



# Efficient Memory Management for Large Language Model Serving with PagedAttention

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng,  
Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, Ion Stoica

Presented by Yuqi Xue

Feb 15, 2024

# The Era of LLMs

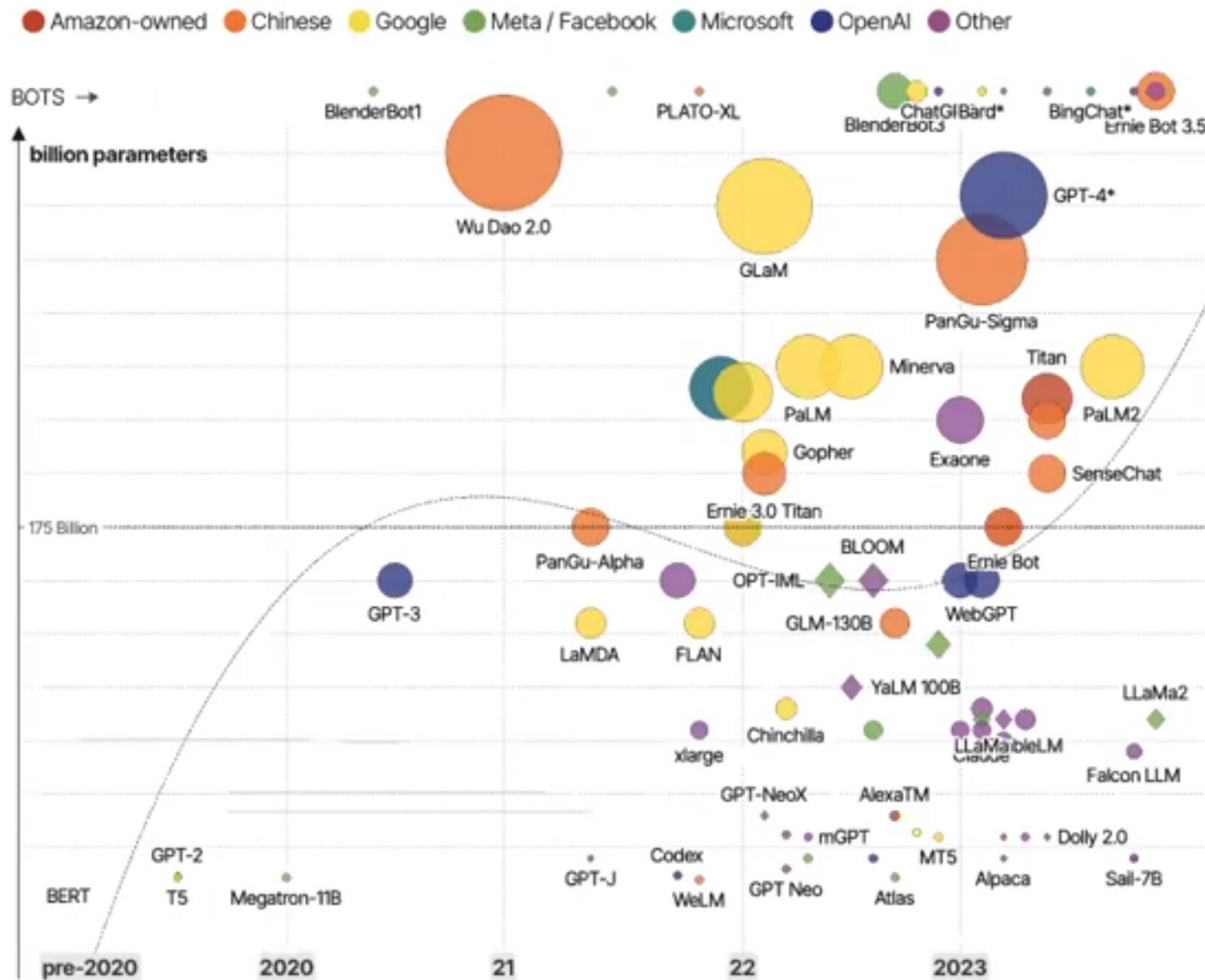


**ChatGPT**



**GitHub Copilot**

## The Rise and Rise of A.I. Large Language Models (LLMs) & their associated bots like ChatGPT



# LLM-Powered Services

Programming




GitHub Copilot




tabnine




Dev tools



Debuild




warp




cogram


Chat



ChatGPT



MessageBird



Sapling

Copywriting



cohere



copy.ai



HyperWrite

BizOps



viable

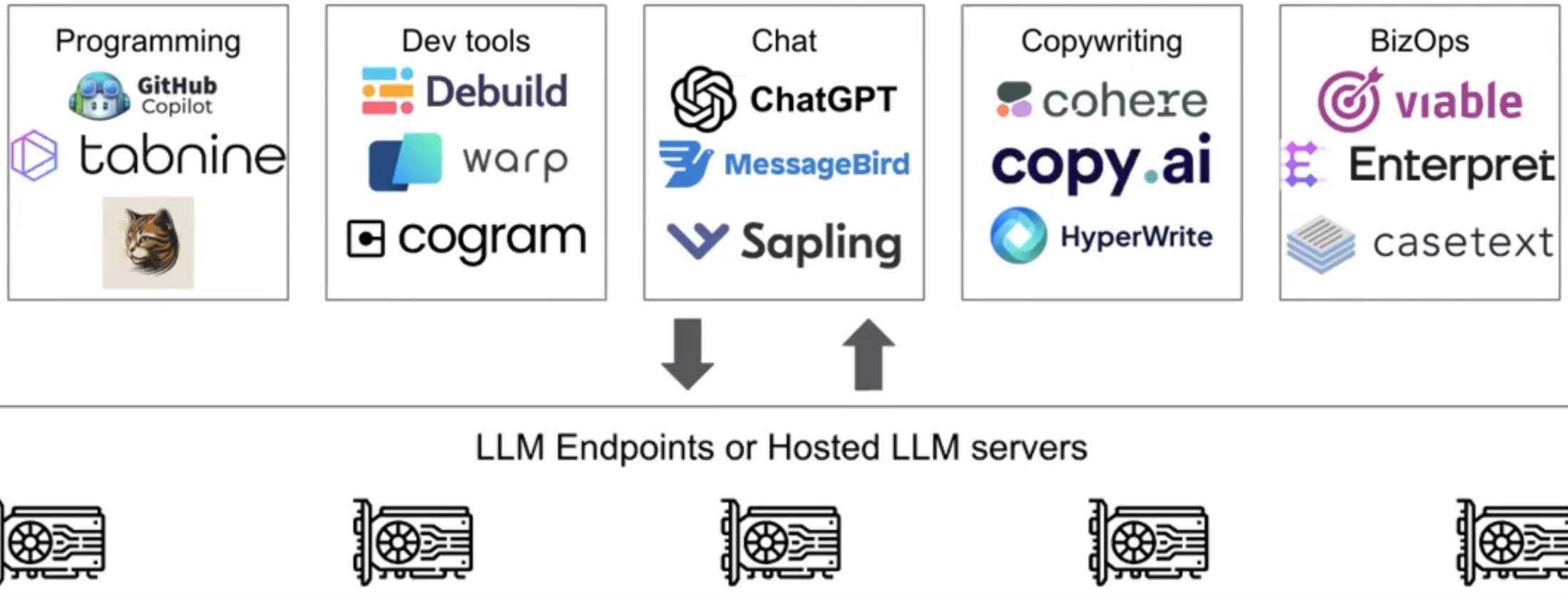


Enterpret



casetext

# LLM-Powered Services



# Serving LLMs is extremely expensive

- LLMs run on high-end GPUs such as NVIDIA A100
- Each GPU can only serve a handful of requests per second
  - For LLaMA-13B and moderate-size inputs, one A100 can process  $< 1$  requests per second
- A ton of GPUs are required for production-scale LLM services



Inference on LLMs is slow (frustratingly high latency), expensive (multiple GPUs or TPUs), & engineering intensive (requires specialized skills to do it well)



## Is local LLM cheaper than ChatGPT API?

ChatGPT api only costs 0.002 dollar for 1k token. I found that LLMs like llama output only 10-20 tokens per second, which is very slow. And such machines costs over 1 dollar per hour. It seems that using api is much cheaper. Based on these observations, it seems that utilizing the ChatGPT API might be a more affordable option.

