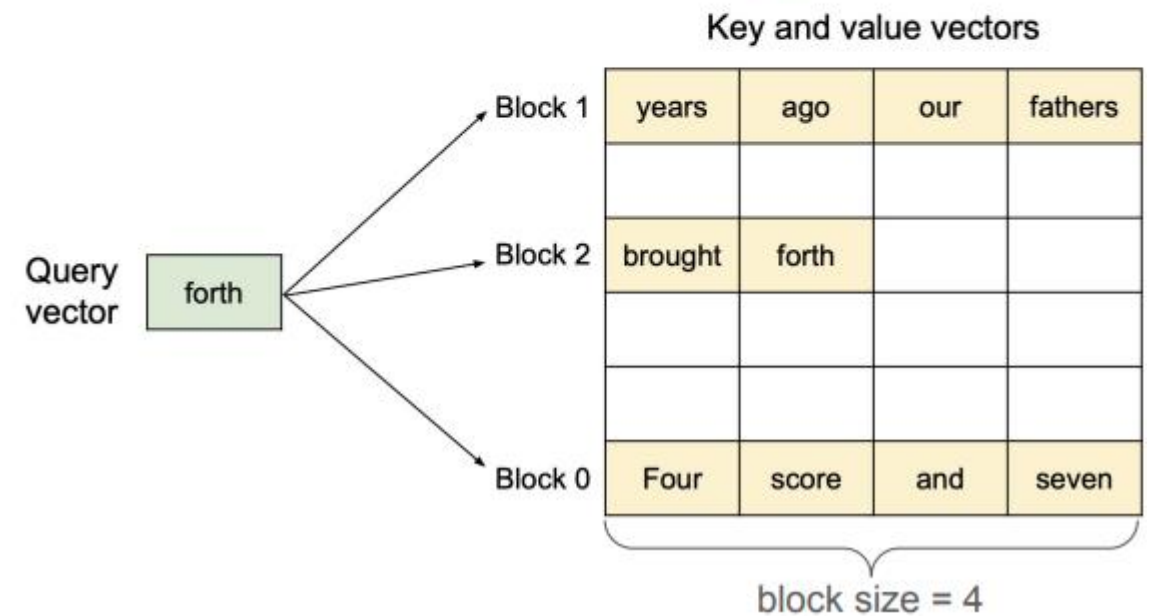


Key and value vectors

Block 1	years	ago	our	fathers
Block 2	brought	forth		
Block 0	Four	score	and	seven

block size = 4

- PagedAttention algorithm allows for storing attention key-value pairs in **non-contiguous blocks** in the memory
- Each KV block is a fixed-size contiguous chunk of memory that stores KV cache from left to right
- Attention for each new token is computed across all tokens in the non-contiguous KV block



Block Table For Managing KV Cache Blocks



Request A

Prompt: "Four score and seven years ago our"
Outputs: "fathers" → "brought" → ...

Logical KV blocks

Block 0	① Four	① score	① and	① seven
Block 1	① years	① ago	① our	② fathers
Block 2	③ brought			
Block 3				

Block Table

Physical block number	# filled
① 7	① 4
① 1	① 3 → 4 ②
③ 3	③ 1
-	-

Physical KV blocks
(on GPU DRAM)

Block 0				
Block 1	① years	① ago	① our	② fathers
Block 2				
Block 3	③ brought			
Block 4				
Block 5				
Block 6				
Block 7	① Four	① score	① and	① seven
Block 8				

