# CSE 234: Data Systems for Machine Learning
# Winter 2025

LLMSys

Optimizations and Parallelization

MLSys Basics

# Logistics

- If 80% of you finish the course eval, all get +2 points in final score!
  - Currently: we are 50%
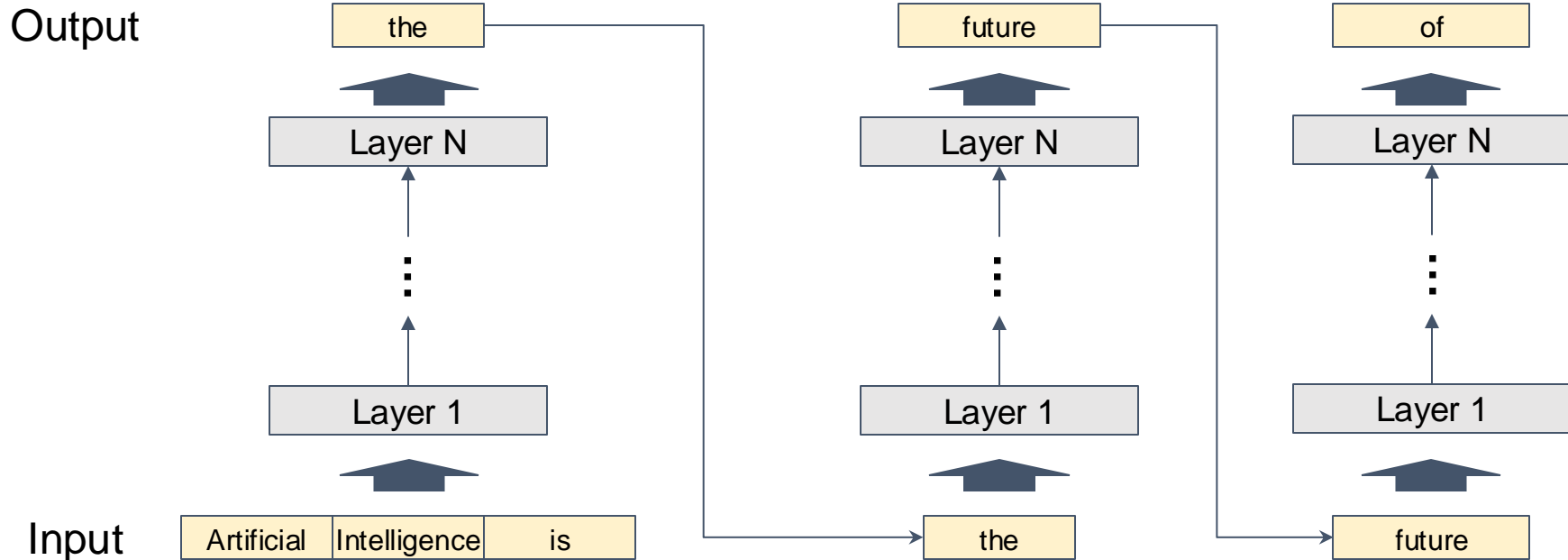- TA will hold a recitation for exam – make sure to attend

# Recap: Next Token Prediction

Probability("San Diego has very nice weather")
= P("San Diego") P("has"|"San Diego")P("very"|"San Diego has")P("city"|...)...P("weather"|...)

$$\text{Max} Prob(x_{1:T}) = \boxed{\prod_{t=1}^{T} P(x_{t+1}|x_{1...t})}$$

This is model we got – capable of "predicting the next token".

# Inference process of LLMs

Output

| the | | future | | of |

Layer N — Layer N — Layer N

⋮ — ⋮ — ⋮

Layer 1 — Layer 1 — Layer 1

Input

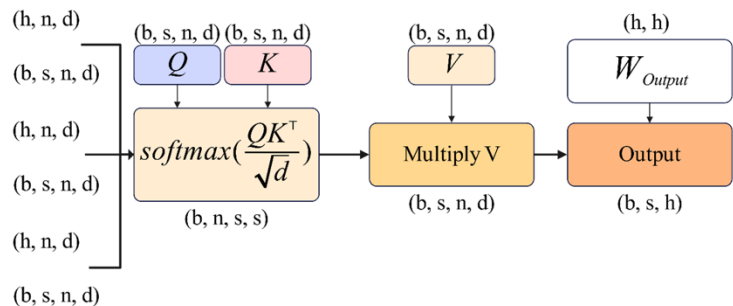| Artificial | Intelligence | is | → | the | → | future |

Repeat until the sequence
- Reaches its pre-defined maximum length (e.g., 2048 tokens)
- Generates certain tokens (e.g., "<|*end of sequence*|>")
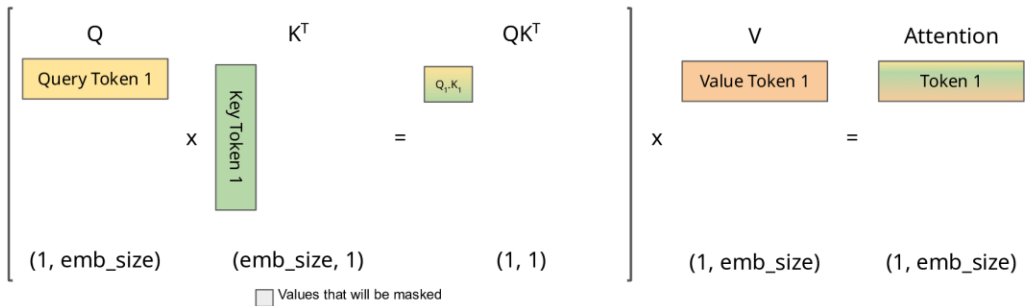
# Generative LLM Inference: Autoregressive Decoding

- <span style="color:red">Pre-filling phase (0-th iteration):</span>
  - Process *all* input tokens at once
- <span style="color:red">Decoding phase (all other iterations):</span>
  - Process a *single* token generated from previous iteration

- <span style="color:red">Key-value cache:</span>
  - Save attention keys and values for the following iterations to avoid recomputation
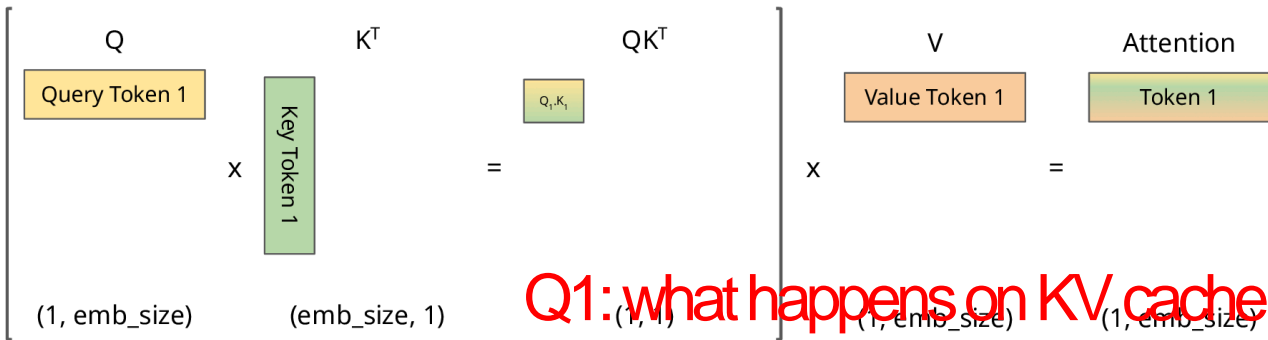  - what is KV cache essentially?

