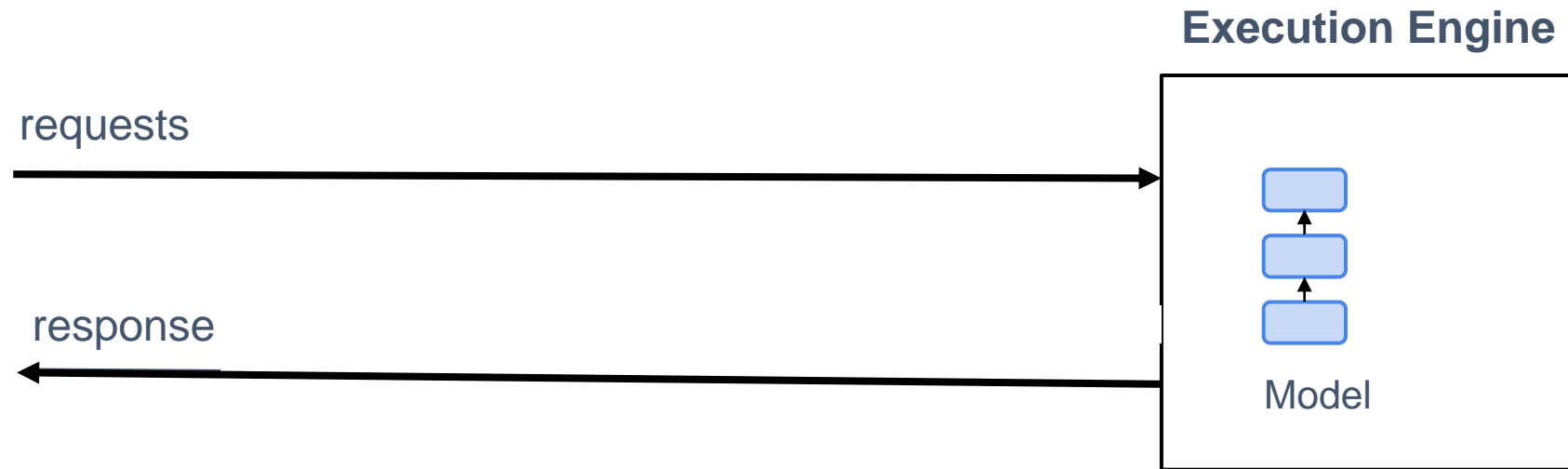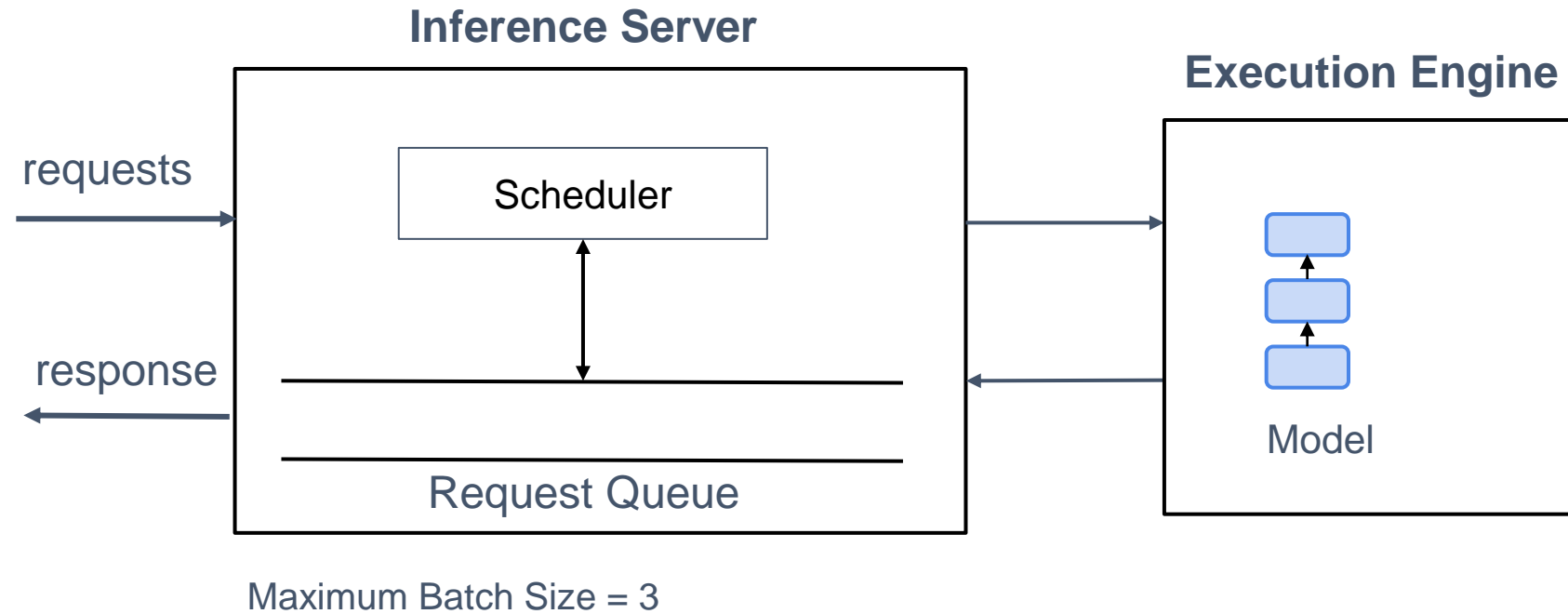# CS 498: Machine Learning System
# Spring 2025