# Inference process of LLMs

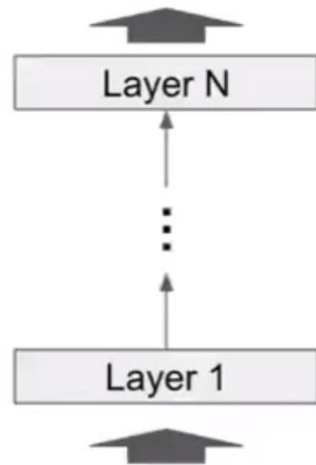
**Output**

**Input**

Artificial	Intelligence	is
------------	--------------	----

# Inference process of LLMs

**Output**

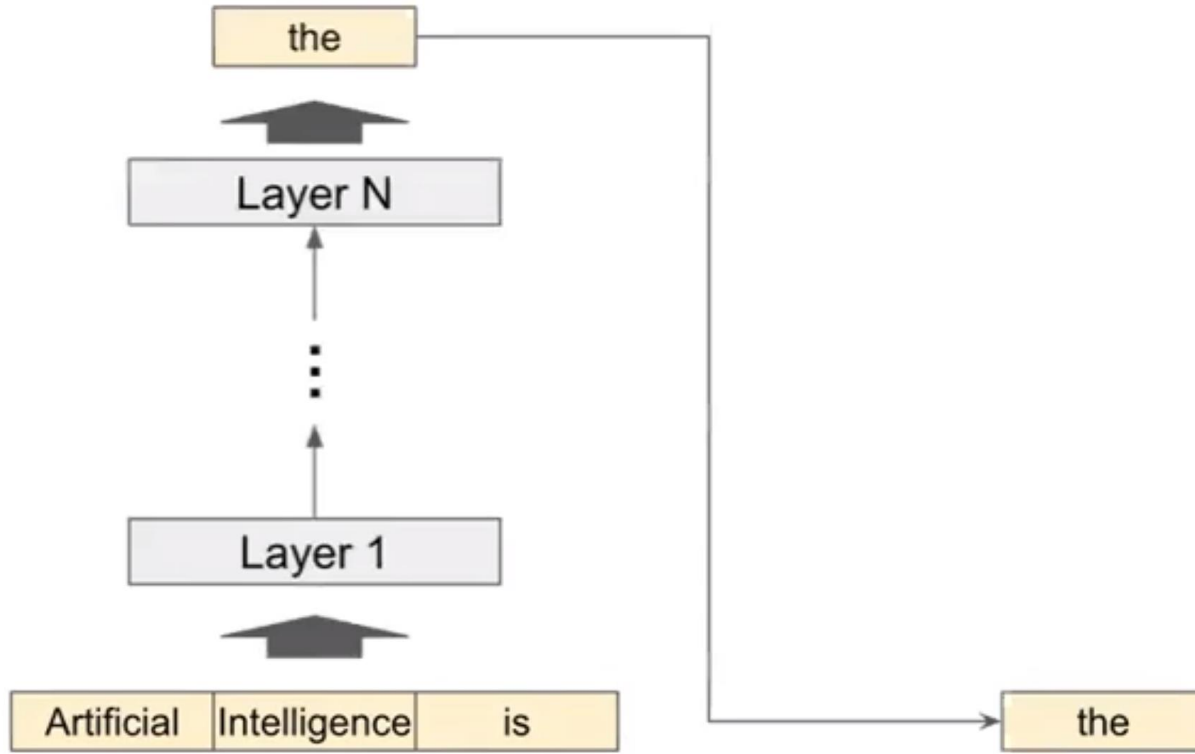


**Input**



# Inference process of LLMs

Output

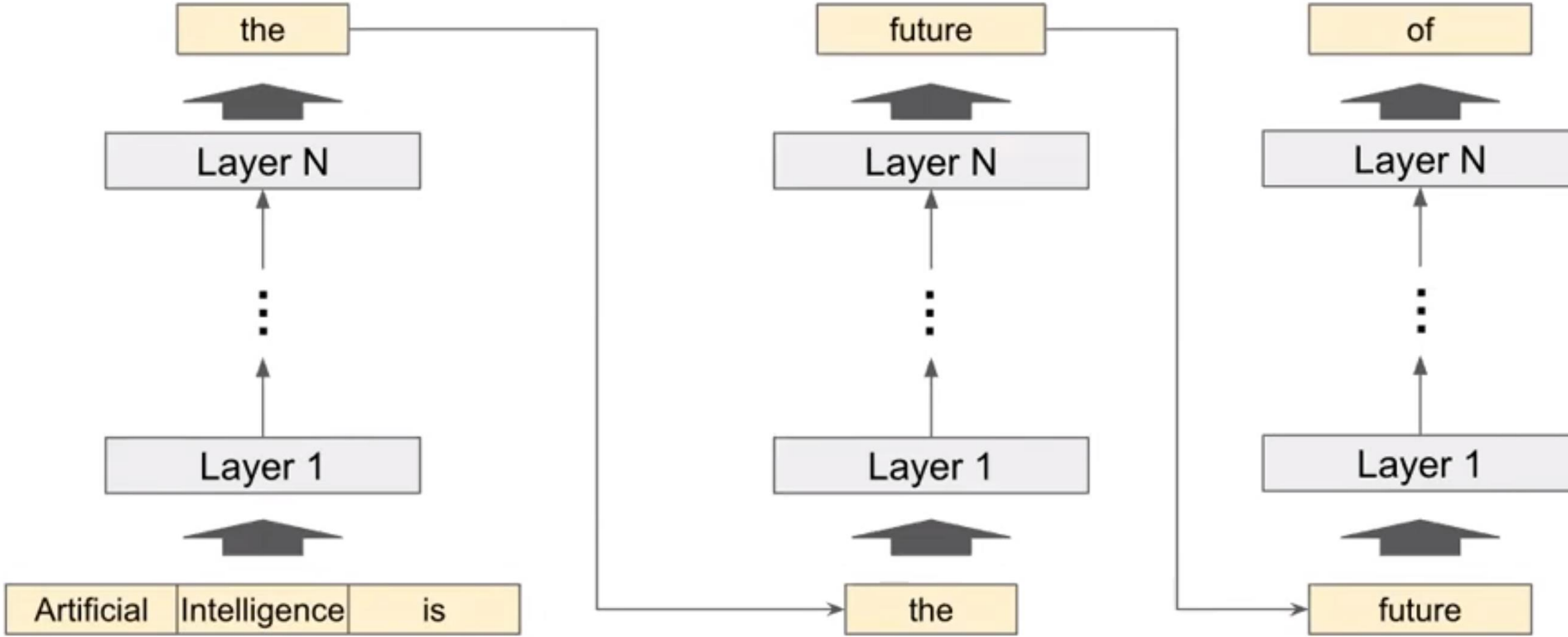


Input



# Inference process of LLMs

Output

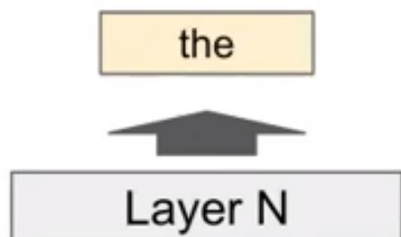


Input

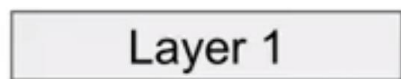
- Repeat until the sequence
  - Reaches its pre-defined max length (e.g., 4K tokens)
  - Generates an EOS (end of sequence) token

# KV Cache

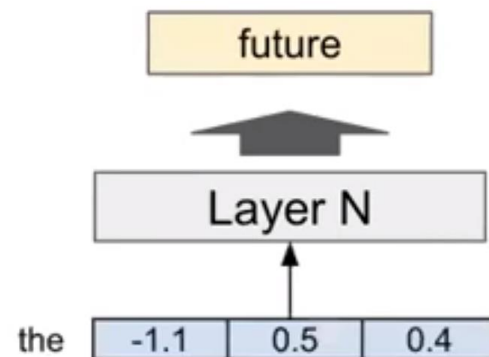
Output



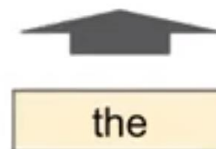
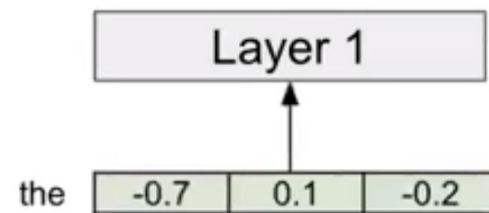
⋮