# w/ KV Cache vs. w/o KV Cache



Zoom-in! (simplified without Scale and Softmax)
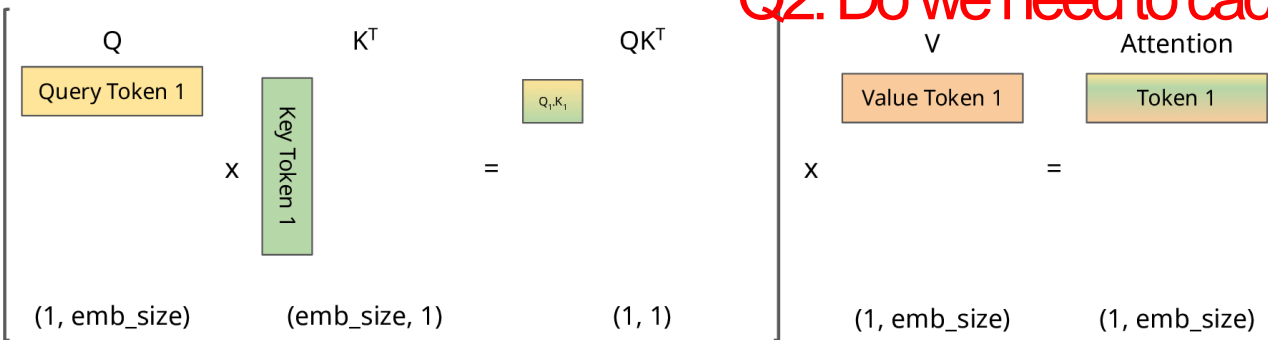
# w/ KV Cache vs. w/o KV Cache

**Step 1**

Without cache

| Q | | $K^T$ | | $QK^T$ |
|---|---|---|---|---|

Query Token 1

Key Token 1

$Q_1.K_1$

x

=

(1, emb_size)    (emb_size, 1)    (1, 1)

V           Attention

Value Token 1        Token 1

x

=

(1, emb_size)    (1, emb_size)

With cache

Q           $K^T$           $QK^T$

Query Token 1

Key Token 1

$Q_1.K_1$

x

=

(1, emb_size)    (emb_size, 1)    (1, 1)

V           Attention

Value Token 1        Token 1

x

=

(1, emb_size)    (1, emb_size)

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

Q1: what happens on KV cache in prefill phase?
Q2: Do we need to cache Q?

# Potential Bottleneck of LLM Inference?



- Compute:
  - Prefill: largely same with training
  - Decode: $s = 1$
- Memory
  - New: KV cache
- Communication
  - mostly same with training

Q? how about batch size b?

# Serving vs. Inference



large b

b = 1

**Serving**: many requests, online traffic, emphasize cost-per-query.

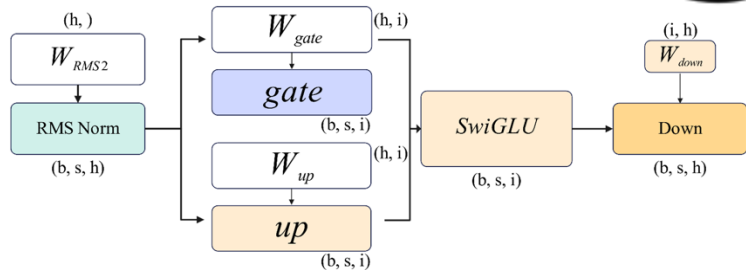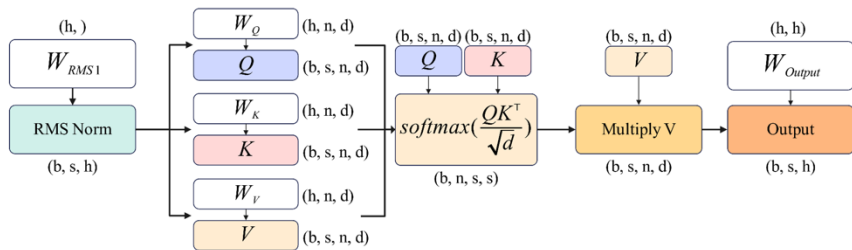s.t. some mild latency constraints

emphasize **throughput**

**Inference**: fewer request, low or offline traffic,

emphasize **latency**
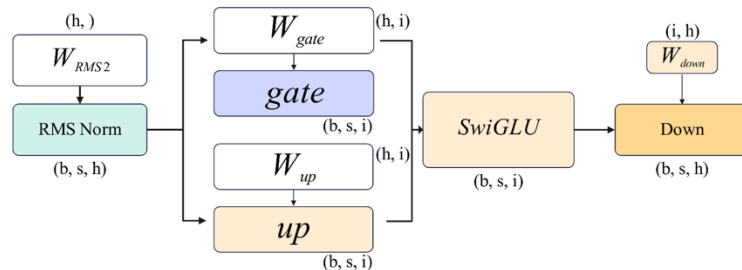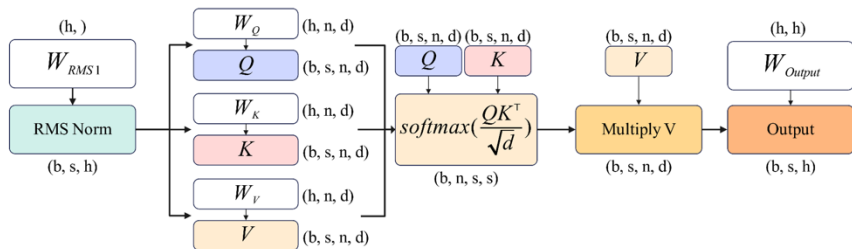
# Potential Bottleneck of LLM Inference in Serving



- Compute:
  - Prefill:
    - Different prompts have different length: how to batch?
  - Decode
    - Different prompts have different, unknown #generated tokens
    - s = 1, b is large
- Memory
  - New: KV cache
  - b is large -> KV is linear with b -> will KVs be large?
- Communication
  - mostly same with training

b=1

# Potential Bottleneck of LLM Inference in Serving
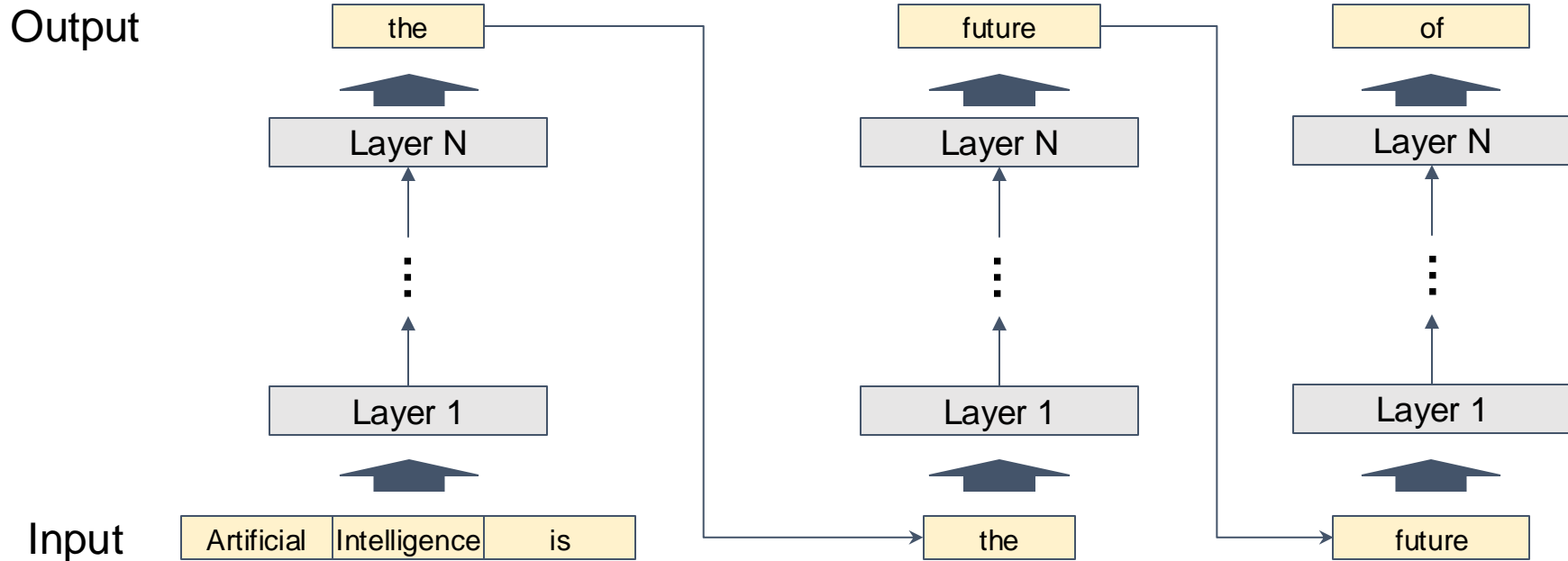


- Compute:
  - Prefill:
    - ~~Different prompts have different length: how to batch?~~
  - Decode
    - Different prompts have different, unknown #generated tokens
    - s = 1, b=1
- Memory
  - New: KV cache
  - ~~b =1 -> KV is linear with b -> will KVs be large?~~
- Communication
  - mostly same with training