Minjia Zhang
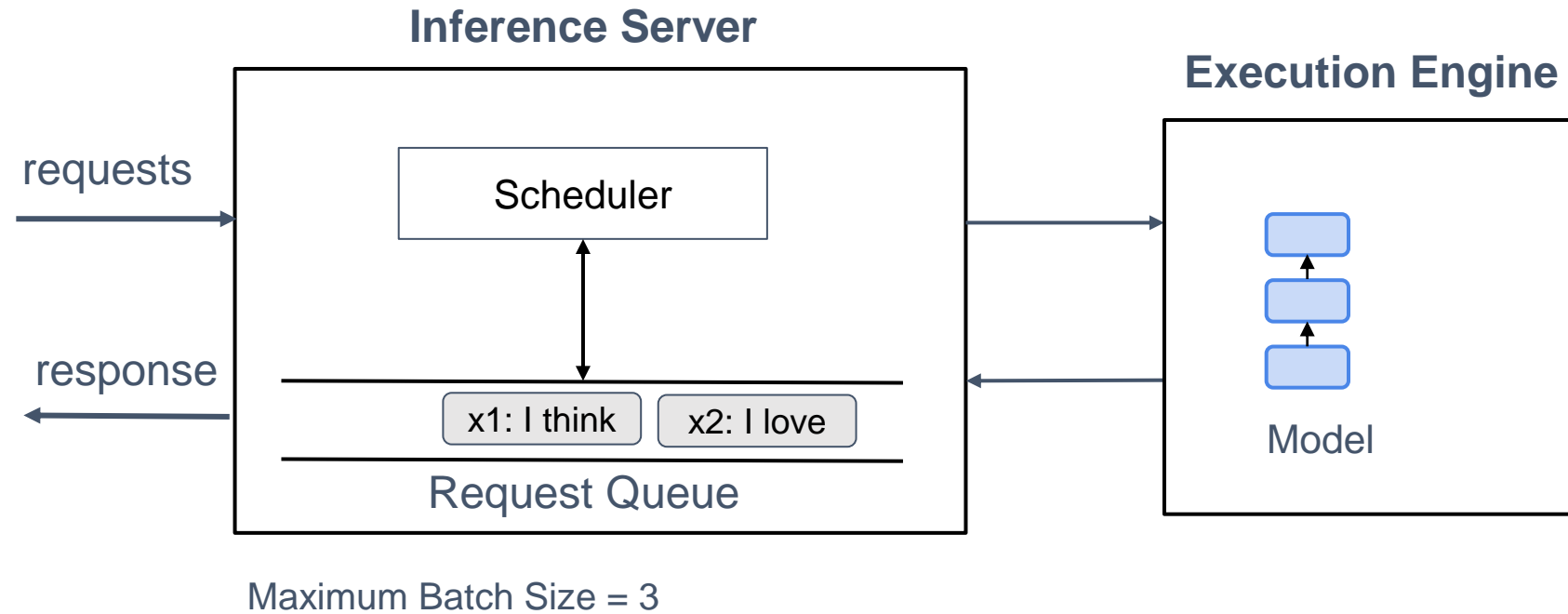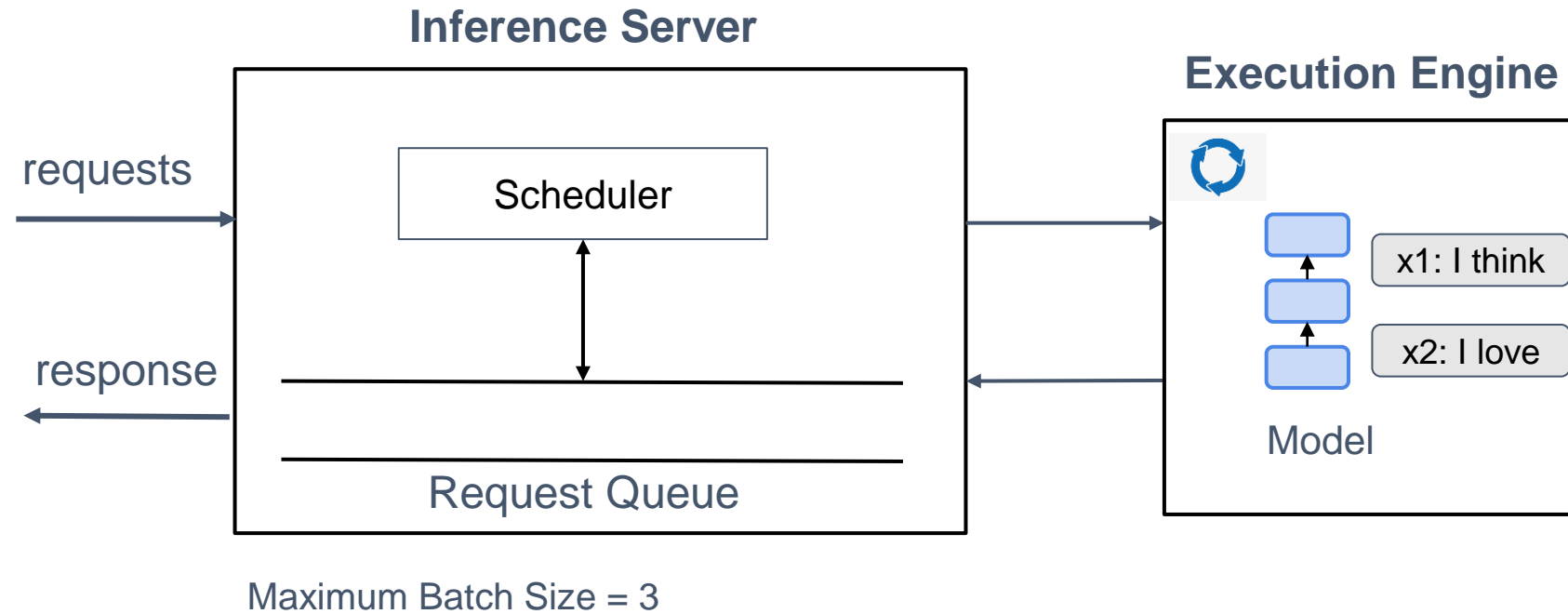