Memory Efficiency of PagedAttention



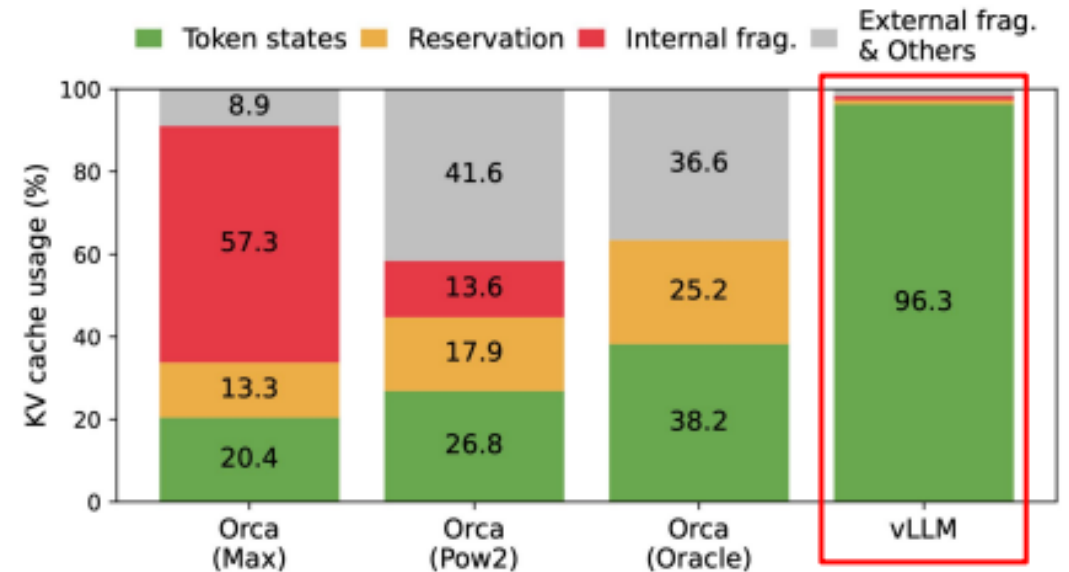
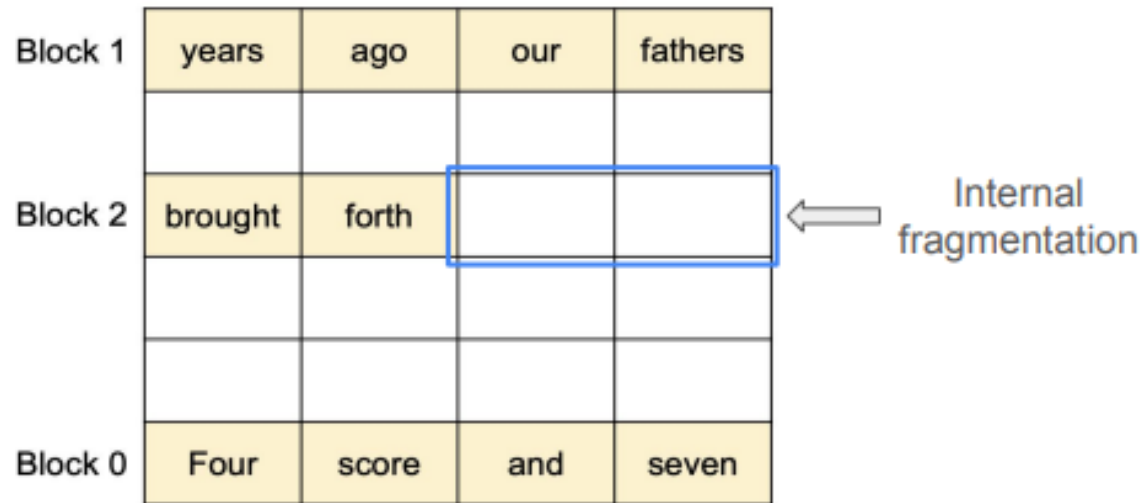
- Minimal internal fragmentation
 - # of wasted tokens per sequence < block size



Memory Efficiency of PagedAttention



- Minimal internal fragmentation
 - # of wasted tokens per sequence < block size
- No external fragmentation

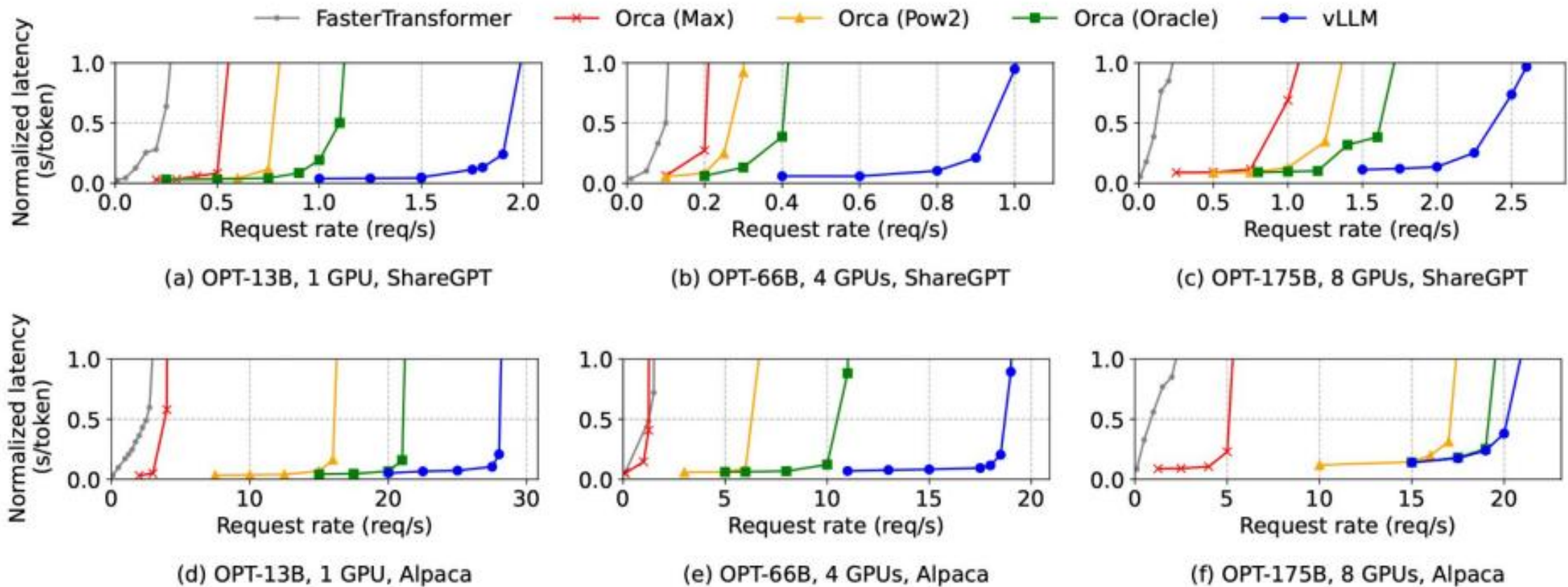


- With PagedAttention, wasted KV cache space is < 4% (3-5x improved memory utilization)