Problems of bs = 1

max AI = #ops / #bytes

# Recap: Inference process of LLMs

Output

| the | | future | | of |



Input

| Artificial | Intelligence | is | | the | | future |

Repeat until the sequence
- Reaches its pre-defined maximum length (e.g., 2048 tokens)
- Generates certain tokens (e.g., "<|*end of sequence*|>")

Problem of bs = 1

b=1

Latency = step latency * # steps

Speculative decoding reduces this, hence amortize the
memory moving cost (but it may increase compute cost)

# Large Language Models

b=1

- Transformers, Attentions
- Scaling Law
  - MoE
- Connecting the dots: Training Optimizations
  - Flash attention ← come back to this later next week
- Serving and inference optimization
  - Continuous batching and Paged attention
  - Speculative decoding (Guest Lecture)
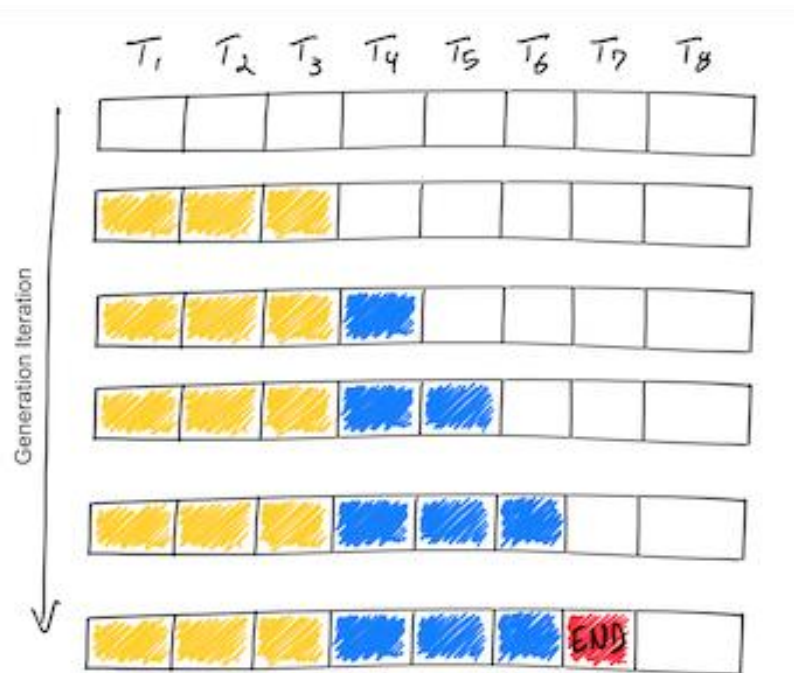- Connecting the dots: Deepseek-v3
- Hot topics

# Large Language Models

- Transformers, Attentions
- Scaling Law
  - MoE
- Connecting the dots: Training Optimizations
  - Flash attention ← come back to this later next week
- Serving and inference optimization
  - **Continuous batching and Paged attention**
  - Speculative decoding (Guest Lecture)
- Connecting the dots: Deepseek-v3
- Hot topics

# LLM Decoding Timeline

# Batching Requests to Improve GPU Performance



Issues with static batching:

- Requests may complete at different iterations
- Idle GPU cycles
- New requests cannot start immediately

# Continuous Batching



Benefits:

- Higher GPU utilization
- New requests can start immediately

Orca: A Distributed Serving System for Transformer-Based Generative Models. OSDI'22

# Continuous Batching Step-by-Step

- Receives two new requests R1 and R2

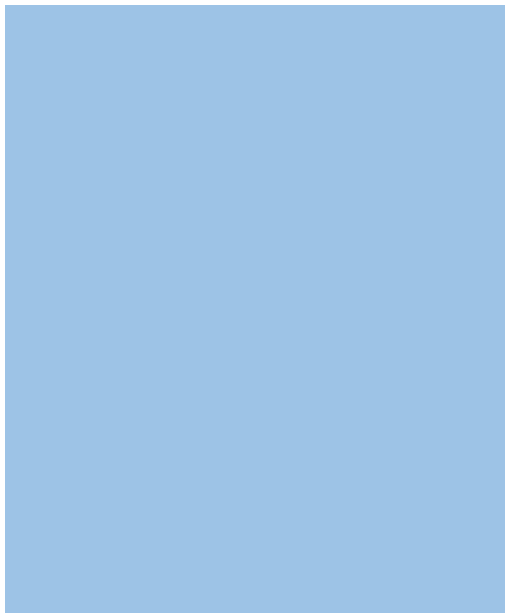R1: optimizing ML systems

R2: LLM serving is

Maximum serving batch size = 3

**Request Pool (CPU)**

**Execution Engine (GPU)**

# Continuous Batching Step-by-Step

- Iteration 1: decode R1 and R2

Maximum serving batch size = 3

R1: optimizing ML systems

R2: LLM serving is

Iteration 1

**Request Pool (CPU)**

**Execution Engine (GPU)**

21

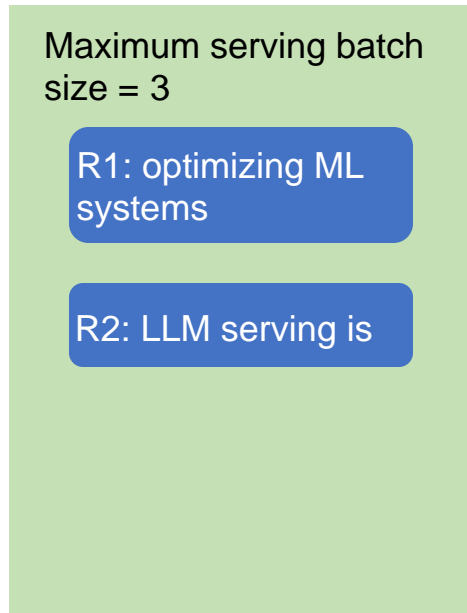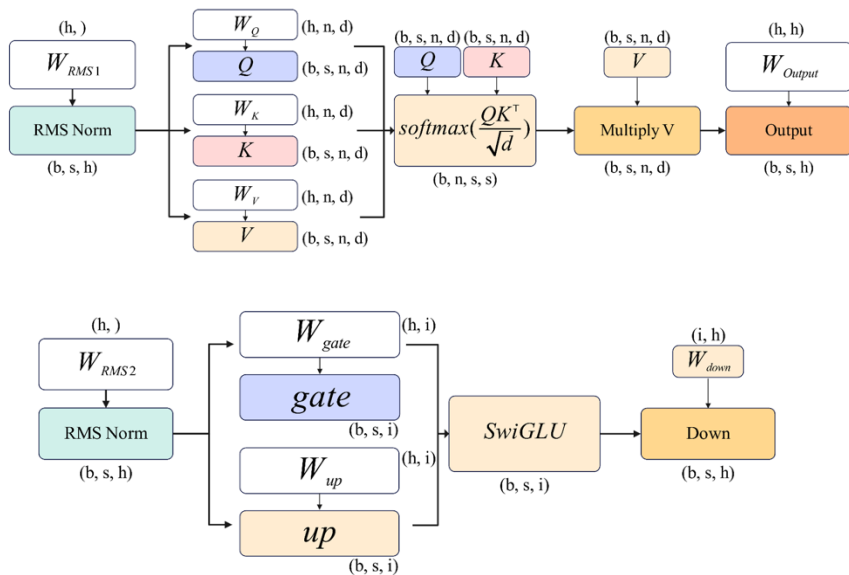# Continuous Batching Step-by-Step

- Iteration 1: decode R1 and R2



Q: How to batch these?