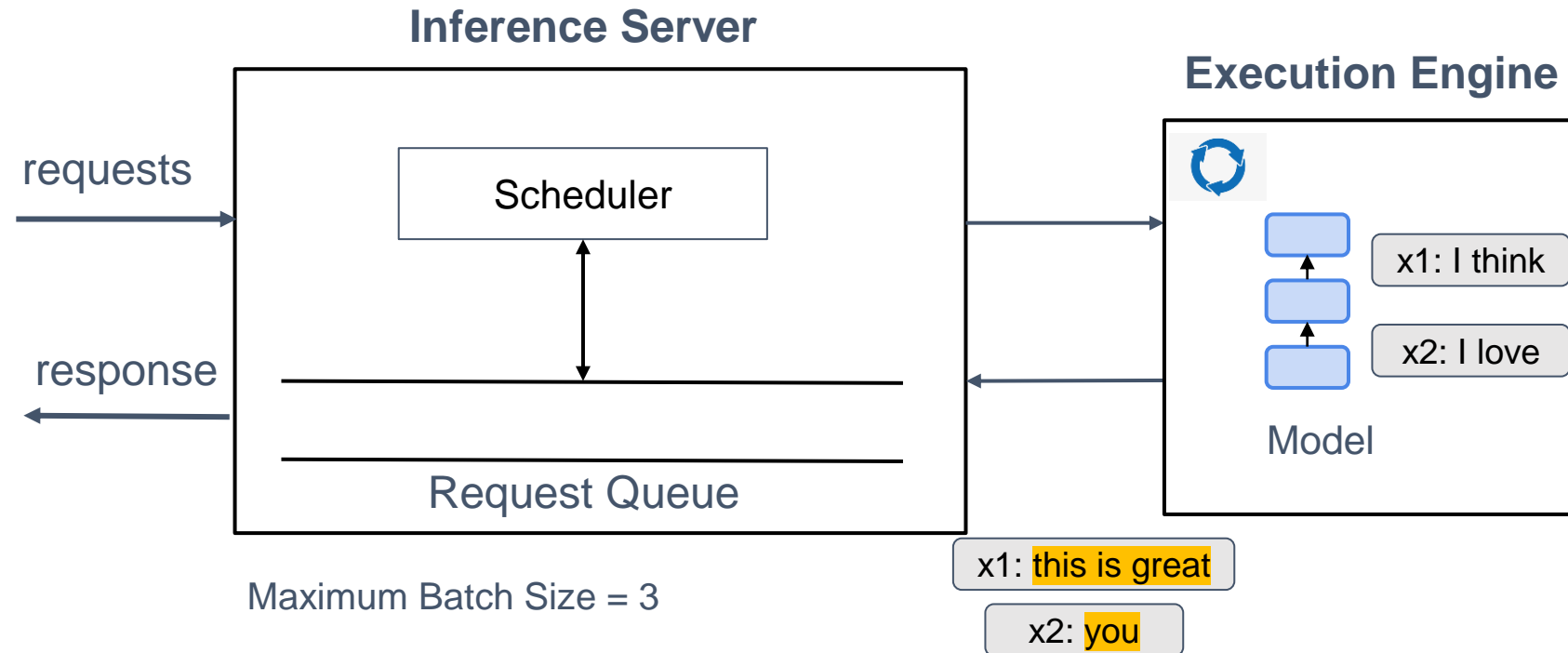The Grainger College of Engineering