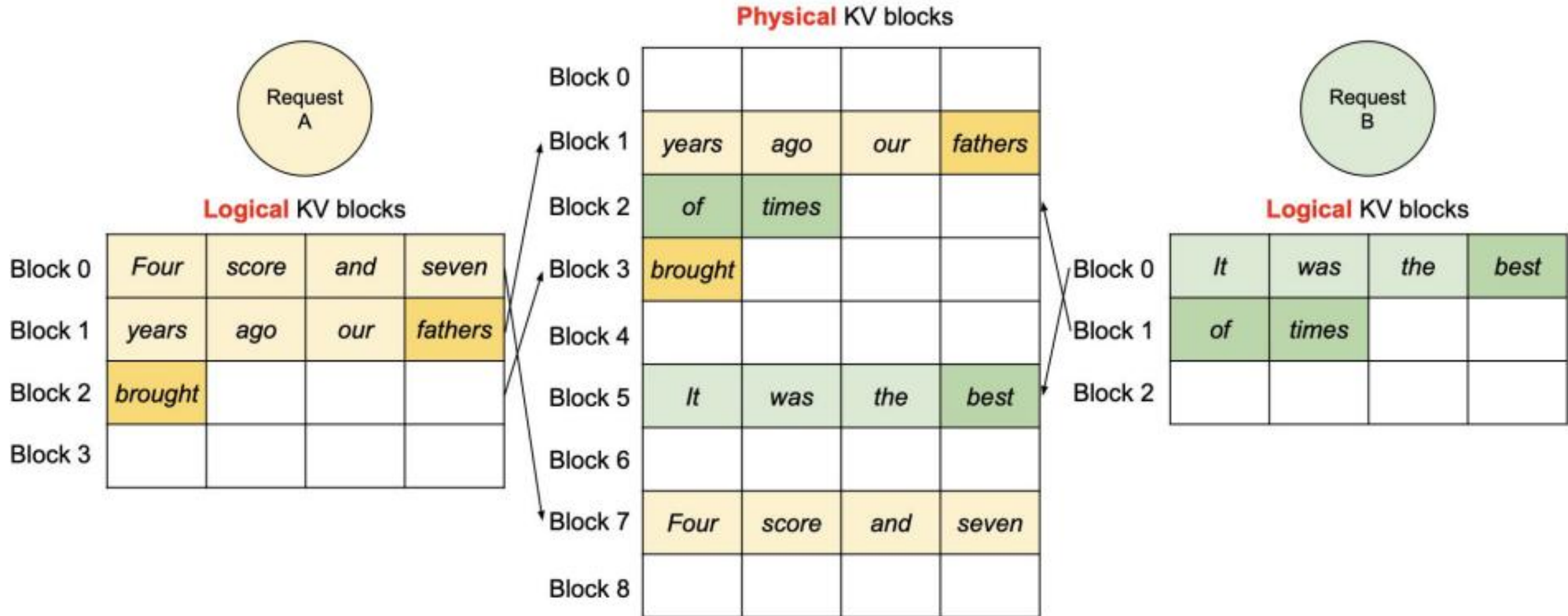
vLLM Performance with Basic Generation



- one sample per request on three models and two datasets