Maximum serving batch size = 3

R1: optimizing ML systems
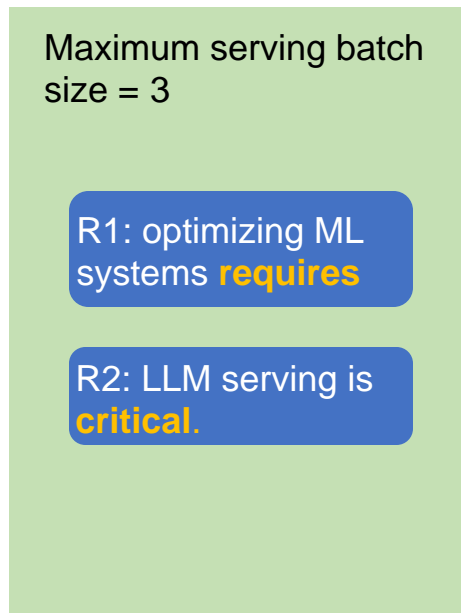
R2: LLM serving is

Iteration 1

**Execution Engine (GPU)**
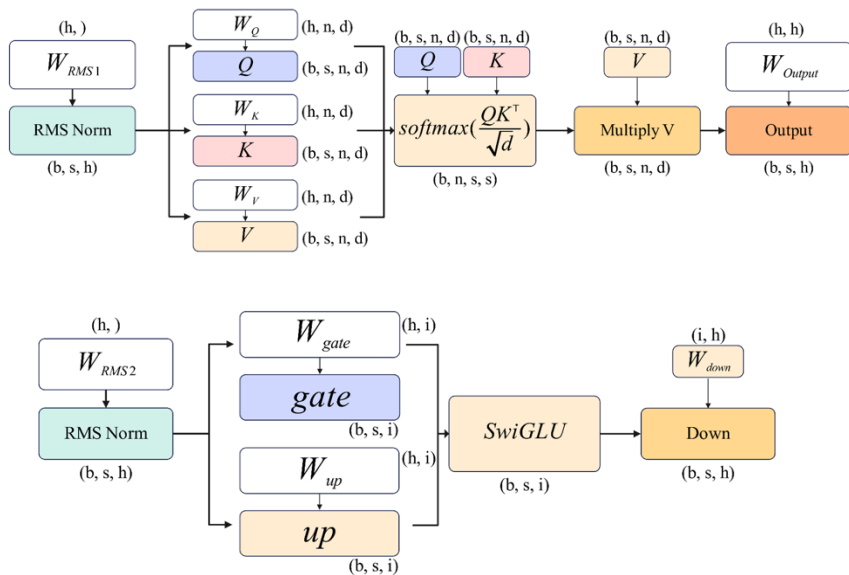
# Continuous Batching Step-by-Step

- Receive a new request R3; finish decoding R1 and R2

R3: A man

Maximum serving batch size = 3

R1: optimizing ML systems **requires**

R2: LLM serving is **critical**.

Iteration 1

**Request Pool
(CPU)**

**Execution Engine
(GPU)**

# Continuous Batching Step-by-Step

**Q: How to batch these?**

- Receive a new request R3; finish decoding R1 and R2



Maximum serving batch size = 3

R1: optimizing ML systems **requires**

R2: LLM serving is **critical.**
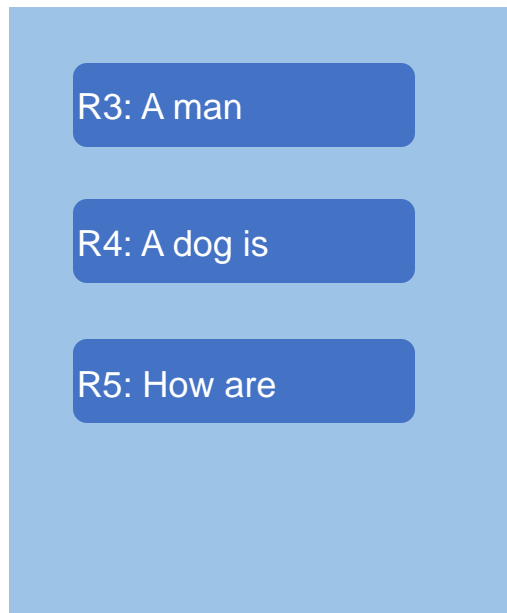
Iteration 1

**Execution Engine (GPU)**
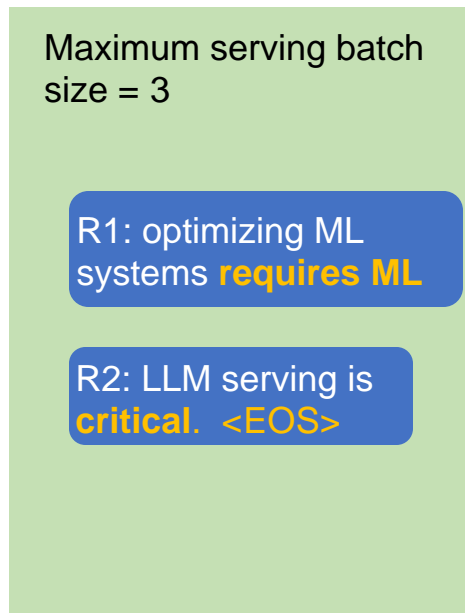
# Traditional Batching

- Receive a new request R3; finish decoding R1 and R2
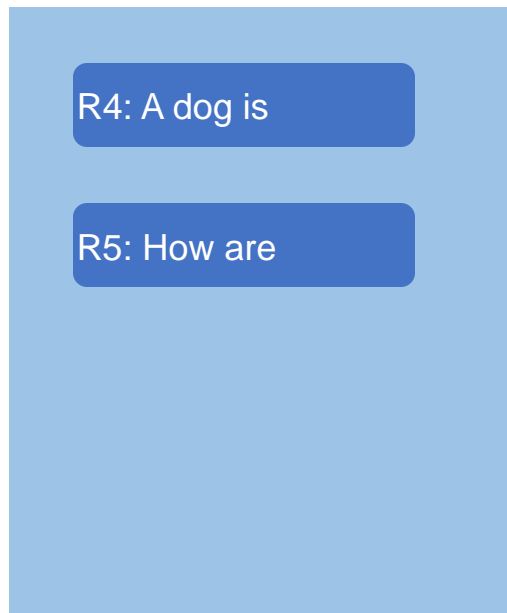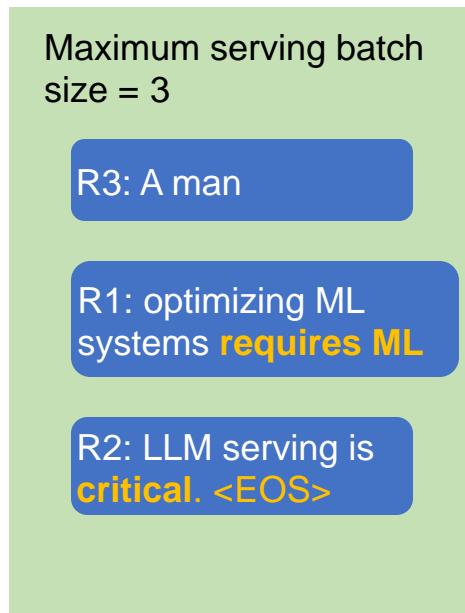
**Request Pool (CPU)**

R3: A man

R4: A dog is

R5: How are

**Execution Engine (GPU)**

Maximum serving batch size = 3

R1: optimizing ML systems **requires ML**

R2: LLM serving is **critical**.  <EOS>

Iteration 2

25

# Continuous Batching

- Iteration 2: decode R1, R2, R3; receive R4, R5; R2 completes

**Request Pool (CPU)**

R4: A dog is

R5: How are

**Execution Engine (GPU)**

Maximum serving batch size = 3

R3: A man

R1: optimizing ML systems **requires ML**

R2: LLM serving is **critical**. <EOS>
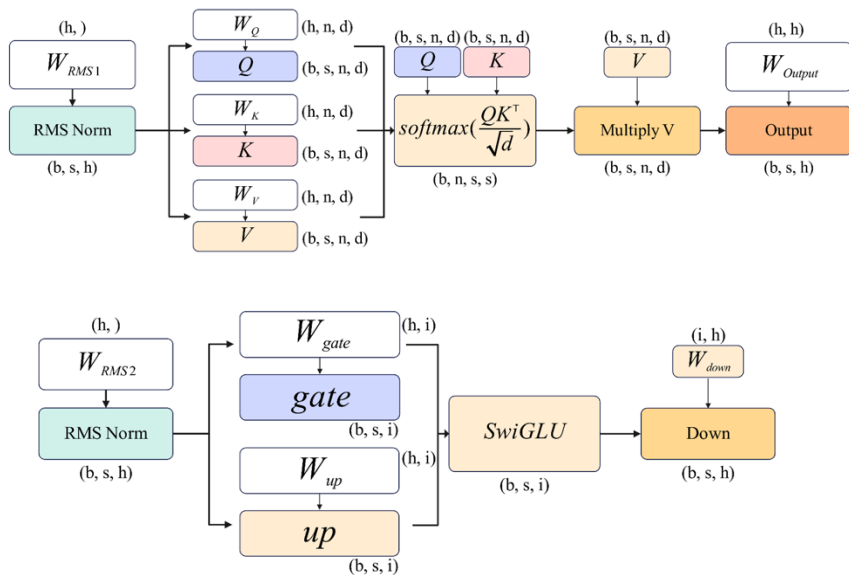
Iteration 2