DL Inference

- LLM Serving System
- Continuous Batching

# LLM Inference

**Execution Engine**

requests

response

Model

# LLM Serving System

**Inference Server**

**Execution Engine**

requests

Scheduler

response

Request Queue

Model

Maximum Batch Size = 3

# LLM Serving System

**Inference Server**

**Execution Engine**

requests

Scheduler

response

x1: I think     x2: I love

Request Queue

Model

Maximum Batch Size = 3

# LLM Serving System

**Inference Server**

**Execution Engine**

requests

Scheduler

response

Request Queue

x1: I think

x2: I love

Model

Maximum Batch Size = 3

# LLM Serving System



**Inference Server**

**Execution Engine**

requests

Scheduler

response

Request Queue

Maximum Batch Size = 3

Model

x1: I think

x2: I love

x1: this is great

x2: you

# LLM Serving System

**Inference Server**

**Execution Engine**

requests

Scheduler

response

Request Queue

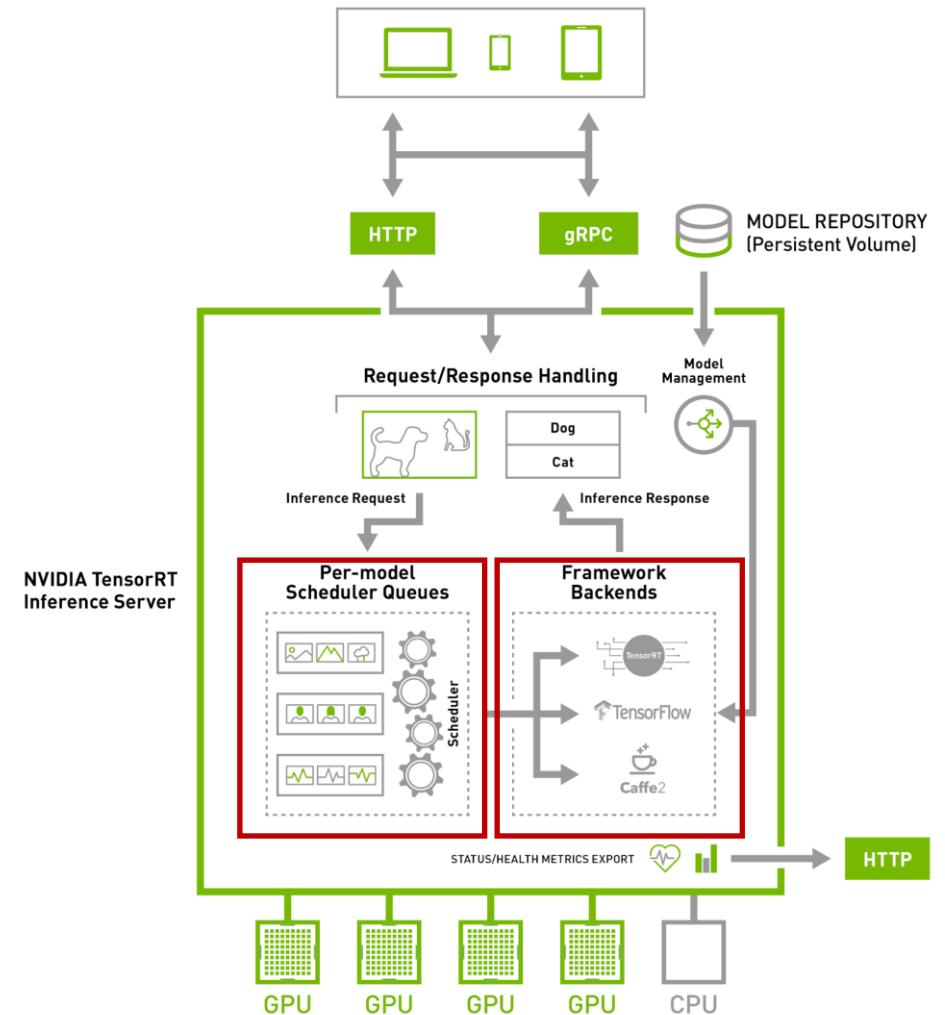Model

x1: I think this is great

X2: I love you

Maximum Batch Size = 3

- Separates implementation of serving layer and execution layer

○ Separates implementation of serving layer and execution layer
○ Implements scheduling and batching algorithms
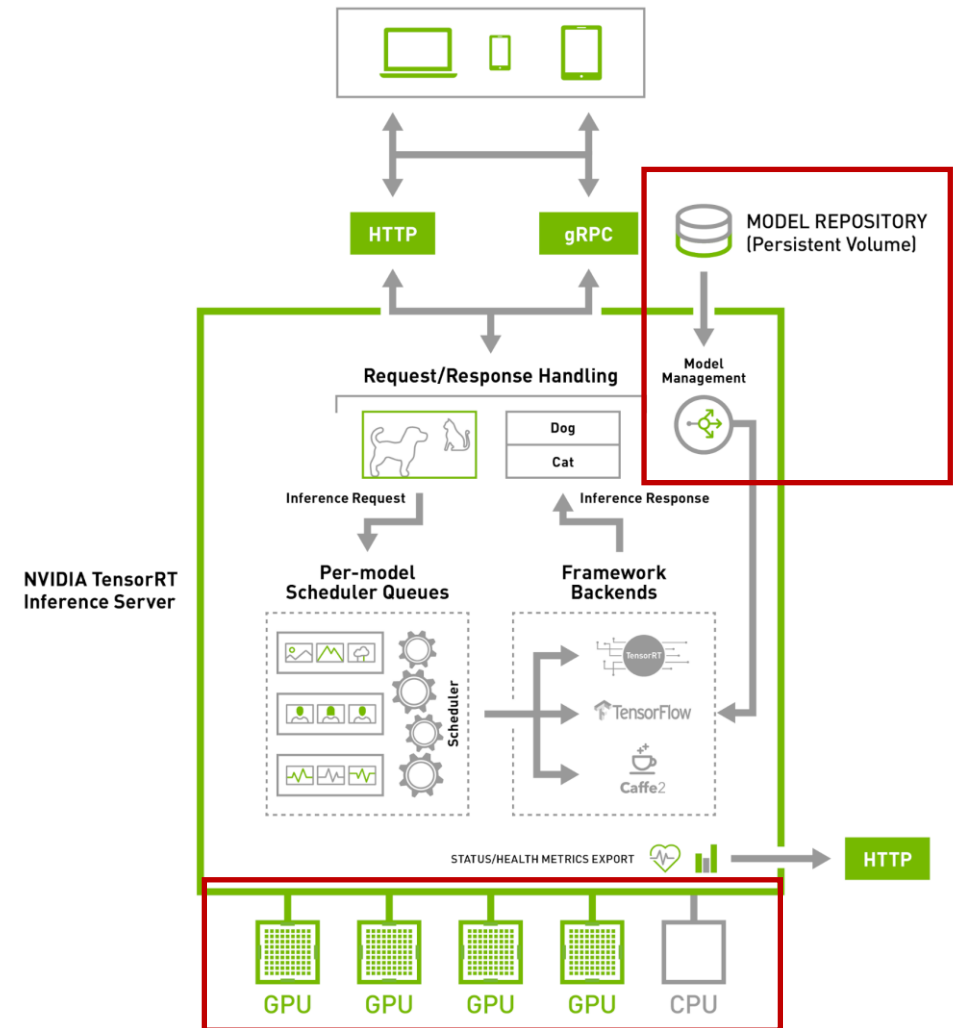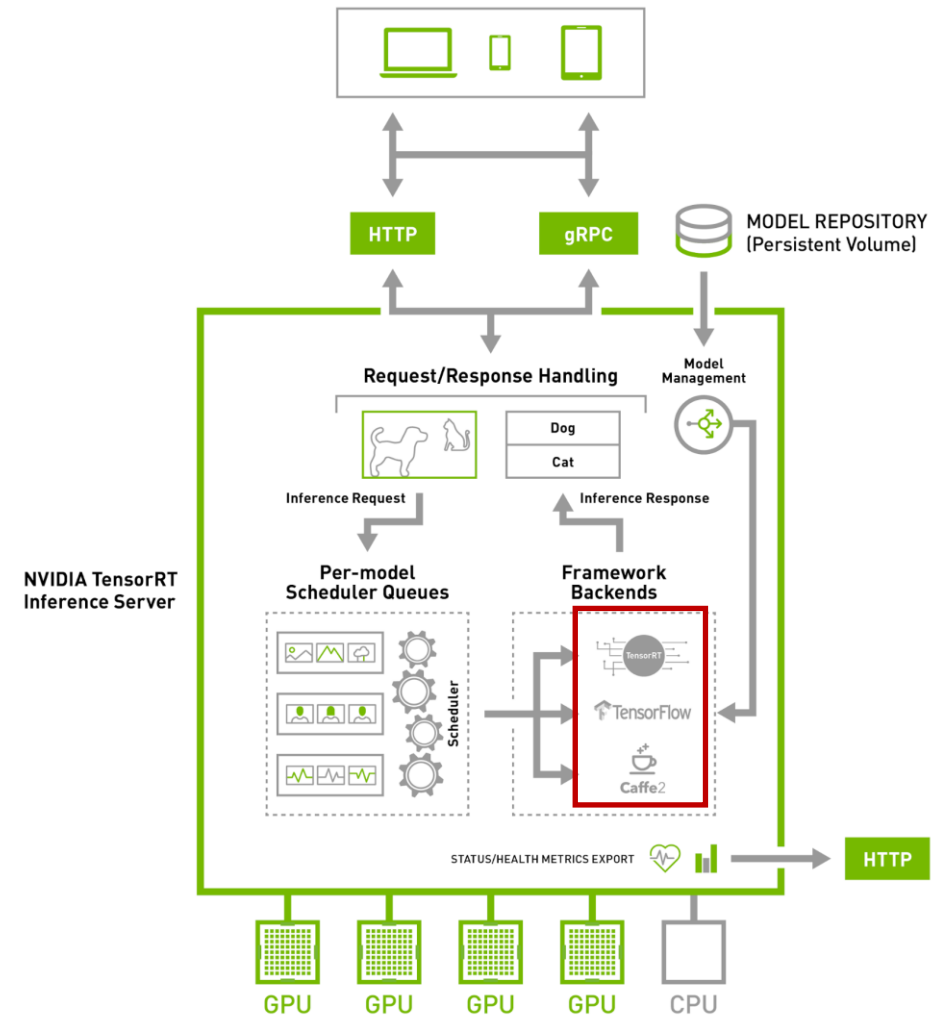  ■ Sequence Batching
  ■ Continuous Batching