Handling Multiple Requests

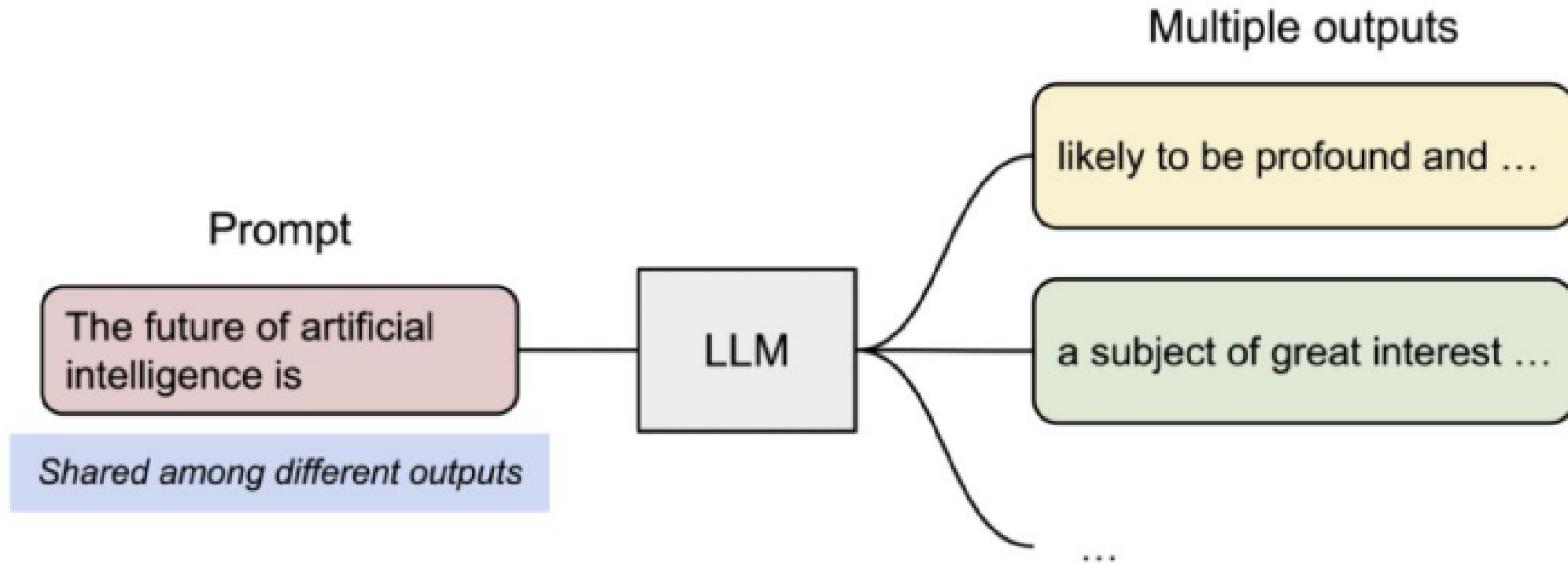


E.g.

- Parallel sampling

- Beam search

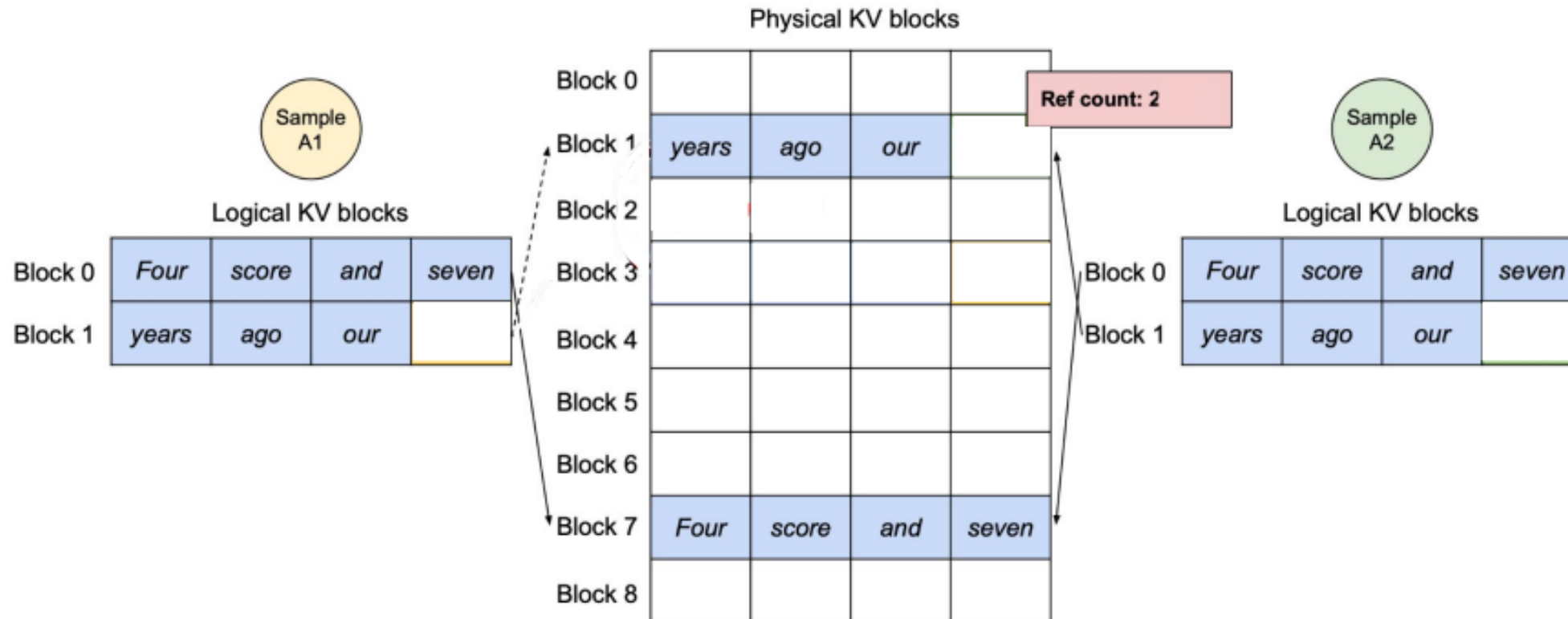
...