# Continuous Batching

- Iteration 2: decode R1, R2, R3; receive R4, R5; R2 completes



Maximum serving batch size = 3

R3: A man

R1: optimizing ML systems **requires** **ML**

R2: LLM serving is **critical** <EOS>
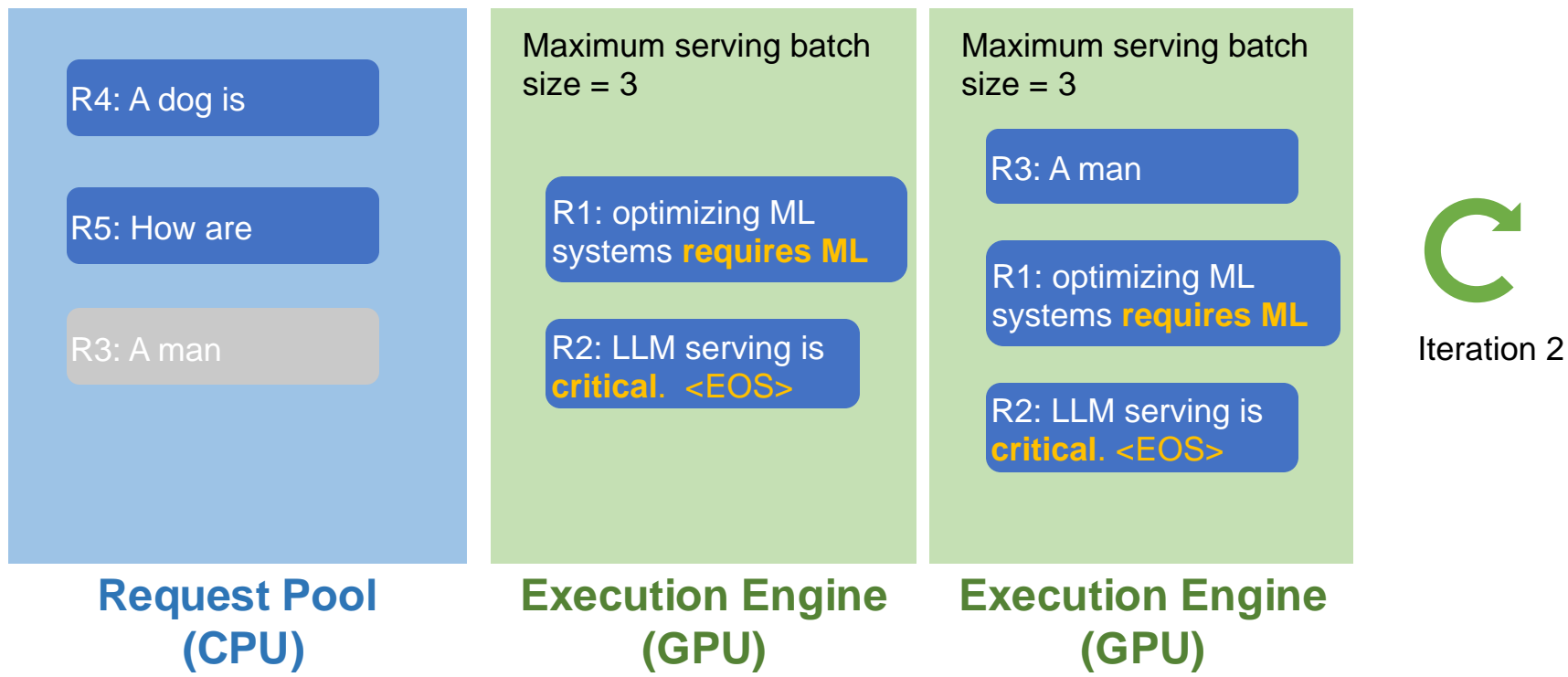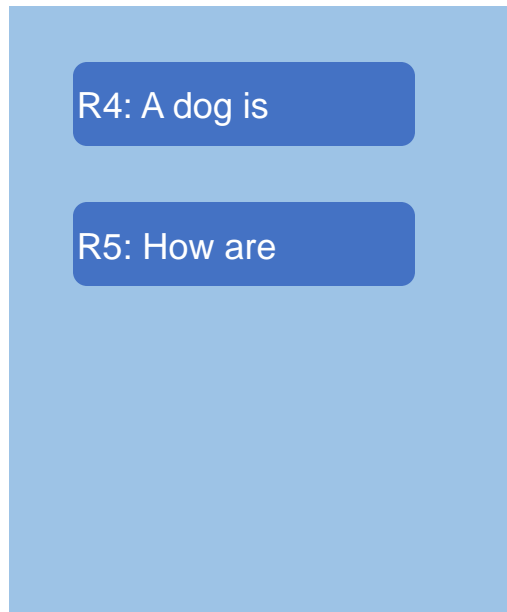
Iteration 2

**Execution Engine (GPU)**

27

# Traditional vs. Continuous Batching

- Iteration 2: decode R1, R2, R3; receive R4, R5; R2 completes



**Request Pool (CPU)**

R4: A dog is

R5: How are

R3: A man

**Execution Engine (GPU)**

Maximum serving batch size = 3

R1: optimizing ML systems **requires ML**

R2: LLM serving is **critical**. <EOS>

**Execution Engine (GPU)**

Maximum serving batch size = 3

R3: A man

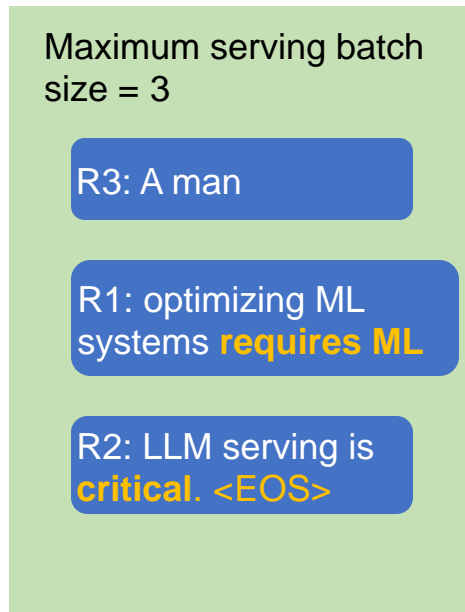R1: optimizing ML systems **requires ML**

R2: LLM serving is **critical**. <EOS>

Iteration 2

28

# Continuous Batching

- Iteration 2: decode R1, R2, R3; receive R4, R5; R2 completes

**Request Pool**
**(CPU)**

R4: A dog is

R5: How are

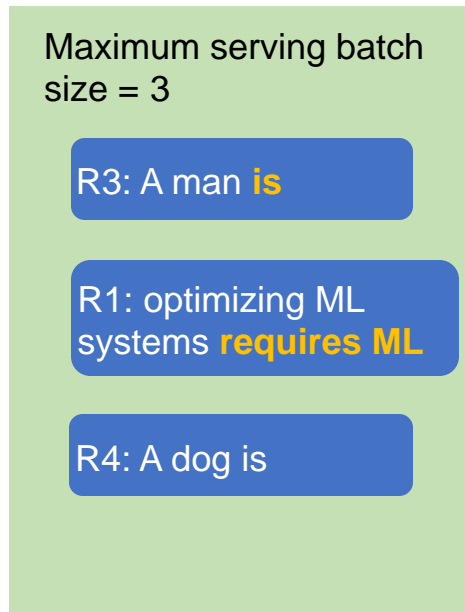**Execution Engine**
**(GPU)**

Maximum serving batch size = 3

R3: A man

R1: optimizing ML systems **requires ML**

R2: LLM serving is **critical**. <EOS>

Iteration 2

# Continuous Batching Step-by-Step

- Iteration 3: decode R1, R3, R4



Maximum serving batch size = 3

R3: A man **is**

R1: optimizing ML systems **requires ML**

R4: A dog is

Iteration 3
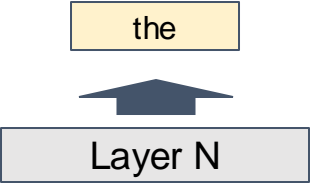
R5: How are

**Request Pool (CPU)**

**Execution Engine (GPU)**

# Summary: Continuous Batching

- Handle early-finished and late-arrived requests more efficiently

- Improve GPU utilization

- Key observation
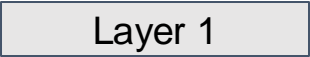
  - MLP kernels are agnostic to the sequence dimension