Input

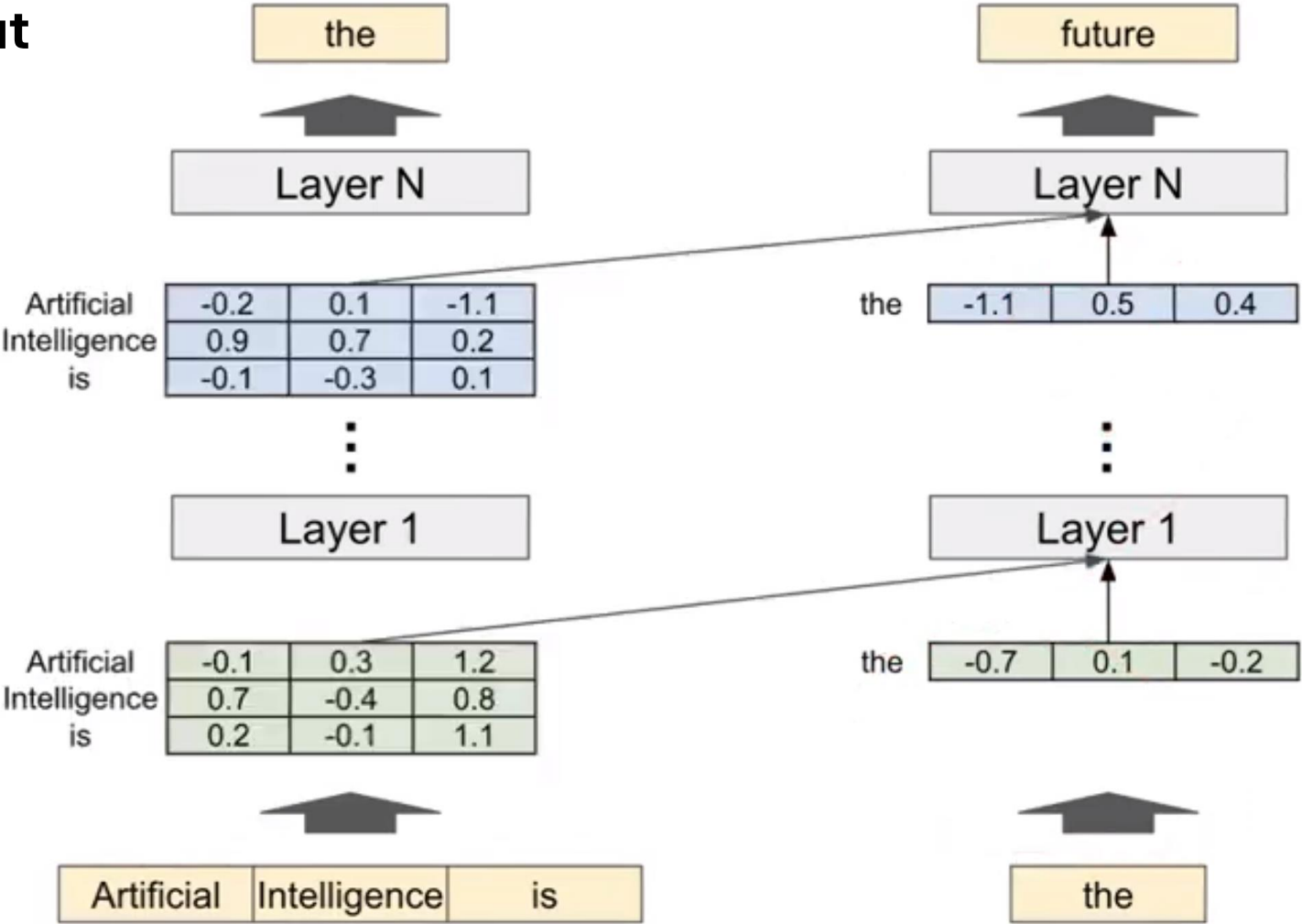


⋮



# KV Cache

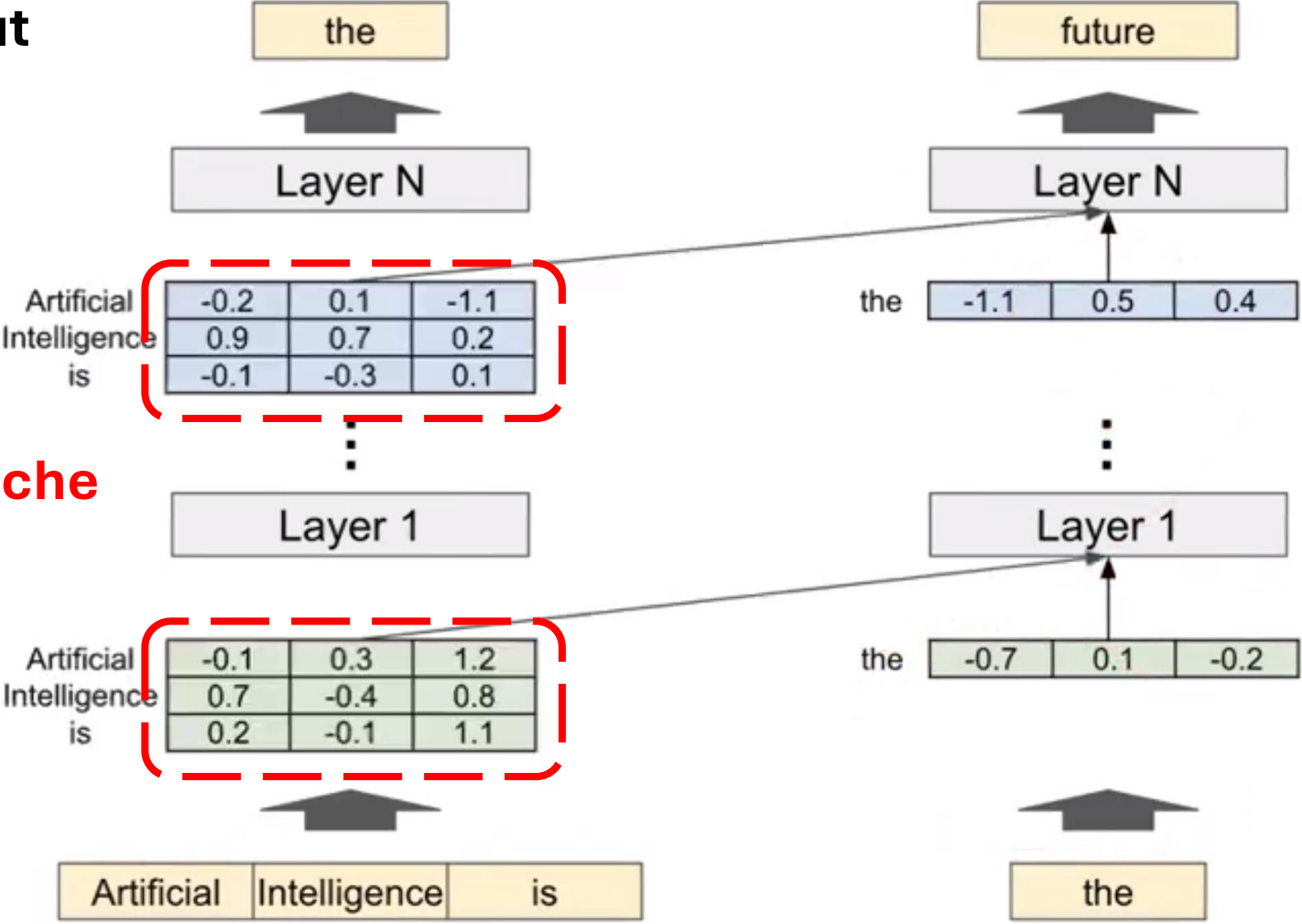
Output



Input

# KV Cache

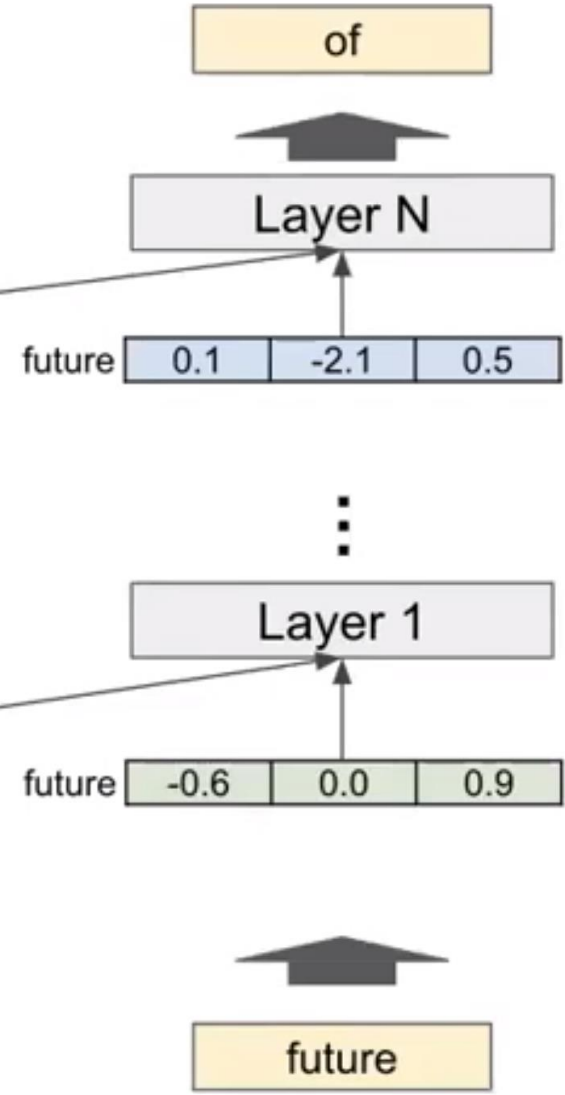
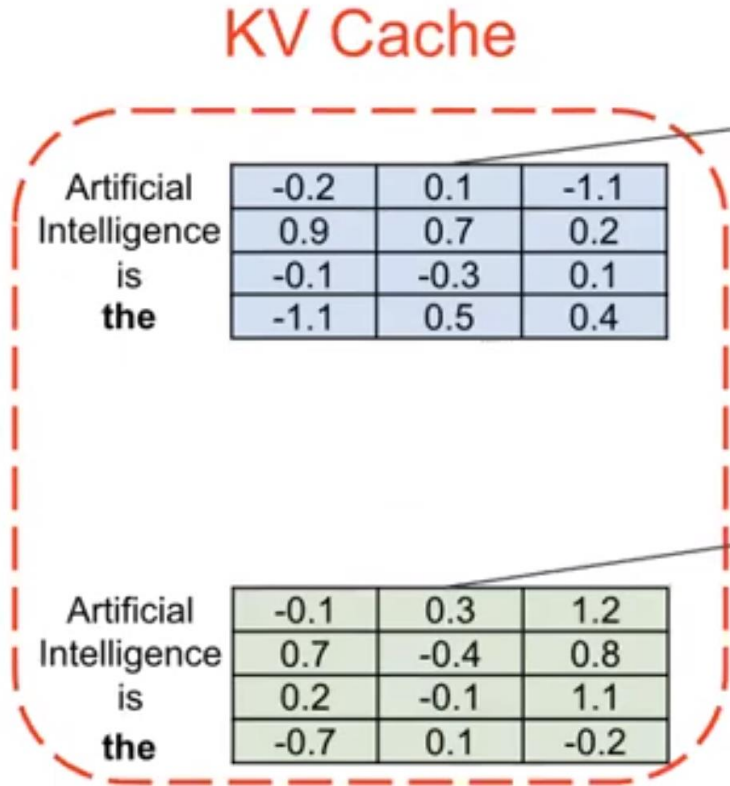
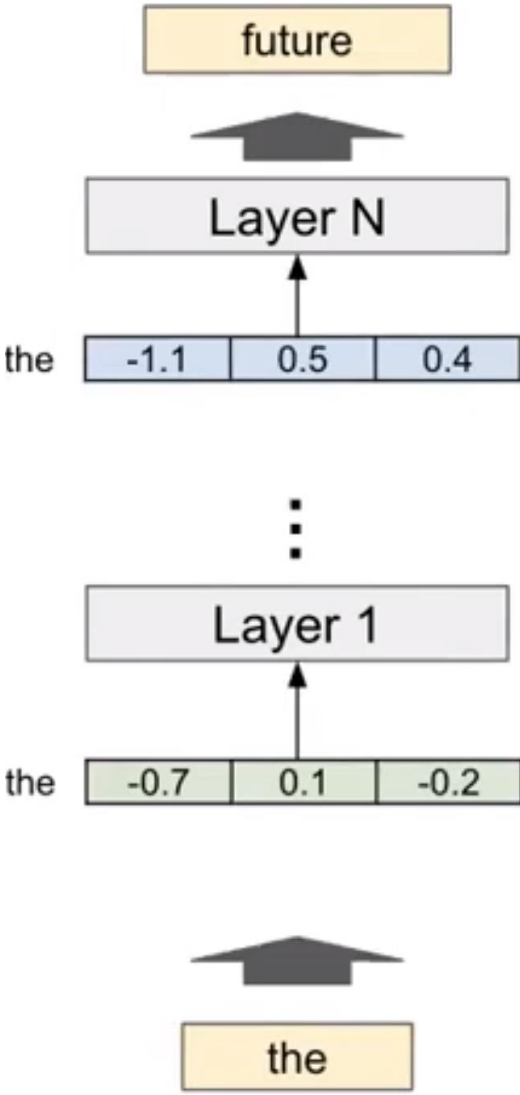
Output