Parallel Sampling with PagedAttention



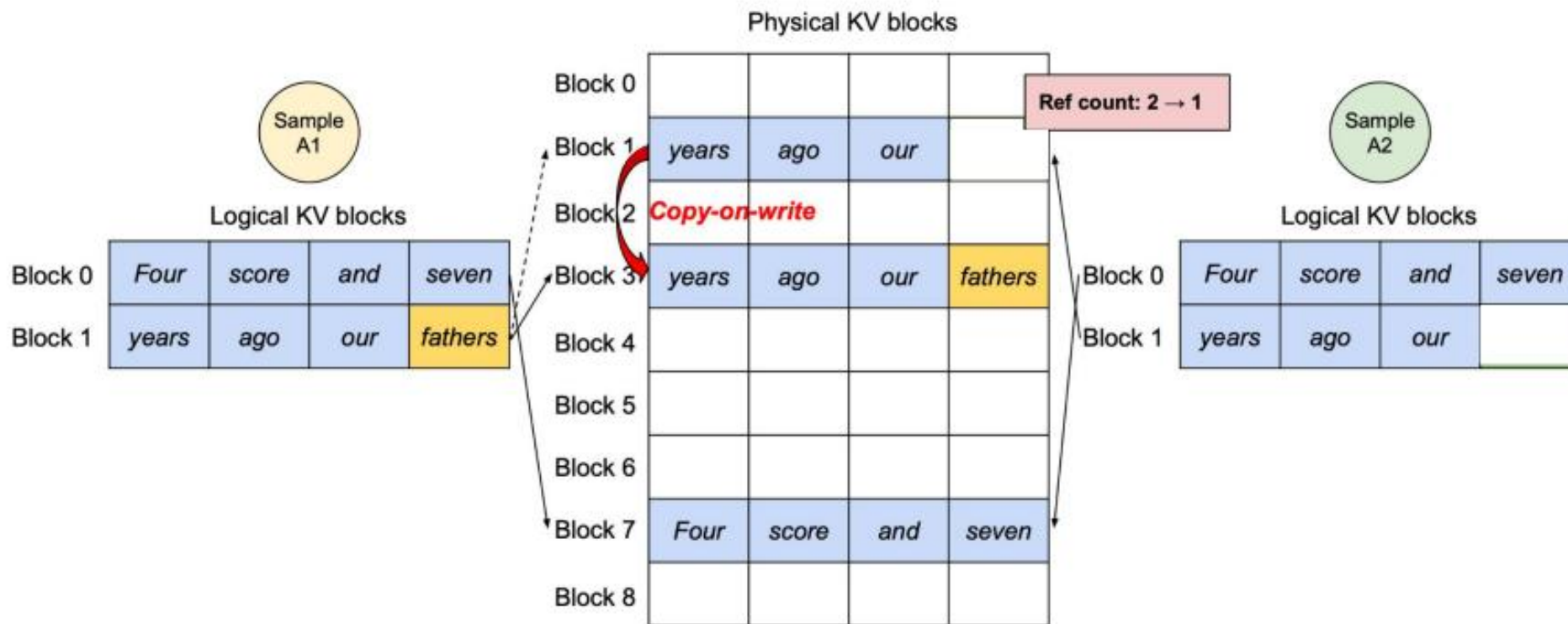
- Dynamic block mapping enables sharing of KV blocks
- Reference count keeps track of the number of logical KV blocks that are mapped to each physical KV block



Parallel Sampling with PagedAttention



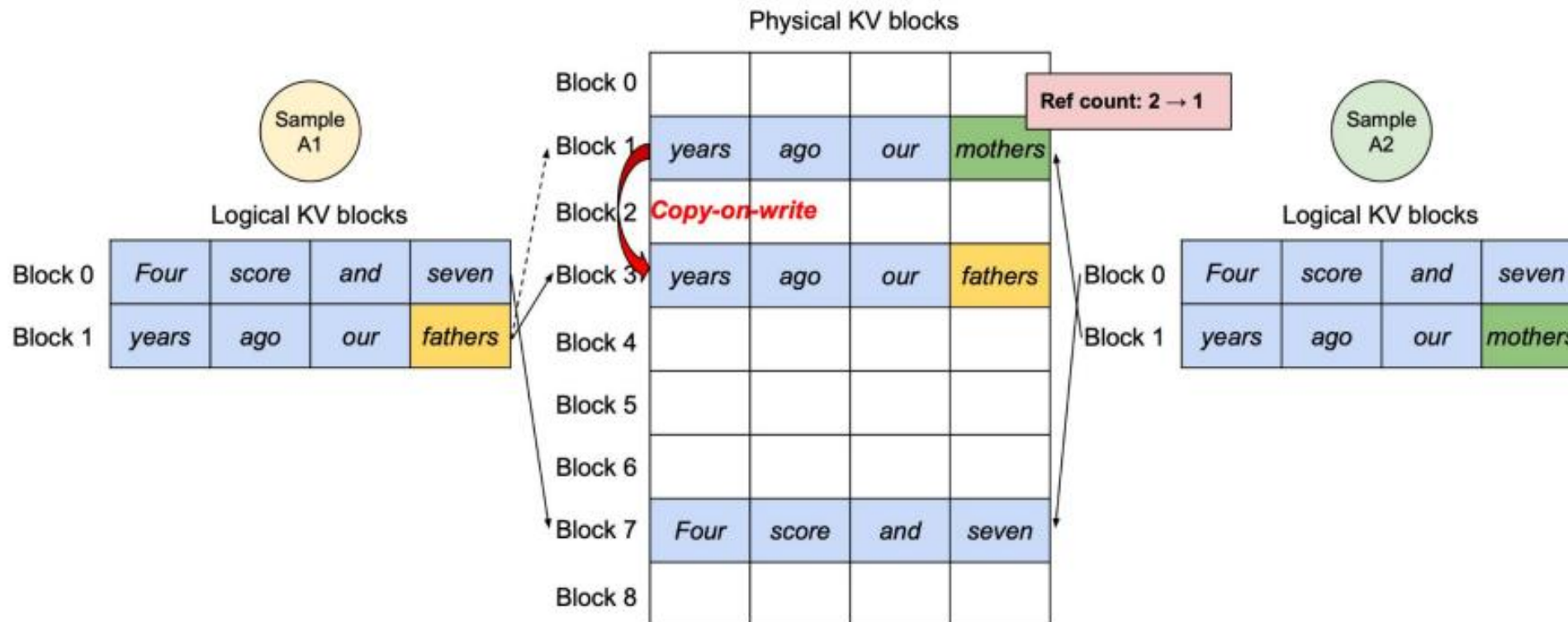
- At the generation phase, copy-on-write copies info to a newly allocated physical KV block when reference count > 1