# KV Cache

# KV Cache

Output

| of |

| | Layer N |

Layer N

| | | | future | 0.1 | -2.1 | 0.5 |

KV Cache

| Artificial | -0.2 | 0.1 | -1.1 |
| Intelligence | 0.9 | 0.7 | 0.2 |
| is | -0.1 | -0.3 | 0.1 |
| **the** | -1.1 | 0.5 | 0.4 |

Layer 1

| future | -0.6 | 0.0 | 0.9 |

| Artificial | -0.1 | 0.3 | 1.2 |
| Intelligence | 0.7 | -0.4 | 0.8 |
| is | 0.2 | -0.1 | 1.1 |
| **the** | -0.7 | 0.1 | -0.2 |

Input

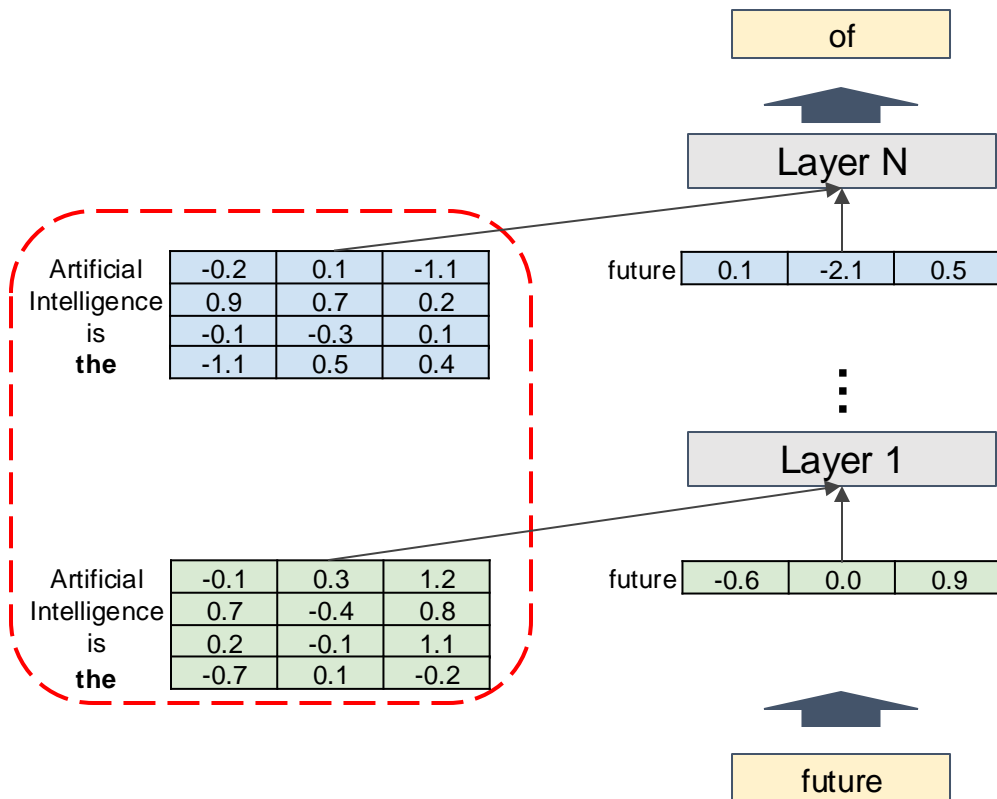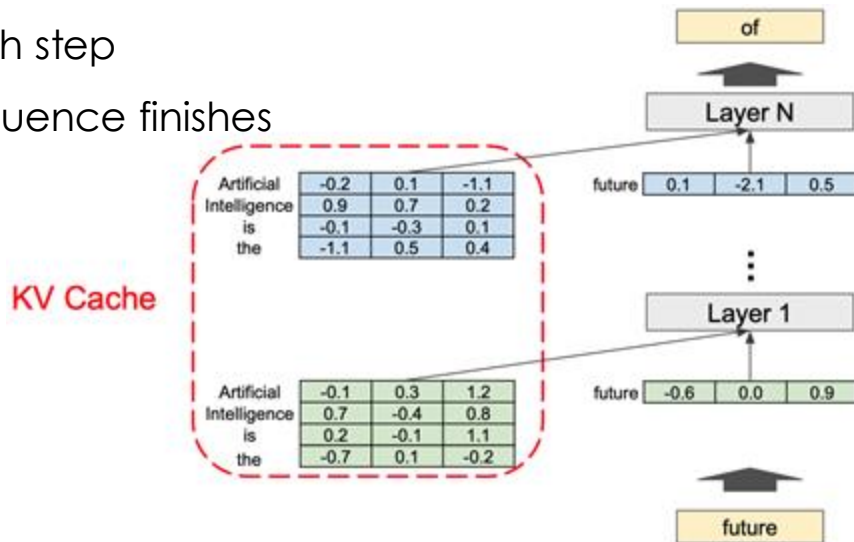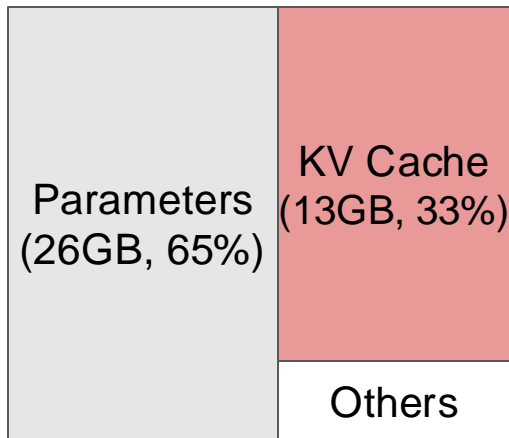| future |

# KV Cache

- Memory space to store intermediate vector representations of tokens
  - **Working set** rather than a "cache"
- The size of KV Cache dynamically grows and shrinks
  - A new token is appended in each step
  - Tokens are deleted once the sequence finishes

# Key insight

Efficient management of KV cache is crucial for high-throughput LLM serving



KV Cache (13GB, 33%)

Parameters (26GB, 65%)

Others

13B LLM on A100-40GB

Existing systems — vLLM

Memory usage (GB)
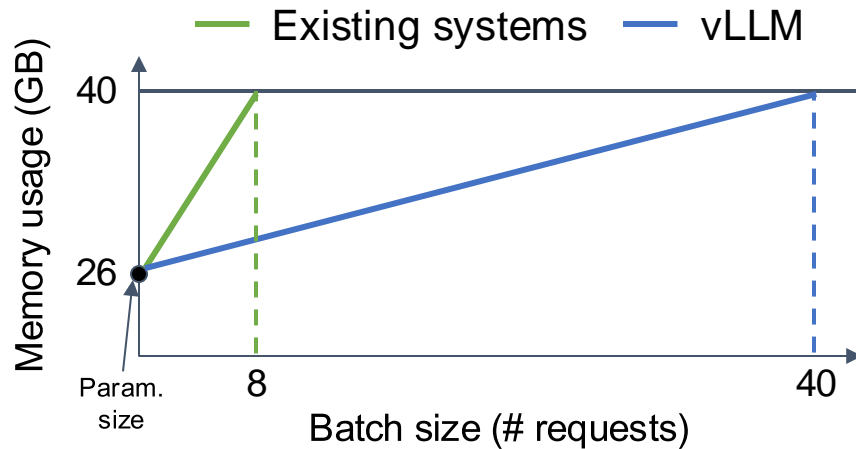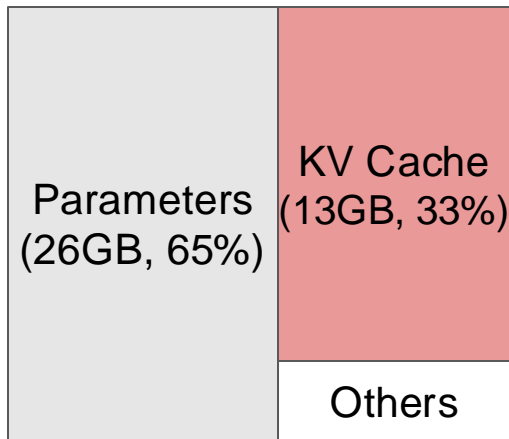
40