KV Cache

Input

# KV Cache

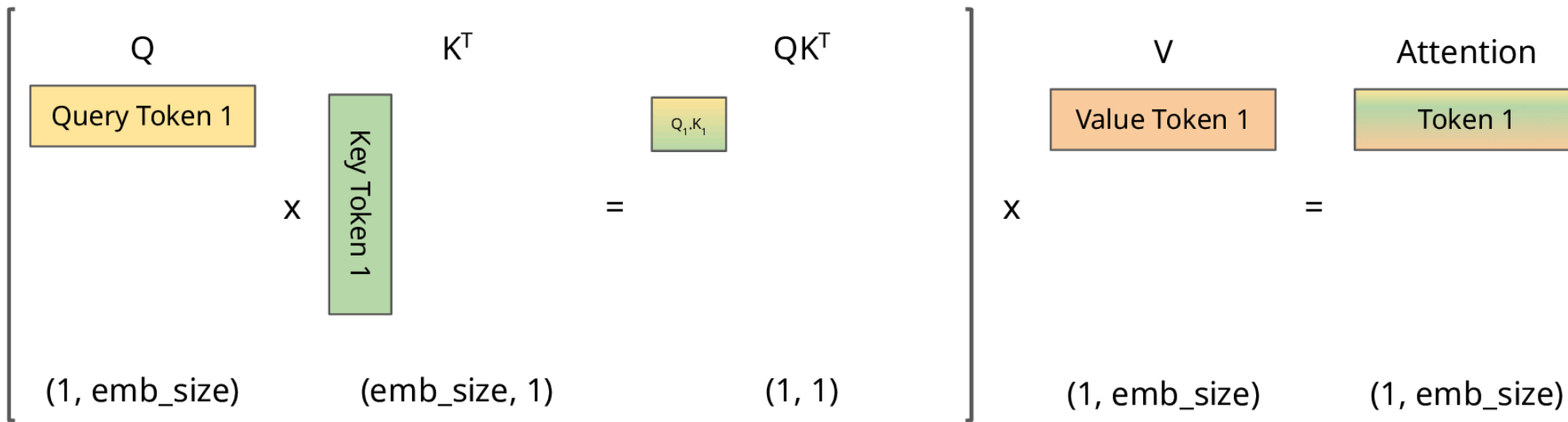


**20 KB / layer / token  
= 800 KB / token  
(LLaMA-13B)**

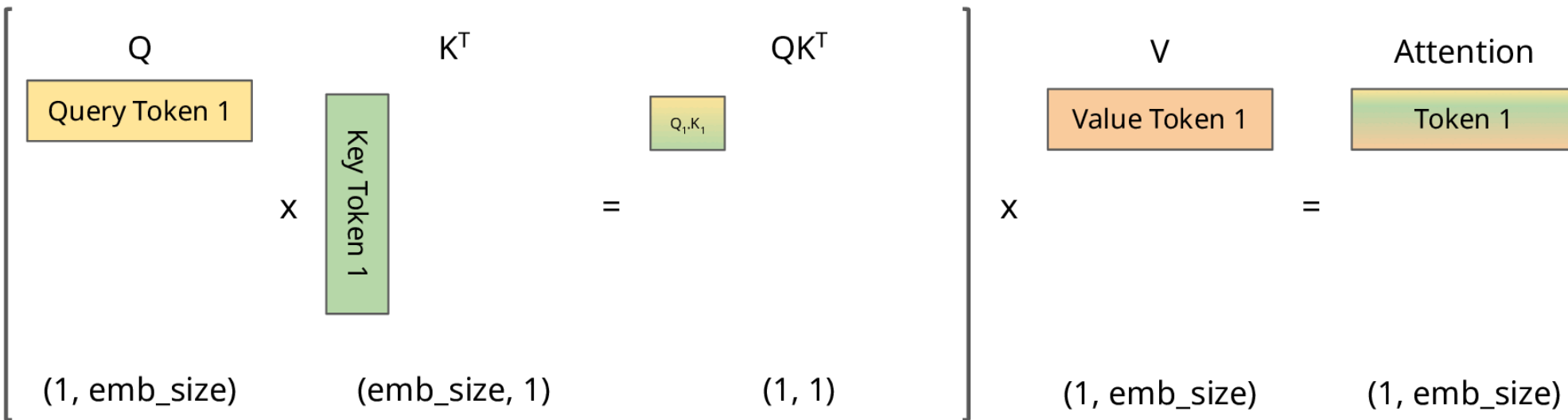
# KV Cache

## Step 1

*Without  
cache*

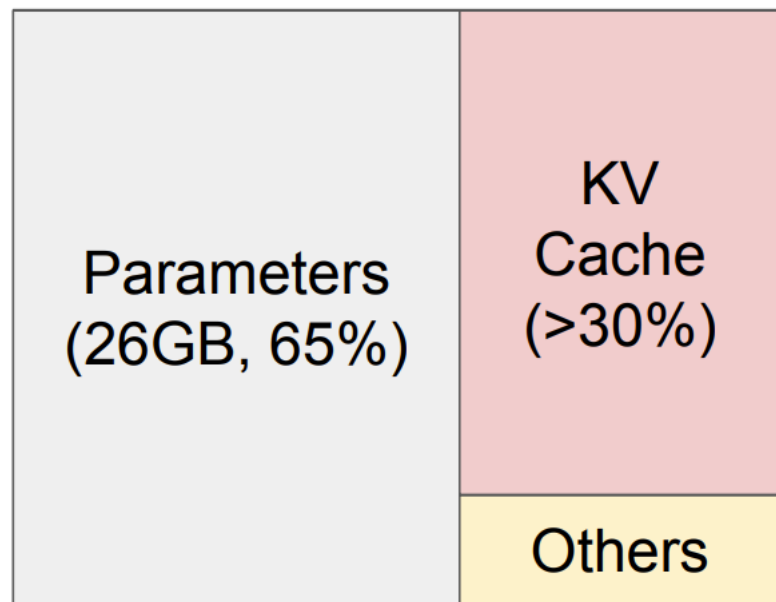


*With  
cache*

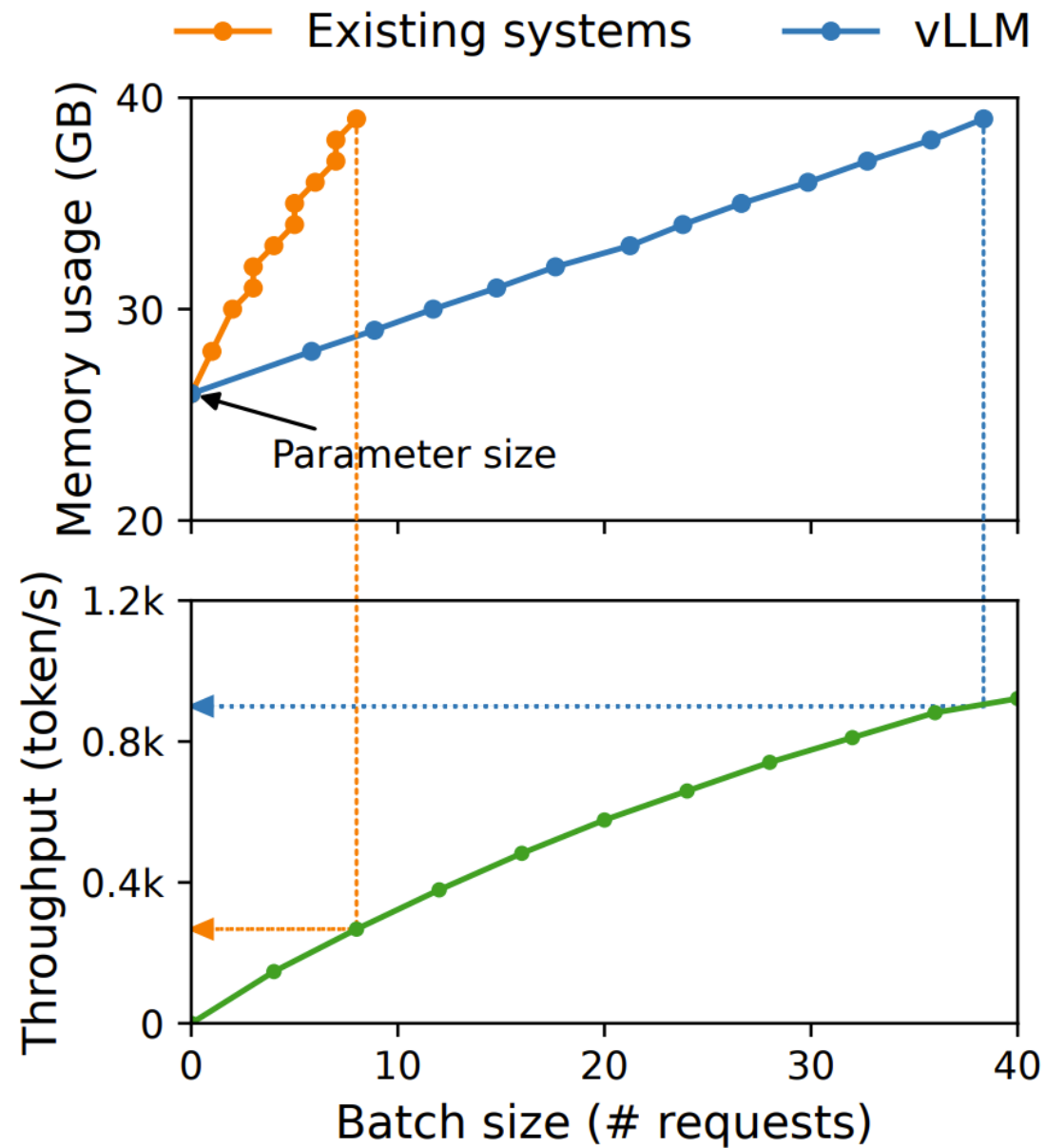


  Values that will be masked    
   Values that will be taken from cache

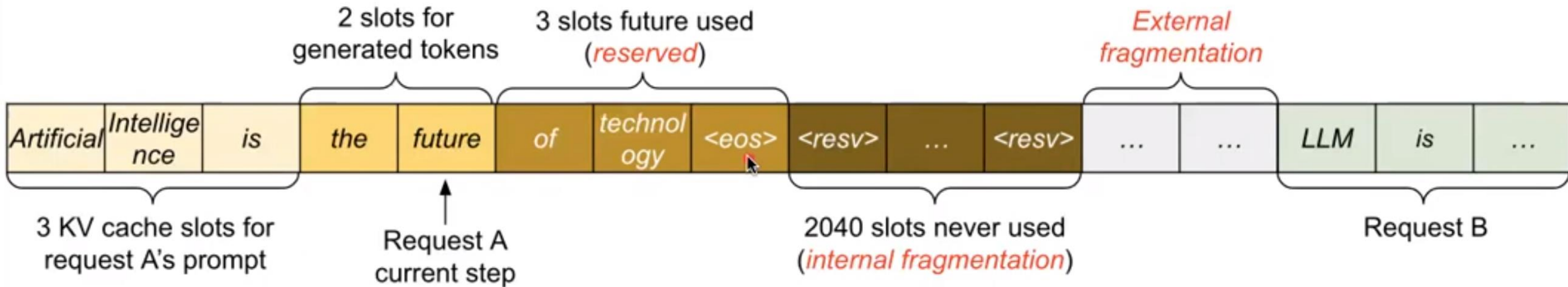
# Key Insight



13B LLM on A100-40GB



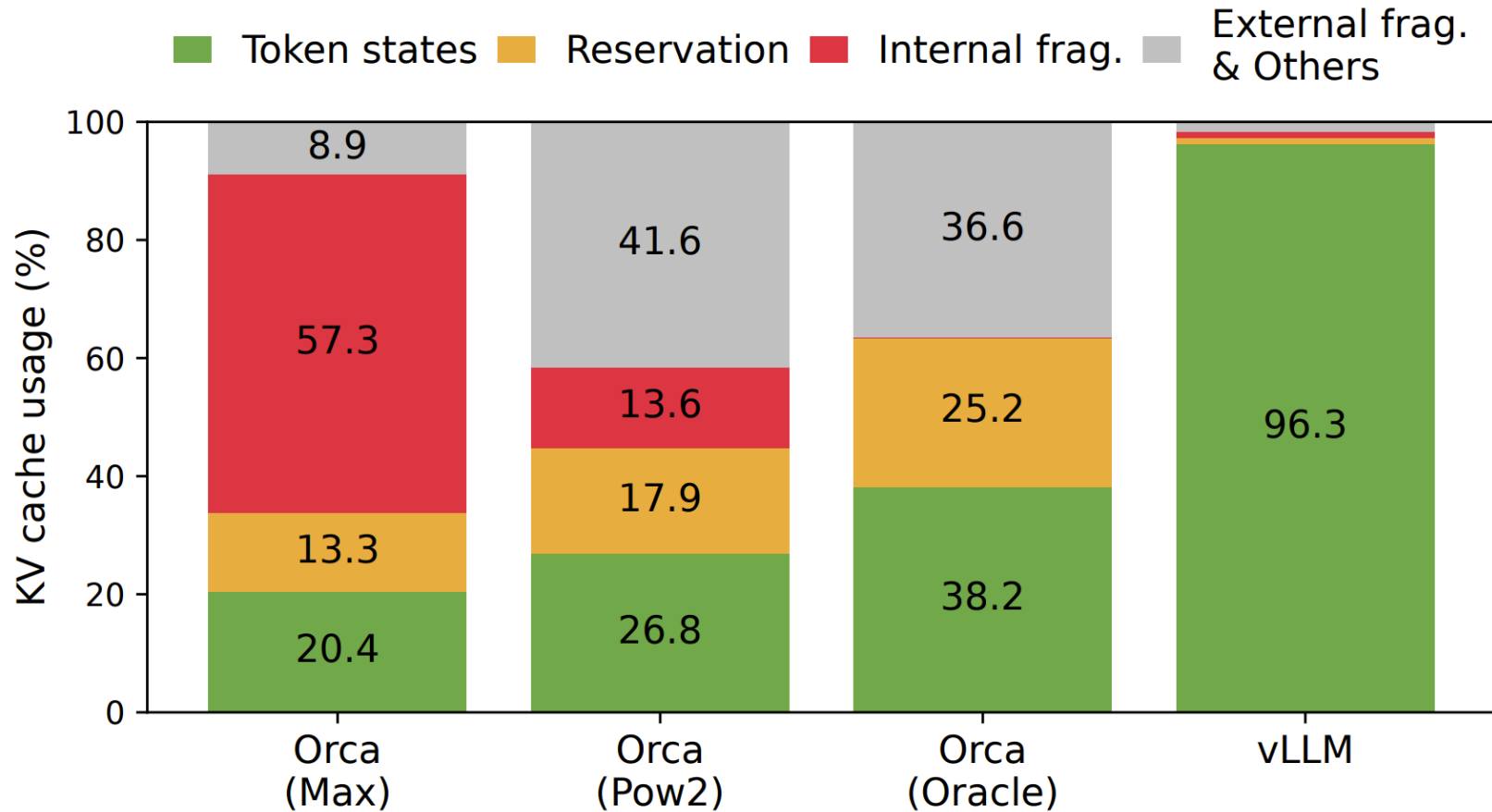
# Memory Waste in KV Cache



- **Internal fragmentation:** over-allocated due to the unknown output length
- **Reservation:** not used at the current step, but used in the future
- **External fragmentation:** due to different sequence lengths



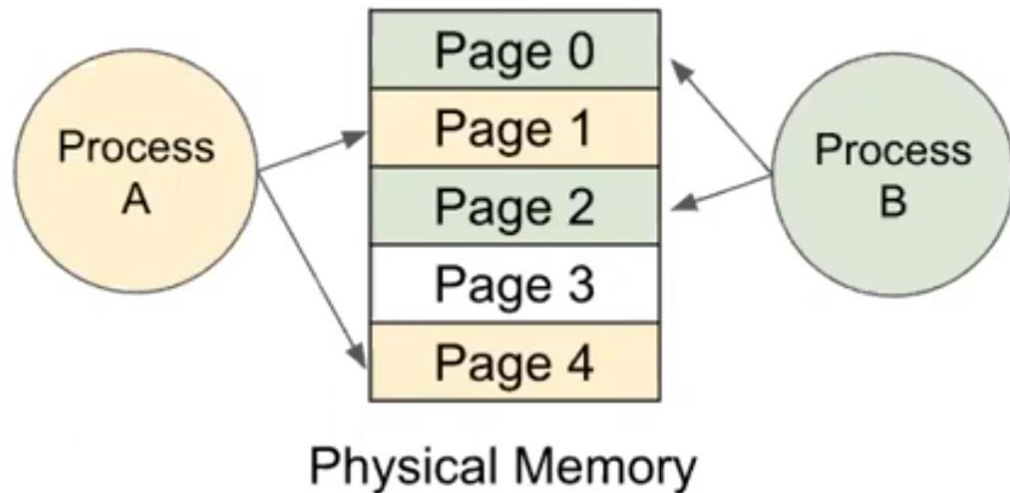
# Memory Waste in KV Cache



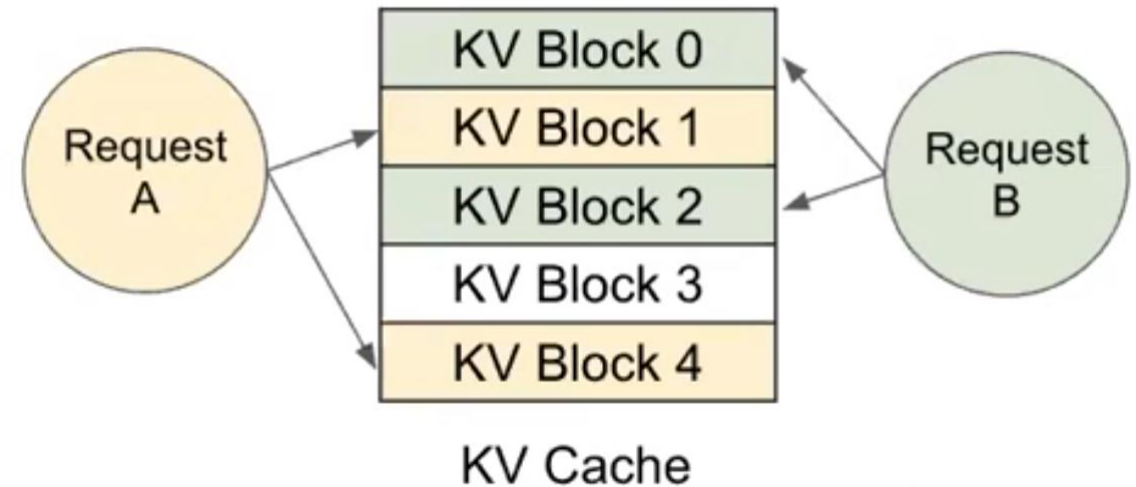
Only **20% - 40%** of KV cache is utilized to store token states