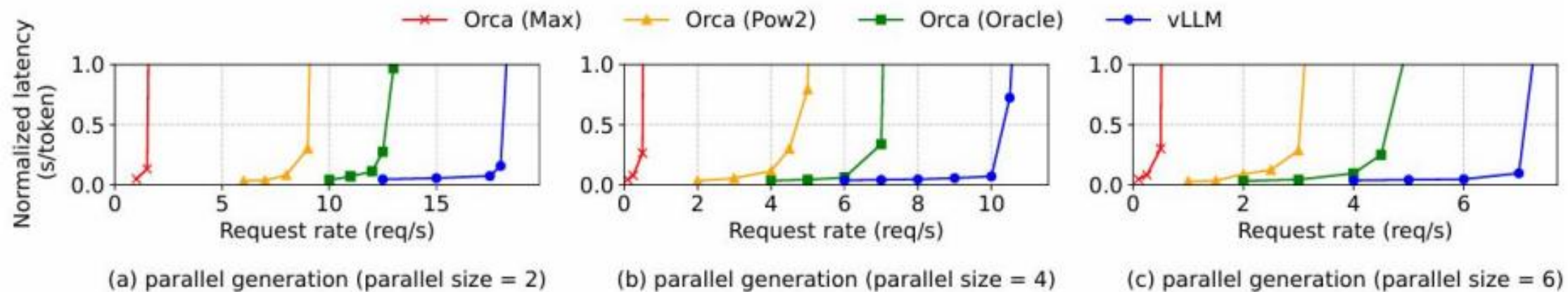
Parallel Sampling with PagedAttention



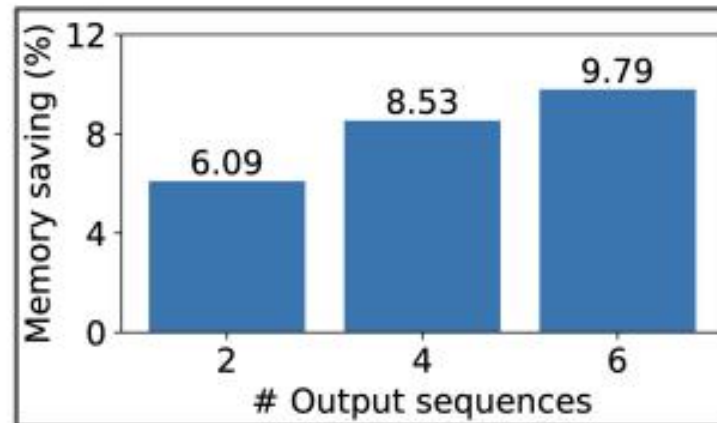
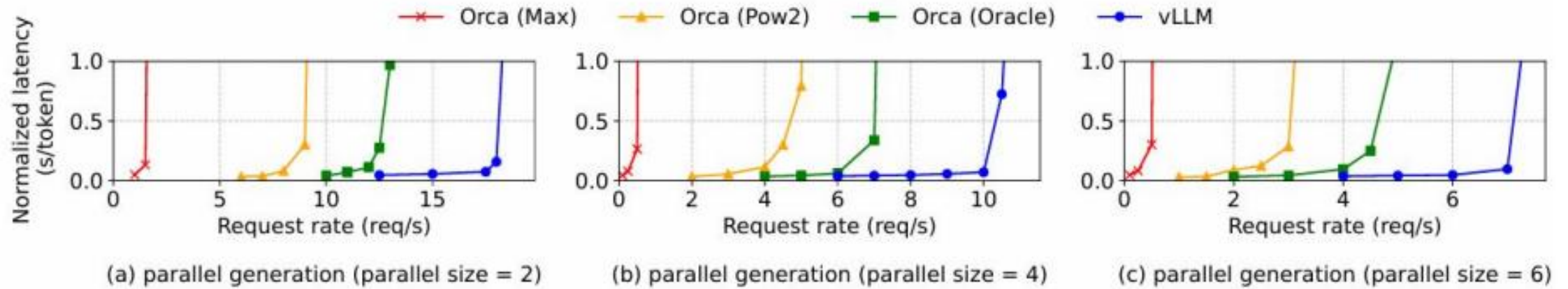
- At the generation phase, copy-on-write copies info to a newly allocated physical KV block when reference count > 1