26

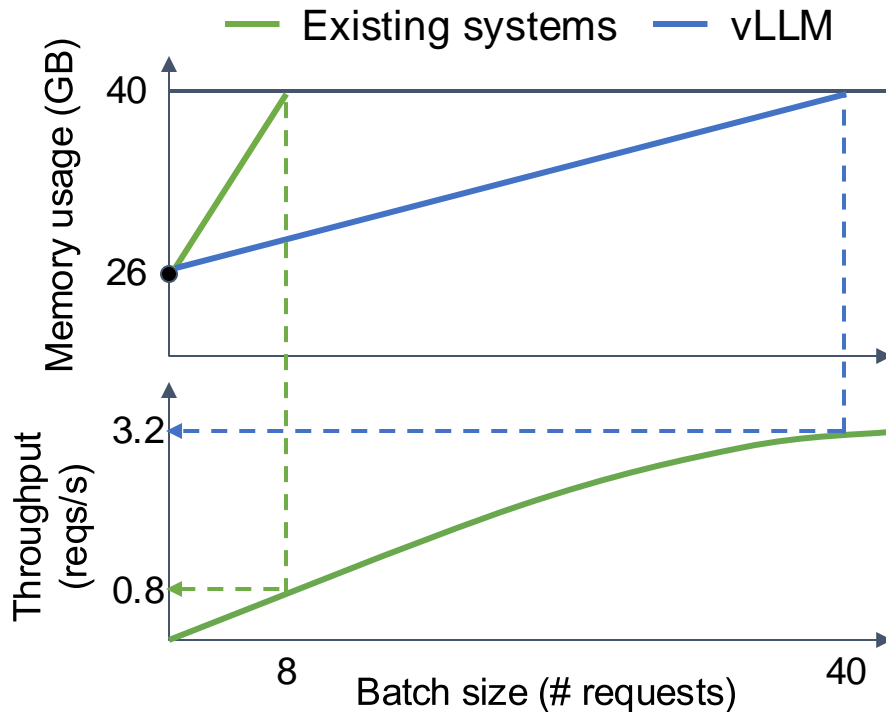Param. size

8

40

Batch size (# requests)

# Key insight

Efficient management of KV cache is crucial for high-throughput LLM serving



13B LLM on A100-40GB

# Memory waste in KV Cache



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

# Memory waste in KV Cache



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

* Yu, G. I., Jeong, J. S., Kim, G. W., Kim, S., Chun, B. G. "Orca: A Distributed Serving System for Transformer-Based Generative Models" (OSDI 22).

# vLLM: Efficient memory management for LLM inference

Inspired by **virtual memory** and **paging**

# Token block

- A **fixed-size** contiguous chunk of memory that can store token states **from left to right**

Token blocks
(KV Cache)

block 0

block 1

block 2

block 3

KV Cache

block 4

block 5

block 6

block 7

Block size = 4
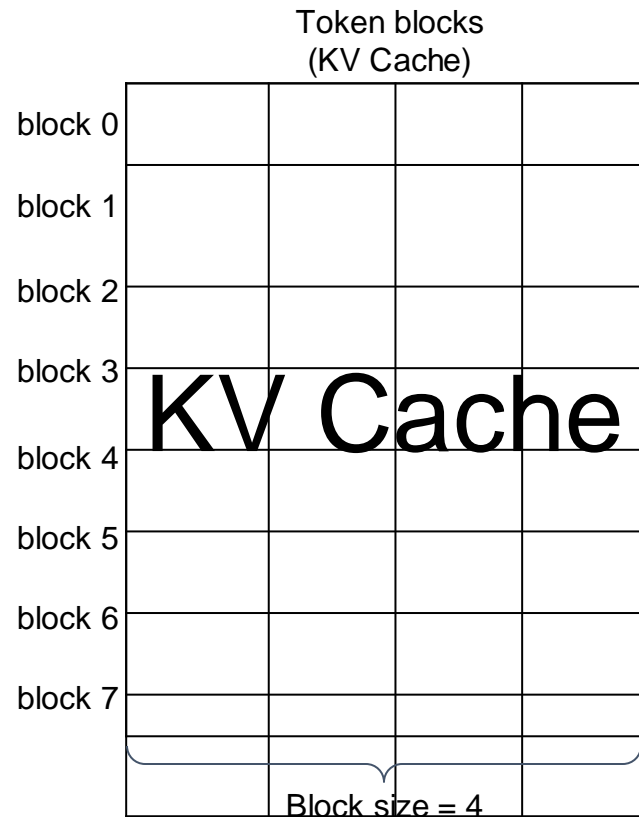
# Token block

- A **fixed-size** contiguous chunk of memory that can store token states **from left to right**

Token blocks
(KV Cache)

| | block 0 | | | |
|---|---|---|---|---|
| | block 1 | | | |
| | block 2 | | | |
| | block 3 | | | |
| | block 4 | | | |
| block 5 | Artificial | Intelligence | is | the |
| | block 6 | | | |
| | block 7 | | | |

Block size = 4

Block 4

| | | | |
|---|---|---|---|
| Artificial | -0.2 | 0.1 | -1.1 |
| Intelligence | 0.9 | 0.7 | 0.2 |
| is | -0.1 | -0.3 | 0.1 |
| the | -1.1 | 0.5 | 0.4 |

820 KB / token
(LLaMA-13B)

# Paged Attention

- An attention algorithm that allows for storing continuous keys and values in non-contiguous memory space

# Logical & physical token blocks

Request A

Prompt: "Alan Turing is a computer scientist"

**Logical** token blocks

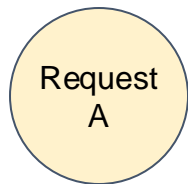| | | | |
|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |

**Physical** token blocks
(KV Cache)

| | | | |
|---|---|---|---|
| block 0 | | | | |
| block 1 | | | | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | | | | |

# Logical & physical token blocks

Request A

Prompt: "Alan Turing is a computer scientist"

**Physical** token blocks (KV Cache)

**Logical** token blocks

| | | | |
|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 2 |
| – | – |
| – | – |

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

# Logical & physical token blocks

Request
A

Prompt: "Alan Turing is a computer scientist"
Completion: "and"

**Logical** token blocks

| | | | |
|---|---|---|---|
| Alan | Turing | is | a |
| computer | scientist | | |
| | | | |
| | | | |

block 0
block 1
block 2
block 3

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 2 |
| – | – |
| – | – |

**Physical** token blocks
(KV Cache)

| | | | |
|---|---|---|---|
| | | | |
| computer | scientist | | |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |
| Alan | Turing | is | a |

block 0
block 1
block 2
block 3
block 4
block 5
block 6
block 7

45

# Logical & physical token blocks

**Request A**

Prompt: "Alan Turing is a computer scientist"
Completion: "<u>and</u>"

**Logical** token blocks

| | | | |
|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | and | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 2 |
| – | – |
| – | – |

**Physical** token blocks
(KV Cache)

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | | |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

# Logical & physical token blocks

Request A

Prompt: "Alan Turing is a computer scientist"
Completion: "_and_"

**Physical** token blocks
(KV Cache)

| | | | |
|---|---|---|---|
| block 0 | | | |
| block 1 | computer | scientist | and |
| block 2 | | | |
| block 3 | | | |
| block 4 | | | |
| block 5 | | | |
| block 6 | | | |
| block 7 | Alan | Turing | is | a |

**Logical** token blocks

| | | | | |
|---|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | and | |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 3 |
| – | – |
| – | – |