# vLLM: Efficient Memory Management for LLM Inference

## Memory management in OS

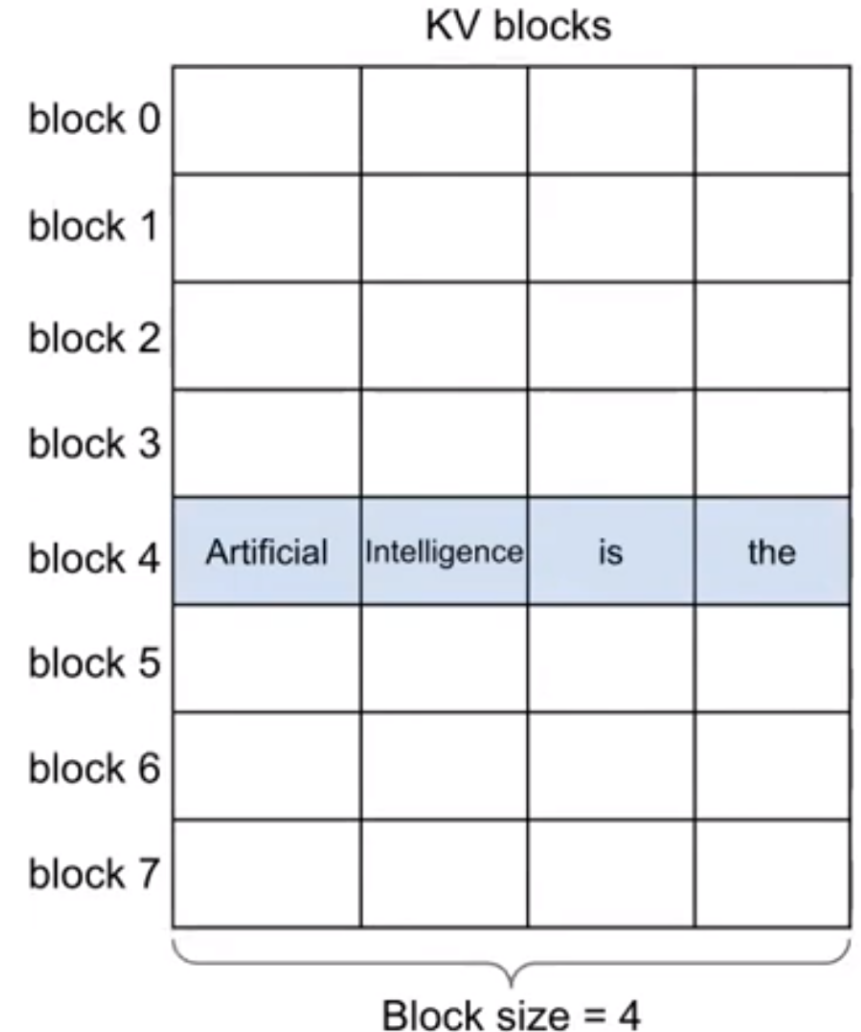


## Memory management in vLLM



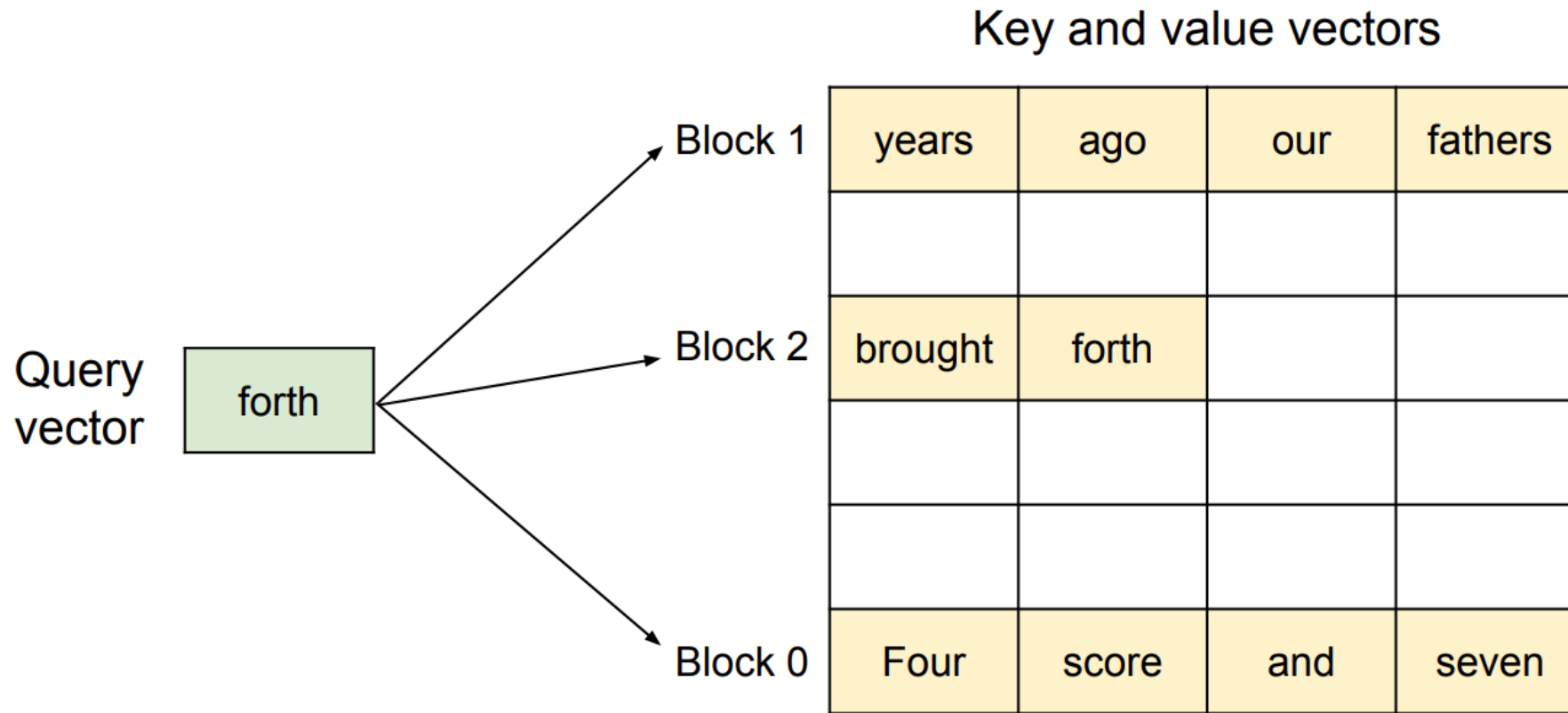
# KV Block

- A fixed-size contiguous chunk of memory that stores the tokens' KV states
- Similar to the concept of a memory page



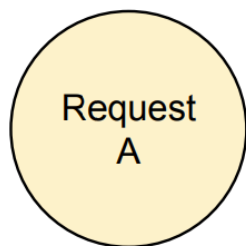
# PagedAttention

- Manages KV cache in block granularity instead of sequence (i.e., request) granularity
- Allows storing logically contiguous KV blocks in non-contiguous physical memory



Example sequence: “Four score and seven years ago our | fathers brought forth”

# Logical-to-Physical KV Block Translation



**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				

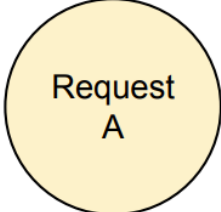
**Block Table**

Physical block number	# filled
①7	①4
①1	①3 → 4②
③3	③1
-	-

**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				

# Serving Multiple Requests

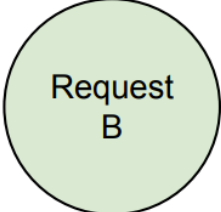


Logical KV blocks

Block 0	<i>Four</i>	<i>score</i>	<i>and</i>	<i>seven</i>
Block 1	<i>years</i>	<i>ago</i>	<i>our</i>	<i>fathers</i>
Block 2	<i>brought</i>			
Block 3				

Physical KV blocks

Block 0				
Block 1	<i>years</i>	<i>ago</i>	<i>our</i>	<i>fathers</i>
Block 2	<i>of</i>	<i>times</i>		
Block 3	<i>brought</i>			
Block 4				
Block 5	<i>It</i>	<i>was</i>	<i>the</i>	<i>best</i>
Block 6				
Block 7	<i>Four</i>	<i>score</i>	<i>and</i>	<i>seven</i>
Block 8				



Logical KV blocks

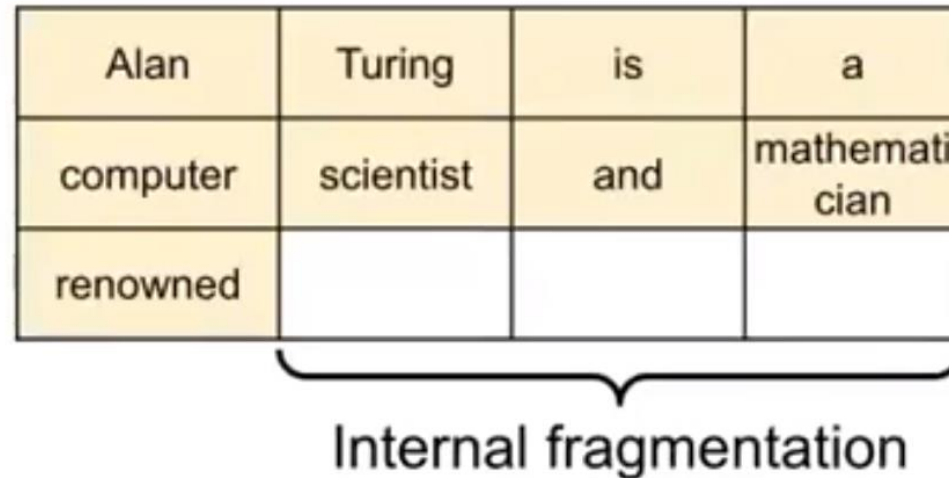
Block 0	<i>It</i>	<i>was</i>	<i>the</i>	<i>best</i>
Block 1	<i>of</i>	<i>times</i>		
Block 2				

Block 0  
Block 1  
Block 2  
Block 3  
Block 4  
Block 5  
Block 6  
Block 7  
Block 8