○ Separates implementation of serving layer and execution layer
○ Implements scheduling and batching algorithms
- Sequence Batching
- Continuous Batching
○ Allows multiple models to concurrently execute

# Example: TensorRT Inference Server

○ Separates implementation of serving layer and execution layer
○ Implements scheduling and batching algorithms
  ▪ Sequence Batching
  ▪ Continuous Batching
○ Allows multiple models to concurrently execute
○ Supports multiple frameworks
  ▪ PyTorch
  ▪ TensorFlow
  ▪ ONNX
  ▪ vLLM backend

# DL Inference

- LLM Serving System
- Continuous Batching
  - Sequence batching
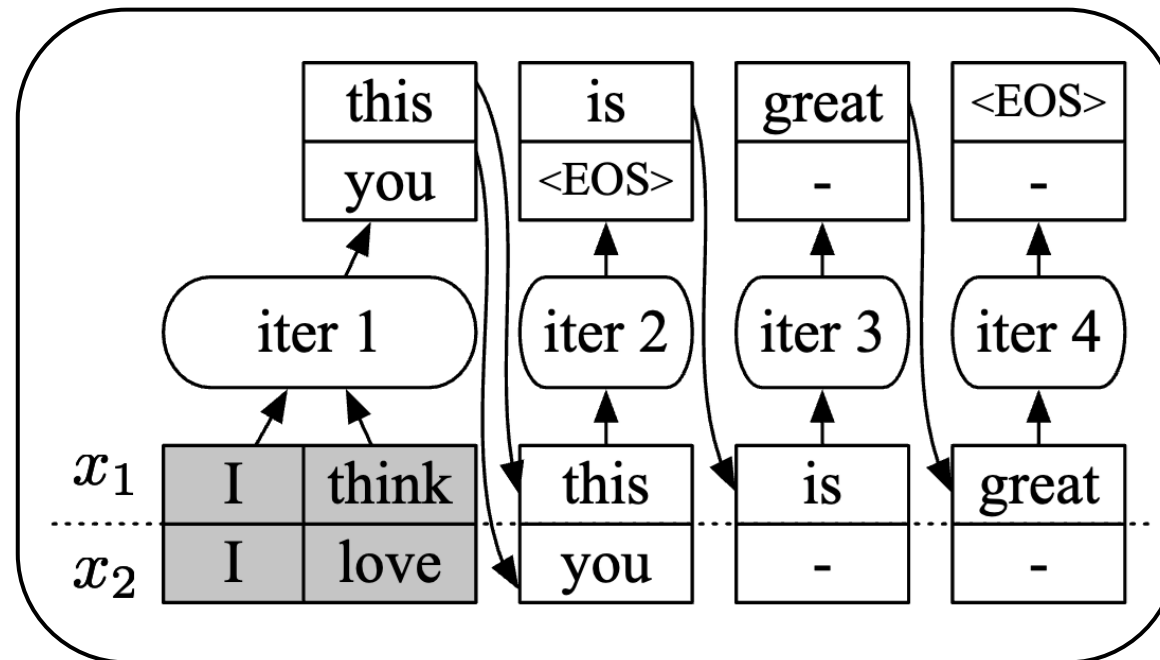  - Continuous batching