# Logical & physical token blocks

**Request A**

Prompt: "Alan Turing is a computer scientist"
Completion: "and mathematician"

**Physical** token blocks
(KV Cache)

| | | | |
|---|---|---|---|
| block 0 | | | |
| block 1 computer | scientist | and | mathematician |
| block 2 | | | |
| block 3 | | | |
| block 4 | | | |
| block 5 | | | |
| block 6 | | | |
| block 7 Alan | Turing | is | a |

**Logical** token blocks

| | | | | |
|---|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | and | mathematician |
| block 2 | | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 4 |
| – | – |
| – | – |

# Logical & physical token blocks

**Physical** token blocks
(KV Cache)

Request A

Prompt: "Alan Turing is a computer scientist"
Completion: "and mathematician <u>renowned</u>"

**Logical** token blocks

| | | | | |
|---|---|---|---|---|
| block 0 | Alan | Turing | is | a |
| block 1 | computer | scientist | and | mathematician |
| block 2 | renowned | | | |
| block 3 | | | | |

**Block table**

| Physical block number | # Filled |
|---|---|
| 7 | 4 |
| 1 | 4 |
| 5 | 1 |
| – | – |

| | | | | |
|---|---|---|---|---|
| block 0 | | | | |
| block 1 | computer | scientist | and | mathematician |
| block 2 | | | | |
| block 3 | | | | |
| block 4 | | | | |
| block 5 | renowned | | | |
| block 6 | | | | |
| block 7 | Alan | Turing | is | a |

Allocated on demand

# Serving multiple requests

**Physical** token blocks
(KV Cache)

| | | | |
|---|---|---|---|
| computer | scientist | and | mathematician |
| | | | |
| Artificial | Intelligence | is | the |
| | | | |
| renowned | | | |
| future | of | technology | |
| Alan | Turing | is | a |

**Request A**

**Block Table**

| | |
|---|---|
| | |
| | |
| | |

**Logical** token blocks

| | | | |
|---|---|---|---|
| Alan | Turing | is | a |
| computer | scientist | and | mathematician |
| renowned | | | |
| | | | |

**Request B**

**Block Table**

| | |
|---|---|
| | |
| | |
| | |

**Logical** token blocks

| | | | |
|---|---|---|---|
| Artificial | Intelligence | is | the |
| future | of | technology | |
| | | | |
| | | | |

# Memory efficiency of vLLM

- Minimal internal fragmentation
  - Only happens at the last block of a sequence
  - **# wasted tokens / seq < block size**
    - Sequence: O(100) – O(1000) tokens
    - Block size: 16 or 32 tokens
- No external fragmentation

| Alan | Turing | is | a |
|------|--------|-----|-------------|
| computer | scientist | and | mathematician |
| renowned | | | |

Internal fragmentation

# Effectiveness of PagedAttention



96.3% KV cache utilization

# Large Language Models

- Transformers, Attentions
- Scaling Law
  - MoE
- Connecting the dots: Training Optimizations
  - Flash attention ← come back to this later next week
- Serving and inference optimization
  - **Continuous batching and Paged attention**
  - Speculative decoding (Guest Lecture)
- Connecting the dots: Deepseek-v3
- Hot topics
  - Prefill-decode disaggregation

# LLM System Today Optimize **Throughput**

# Motivation: Applications have Diverse SLO

- **TTFT**

  Time to first token
  Initial response time

- **TPOT**

  Time per output token
  Average time between two subsequent generated tokens

**Chatbot**
Fast initial response

Human reading speed (P99 latency = 250ms)

**Summarization**
User can tolerate longer initial response

Data output generation (P99 latency = 35ms)

# High Throughput ≠ High Goodput



TTFT

Suppose 10 requests complete within 1 second...

TPOT

**Throughput = 10 rps**
= completed request / time

**High Throughput**
System

...

# High Throughput ≠ High Goodput



TTFT

TPOT

Suppose 10 requests complete within 1 second…

**Throughput = 10 rps**

= completed request / time

under SLO criteria

**Goodput = 3 rps** 😮‍💨

= completed request **within SLO** / time

TTFT

200ms

50ms

TPOT

… but only 3 (out of 10) hold the latency target

**High Throughput**
System

…

can have
**Low Goodput!**

# High Throughput ≠ High Goodput



High Throughput

⇒ Poor UX 💔

**High Throughput** can still have **Low Goodput**

Goodput = 3 rps
= completed request **within SLO** / time

… but only 3 (out of 10) hold the latency target

# Background: Continuous Batching

Disaggregation is a technique that

## Request Arrived

Worker

GPU

Request

Timeline

# Prefill and Decode have Distinct Characteristics

- ## Prefill

  **Compute-bound**
  One prefill saturates compute.

- ## Decode

  **Memory-bound**
  Must batch a lot of requests together to saturate compute

# Continuous Batching Cause Interference

## Continuous Batching
Batch R1 and R2 together in 1 GPU

## Separate prefill / decode
R1 and R2 in separate GPUs

Time wasted for decode

R1

R2

time

Request arrival

R2 arrives

Time wasted for prefill

GPU

GPU

R1

R2

time

Request arrival

R2 arrives

No Interference

wasted time

# Continuous Batching Cause Interference

## Continuous Batching
### Batch R1 and R2 together in 1 GPU

Time wasted for decode

GPU

R1

R2

time

Request
arrival

R2 arrives

Time wasted for prefill

## Continuous Batching
### Batch R1~R4 together in 1 GPU

R1

R2

R3

R4

Request
arrival

R2     R3     R4

time

wasted time

# Colocation → Overprovision Resource to meet SLO



Poor UX 💔

lower cost 💰

add more GPU

Good UX 🥰

Higher cost 💰💰💰💰

# Summary: Problems caused by Colocation



Continuous Batching Cause Interference

Coupled Parallelism Strategy

|  | TP | DP |
|---|---|---|
|  | ❤️ | ❤️ |
|  | 😐 | 😐 |

Is there a better way to achieve better
**Goodput per GPU**?

# Disaggregating Prefill and Decode

Disaggregation is a technique that

## Request Arrived



Timeline

# Disaggregation achieves better goodput

## Colocate

1 GPU for both Prefill and Decode



P90 TTFT <= 400 ms
Prefill = 3 rps

P90 TPOT <= 40 ms
Decode = 1.6 rps

Max System goodput
= Min(Prefill, Decode)
= 1.6 rps / GPU

# Disaggregation achieves better goodput
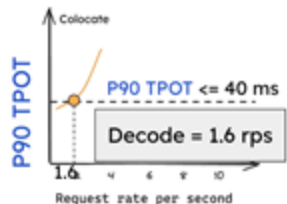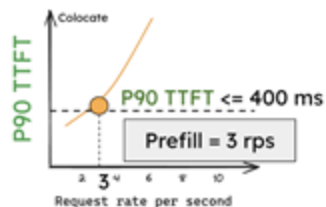
## Colocate
1 GPU for both Prefill and Decode



## Disaggregate (2P1D)
2 GPU for Prefill + 1 GPU for Decode

# Disaggregation achieves better goodput

## Colocate
1 GPU for both Prefill and Decode



## Disaggregate (2P1D)
2 GPU for Prefill + 1 GPU for Decode



Simple Disaggregation achieves **2x** goodput (per GPU)

Max System goodput
= Min(Prefill, Decode)
= 1.6 rps / GPU 😐

Disaggregate (2P1D) goodput
= Min (5.6 x **2**, 10) rps / 3 GPU
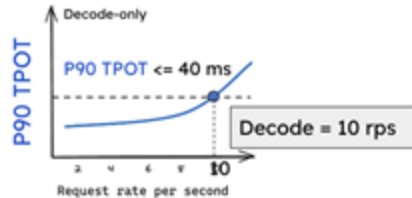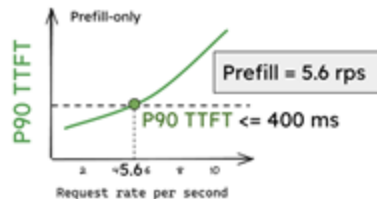= 3.3 rps / GPU 😎

# Disaggregation

- Published in 2024 at UCSD (yes, Hao's lab)
- Soon become the default architecture replacing continuous batching at large scale
- Deepseek v3 uses prefill-decode disaggregation combined with different parallelisms.