Block 0  
Block 1  
Block 2

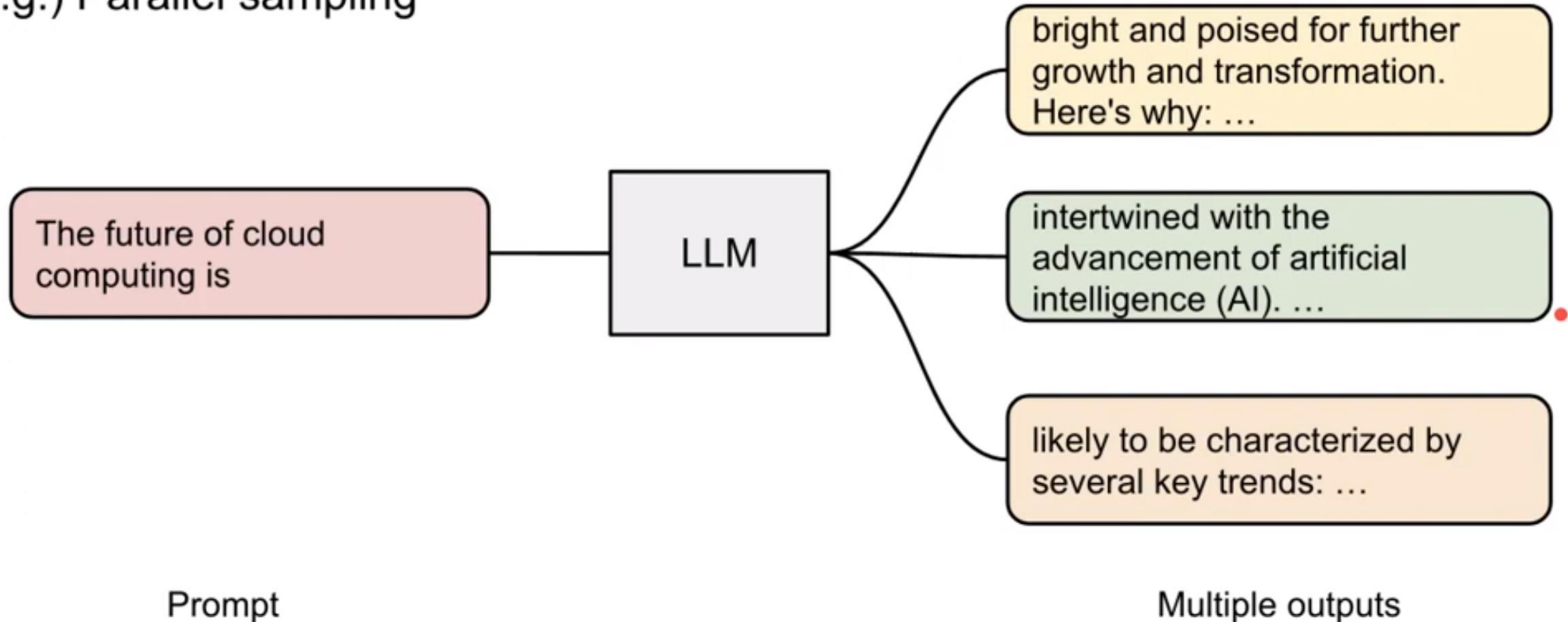
# Memory Efficiency of vLLM

- Minimizes **internal fragmentation**
  - Only happens at the last block of a sequence
  - # wasted tokens per sequence < block size
    - Seq len:  $O(100)$  –  $O(1000)$  tokens
    - Block size: 16 or 32 tokens



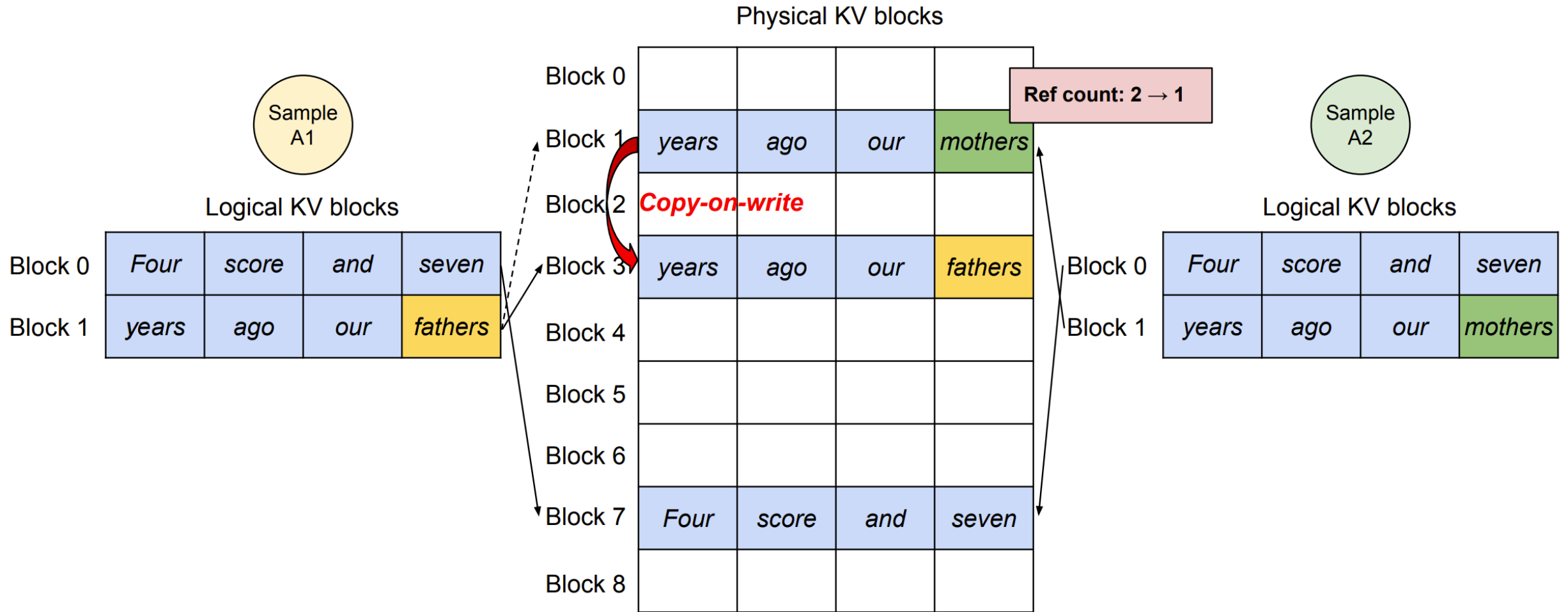
# Dynamic Block Mapping Enables Sharing

E.g.) Parallel sampling

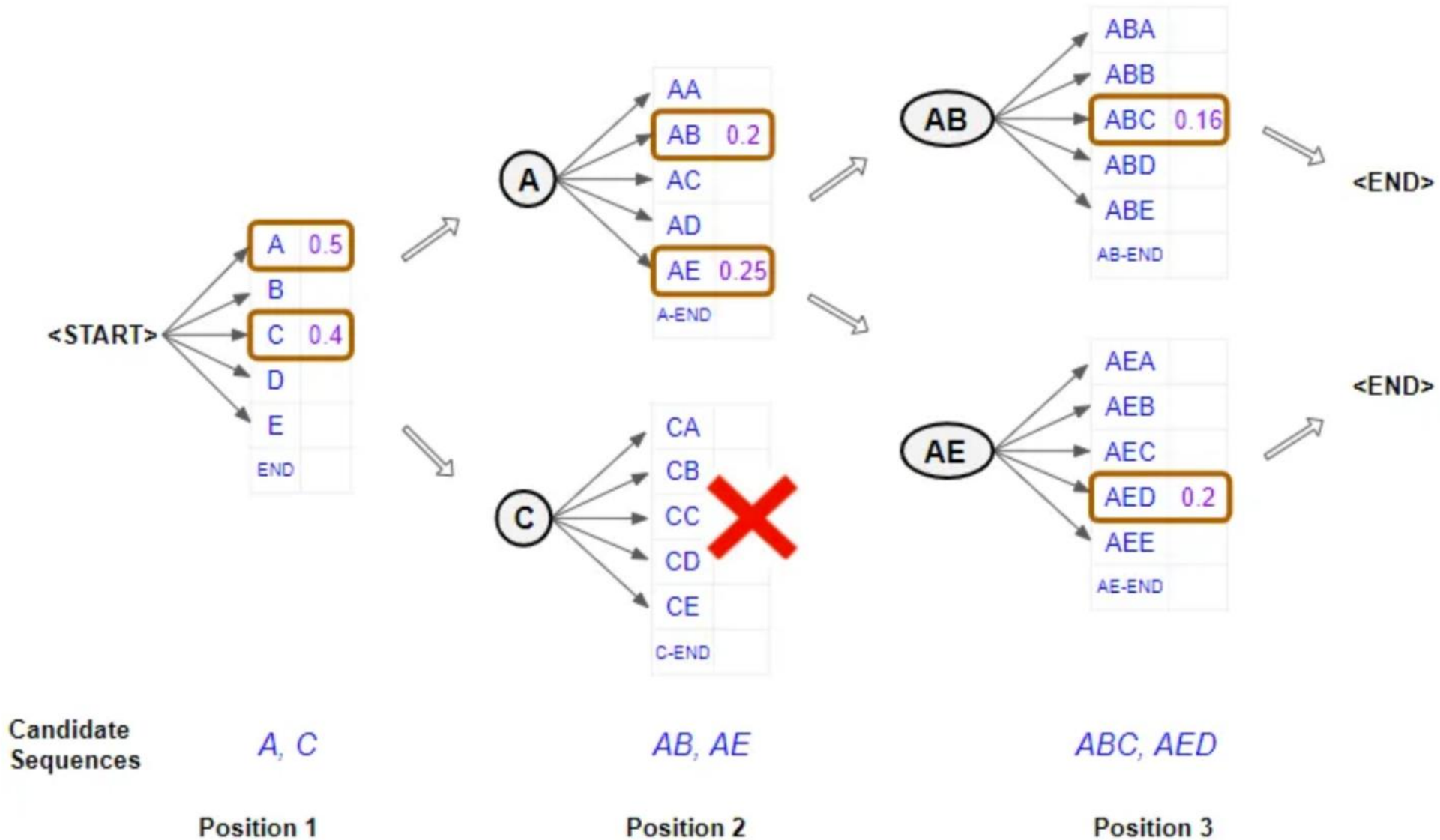




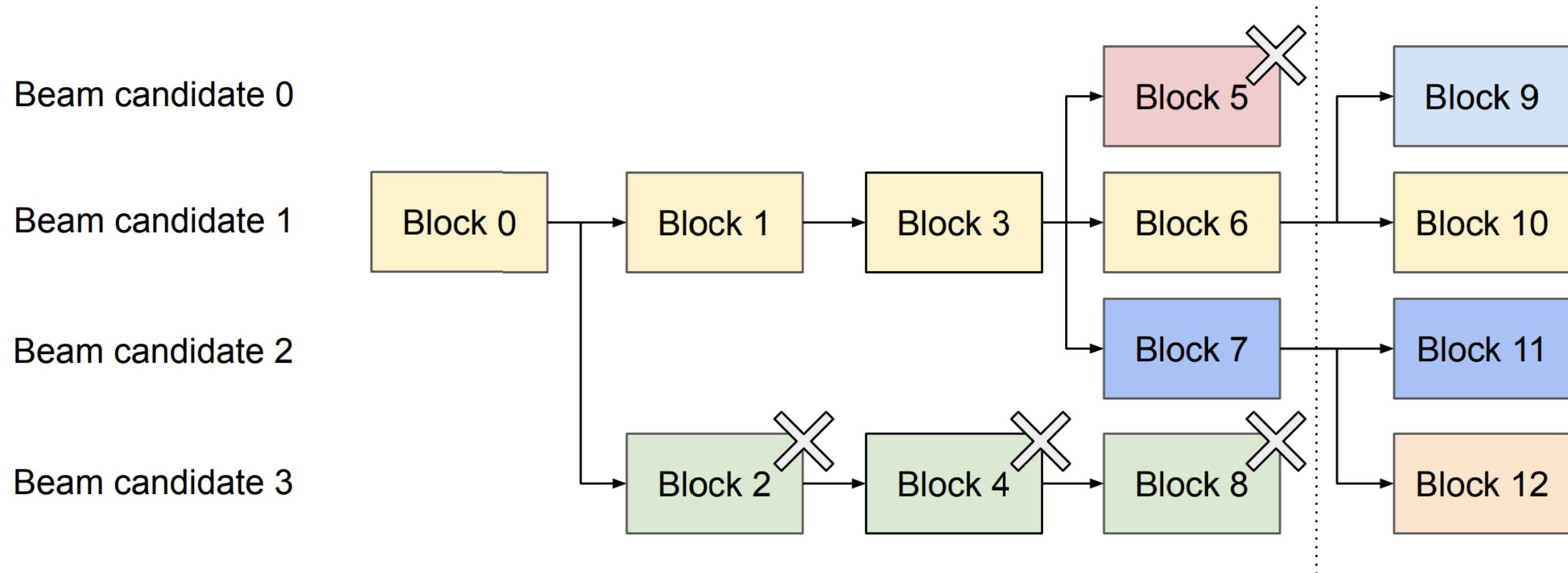
# KV Block Sharing for Parallel Sampling