# DL Inference

- LLM Serving System

- Continuous Batching
  - Sequence batching
  - Continuous batching

Question: Can we use the batching scheme (sequence batching) during training for inference?

**Execution Engine**

**Execution Engine**



Do you see any problem?

**Execution Engine**

x2 generation done

**Execution Engine**

**x2 generation done**

**Early finished requests cannot return to the Inference Server → Latency increase**

# Problem 1: Request Level Scheduling

**Inference Server**

**Execution Engine**

requests

Scheduler

response

Request Queue

x1: I think

x2: I love

Model

Maximum Batch Size = 3

Started processing x1
and x2

# Problem 1: Request Level Scheduling

**Inference Server**

**Execution Engine**

requests →

response ←

Scheduler

New!

x3: A man

Request Queue

Model

x1: I think

x2: I love

Maximum Batch Size = 3

# Problem 1: Request Level Scheduling

**Inference Server**

**Execution Engine**

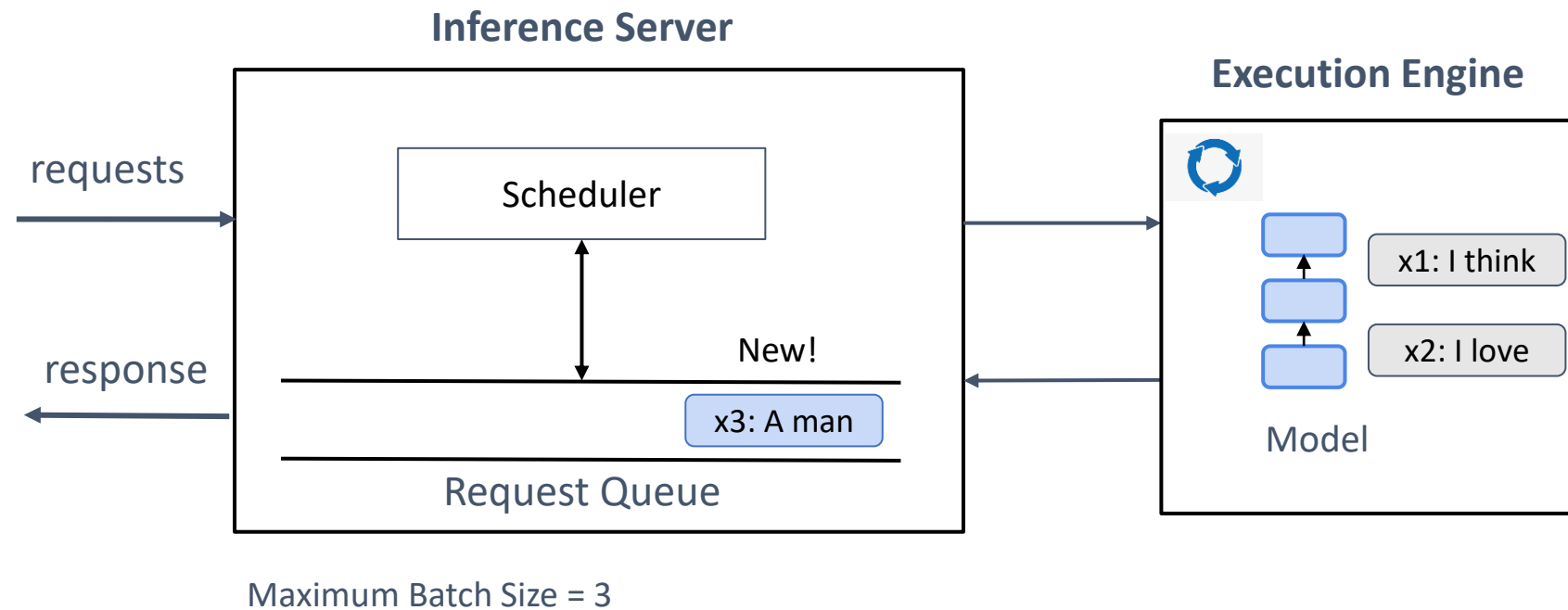requests

Scheduler

response

New!