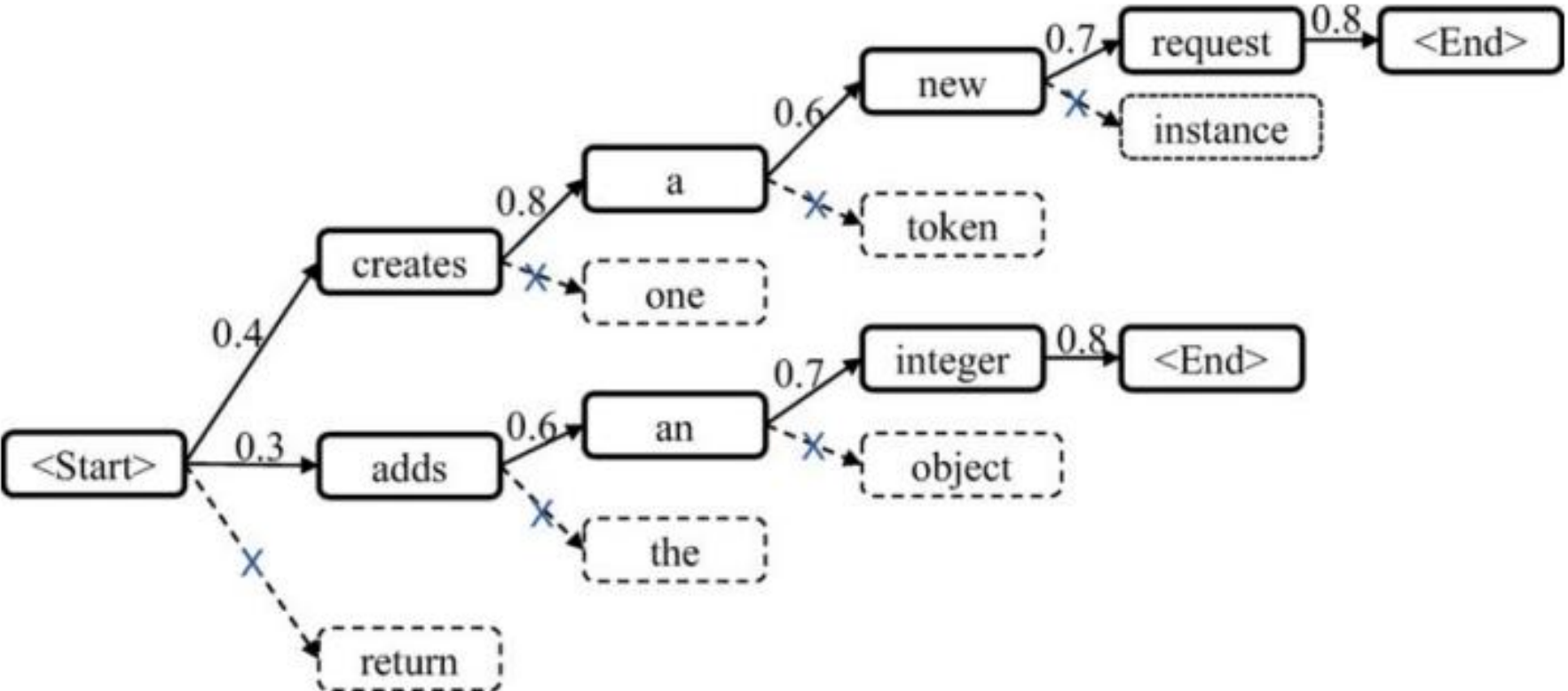
Parallel Sampling with PagedAttention



Parallel Sampling with PagedAttention



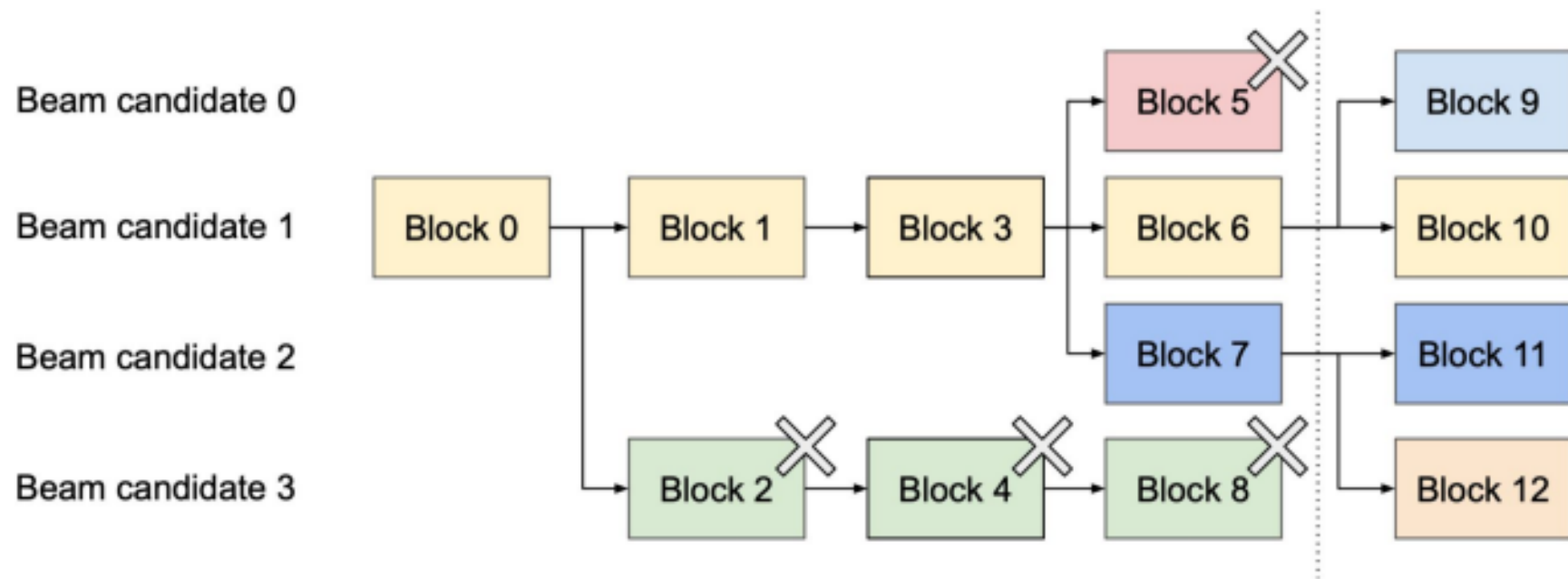
Beam Search



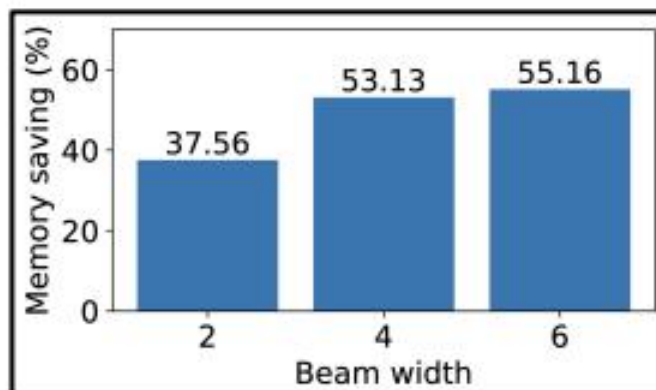
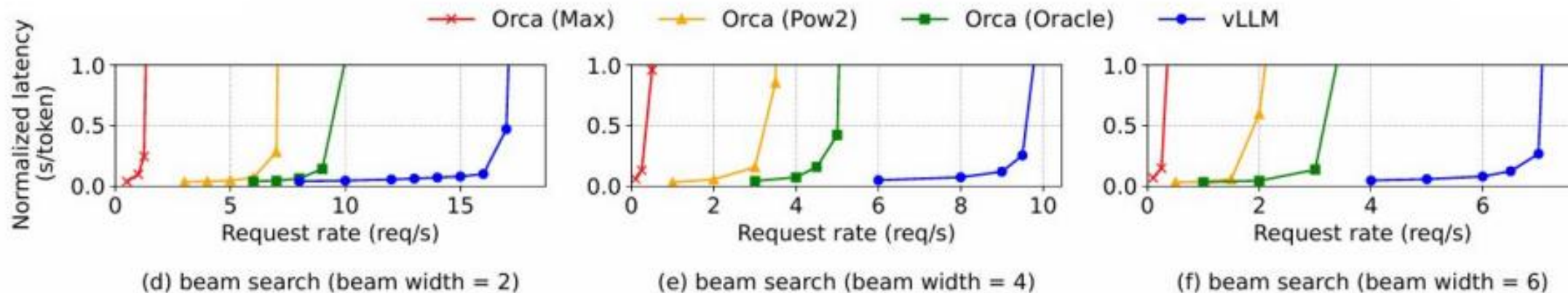
Beam Search with PagedAttention



Efficiently supported by dynamic block mapping and copy-on-write mechanism



Beam Search with PagedAttention



Limitations?



- Reduces memory fragmentation with paging
- Reduces memory usage with KV block sharing
- Enables batching of more requests, increasing the throughput of LLM inference

Questions?