# KV Block Sharing for Beam Search

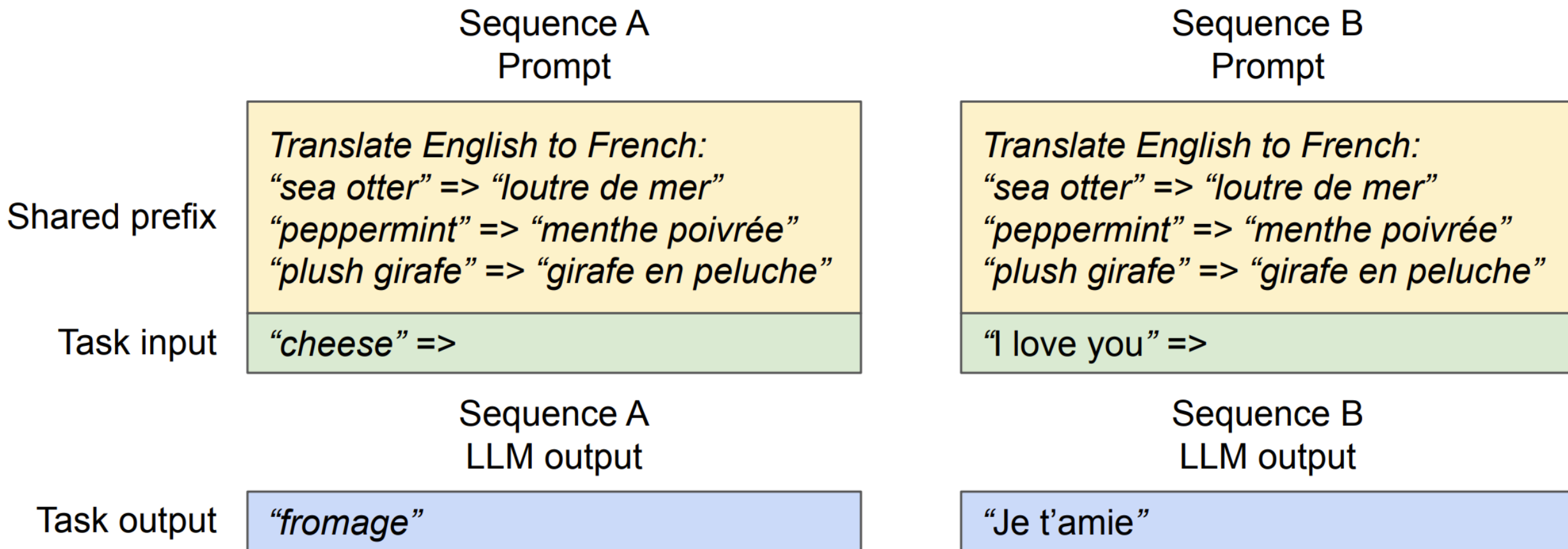


# KV Block Sharing for Beam Search



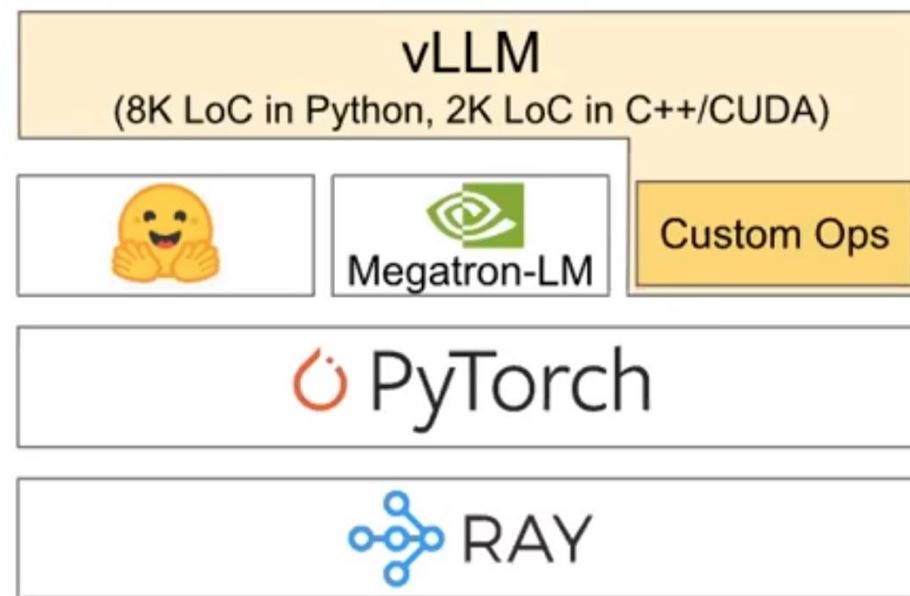
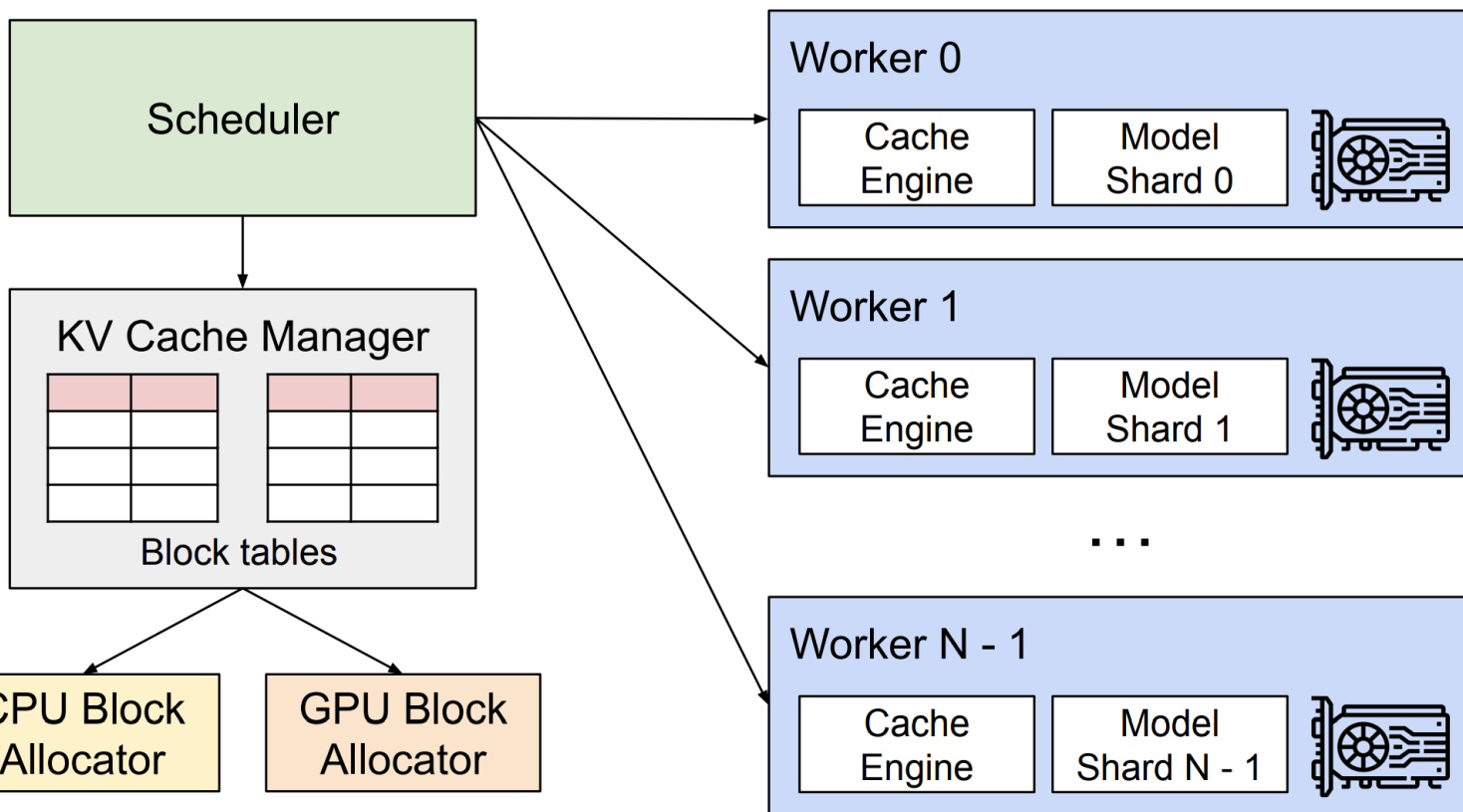
- Similar to process **fork** and **kill**

# Shared Prompt



- Similar to shared libraries in OS

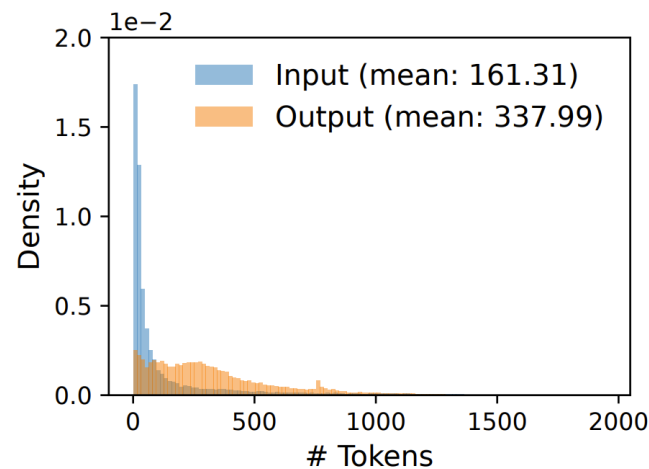
# vLLM System Architecture & Implementation



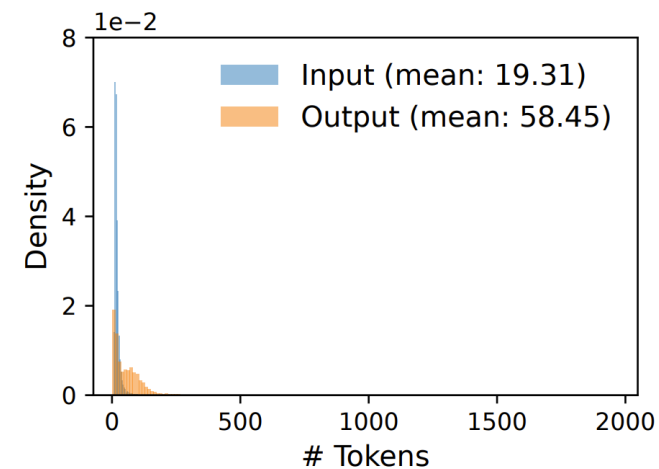
# Evaluation

System configuration.

Model size	13B	66B	175B
GPUs	A100	4×A100	8×A100-80GB
Total GPU memory	40 GB	160 GB	640 GB
Parameter size	26 GB	132 GB	346 GB
Memory for KV cache	12 GB	21 GB	264 GB
Max. # KV cache slots	15.7K	9.7K	60.1K



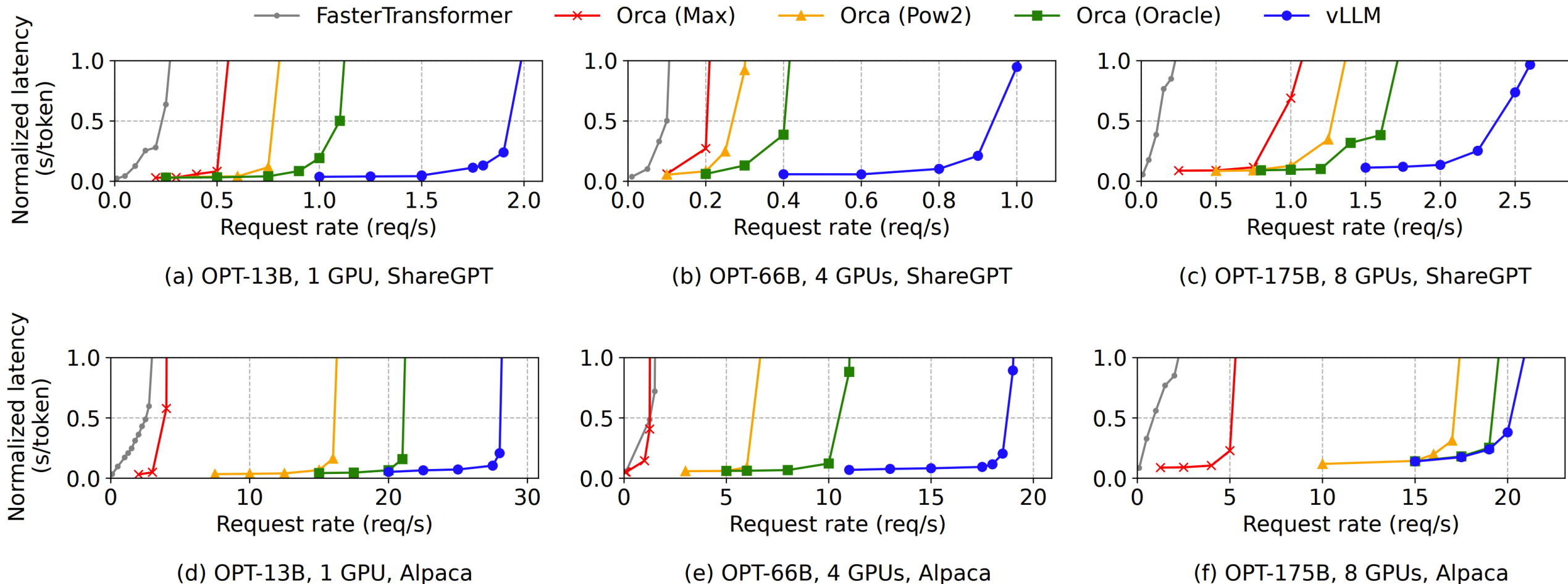
(a) ShareGPT



(b) Alpaca

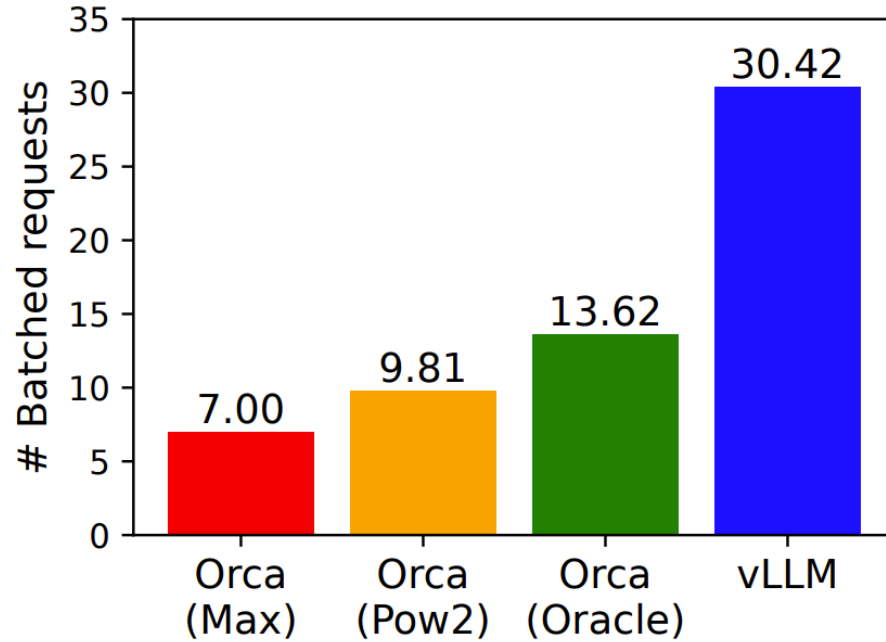
**Figure 11.** Input and output length distributions of the (a) ShareGPT and (b) Alpaca datasets.

# vLLM Improves Inference Throughput by Enabling Larger Batch Size

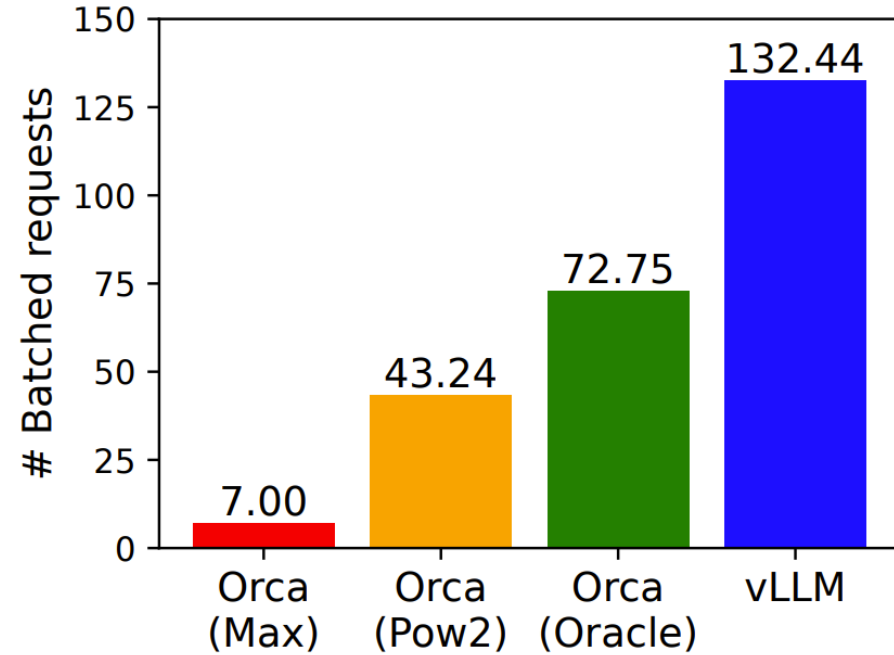


**Figure 12.** Single sequence generation with OPT models on the ShareGPT and Alpaca dataset

# vLLM Improves Inference Throughput by Enabling Larger Batch Size



(a) ShareGPT

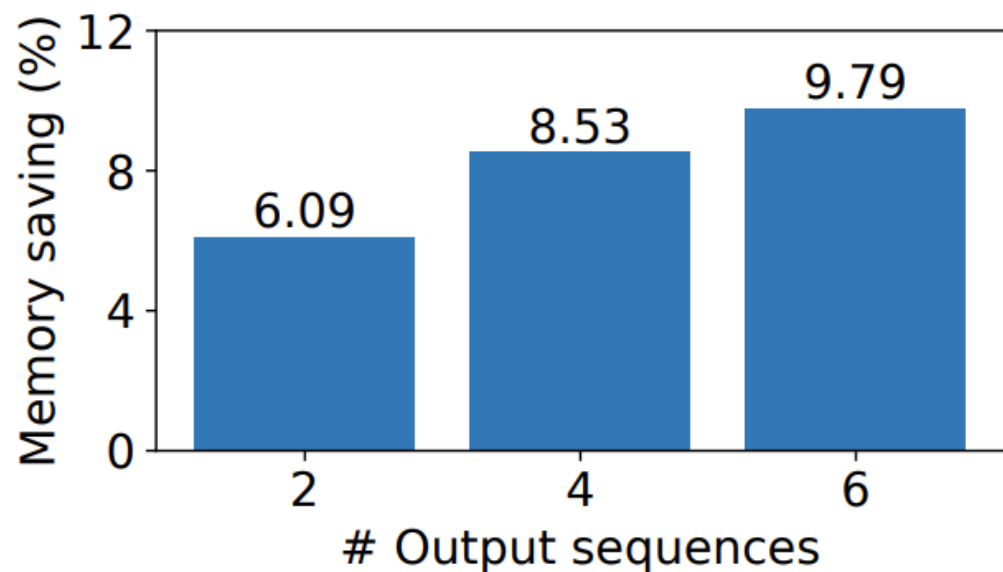


(b) Alpaca

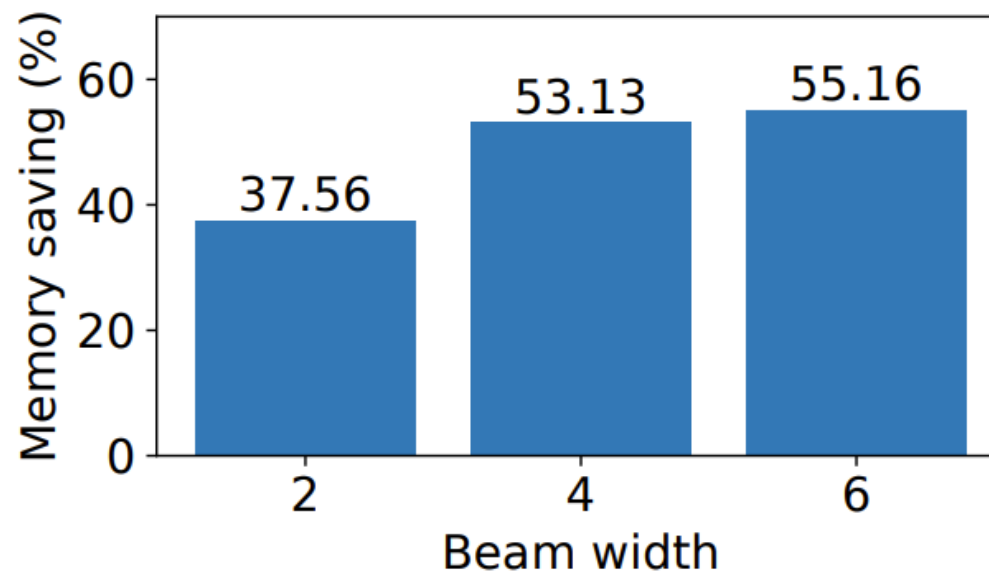
**Figure 13.** Average number of batched requests when serving OPT-13B for the ShareGPT (2 reqs/s) and Alpaca (30 reqs/s) traces.



# Memory Saving of vLLM



(a) Parallel sampling



(b) Beam search

**Figure 15.** Average amount of memory saving from sharing KV blocks, when serving OPT-13B for the Alpaca trace.

# Takeaways

- Strength
  - Interesting observation on the KV cache memory inefficiency
  - Analogy between KV cache management and OS paging
  - Open-source implementation
- Weakness
  - Cannot fundamentally improve inference latency
  - For multi-chip execution, vLLM assumes attention heads are sharded
    - If even a single attention head is too large, or we want to split it across multiple chips to improve latency, how can vLLM support sharding the KV cache?
  - The fundamental bottlenecks faced by LLM serving, **memory capacity** due to large model weights, and **memory bandwidth** due to low compute intensity of auto-regressive decoding, remain unsolved.
    - Speculative decoding?
    - New model architectures? (e.g., [SSM](#), [Mamba](#))
    - New hardware/architectural innovations? (e.g., processing-in-memory, [NVIDIA's new patent](#) that proposes stacking HBM dies on the processor die to expose a wider mem interface)

Q & A