x3: A man

Request Queue

Model

x1: I think

x2: I love

Maximum Batch Size = 3

**Late join requests need to wait until engine finishes execution**
**→ Latency Increase**

**Inference Server**

**Execution Engine**

requests

Scheduler

New!

x3: A man

Request Queue

x1: I think

x2: I love
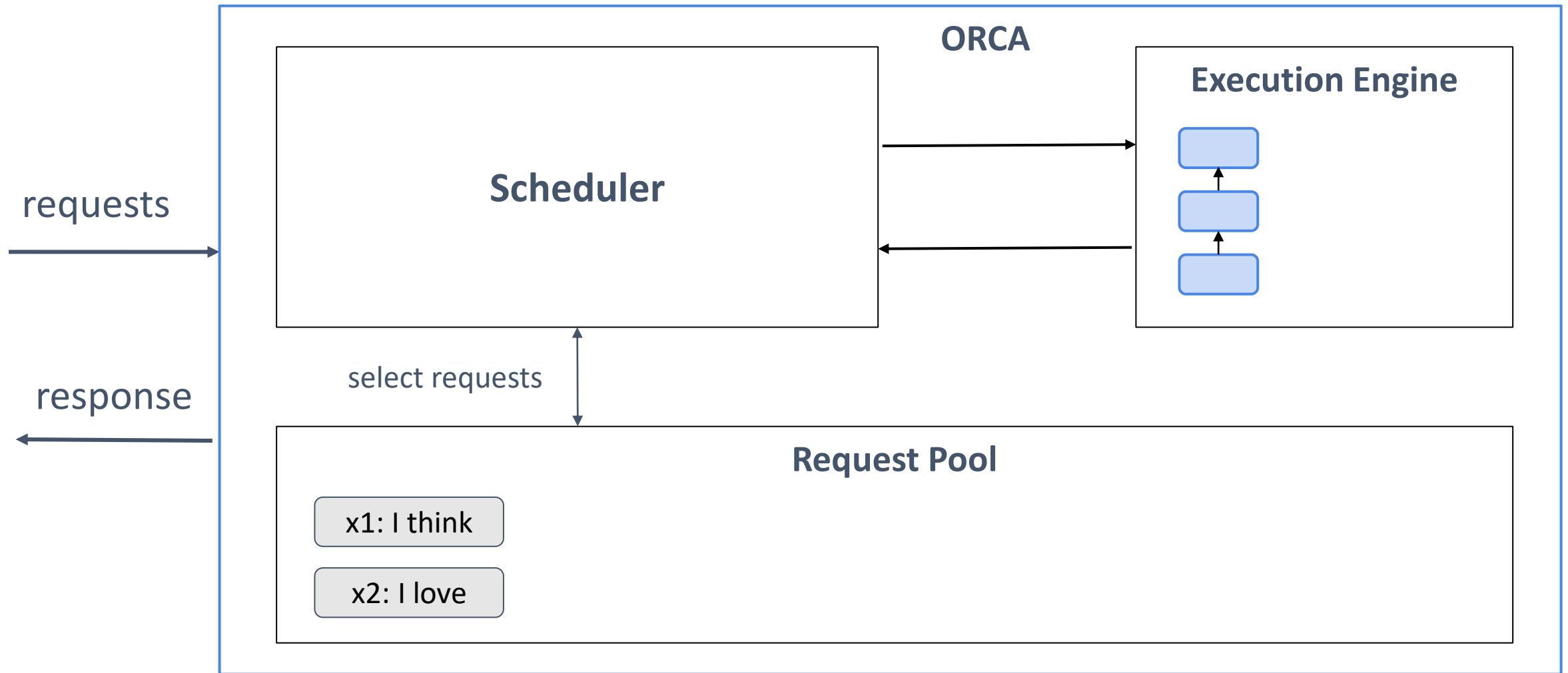
Model

Maximum Batch Size = 3

**Question: How can we avoid redundant computation and ensure late-arriving requests to be processed more promptly?**
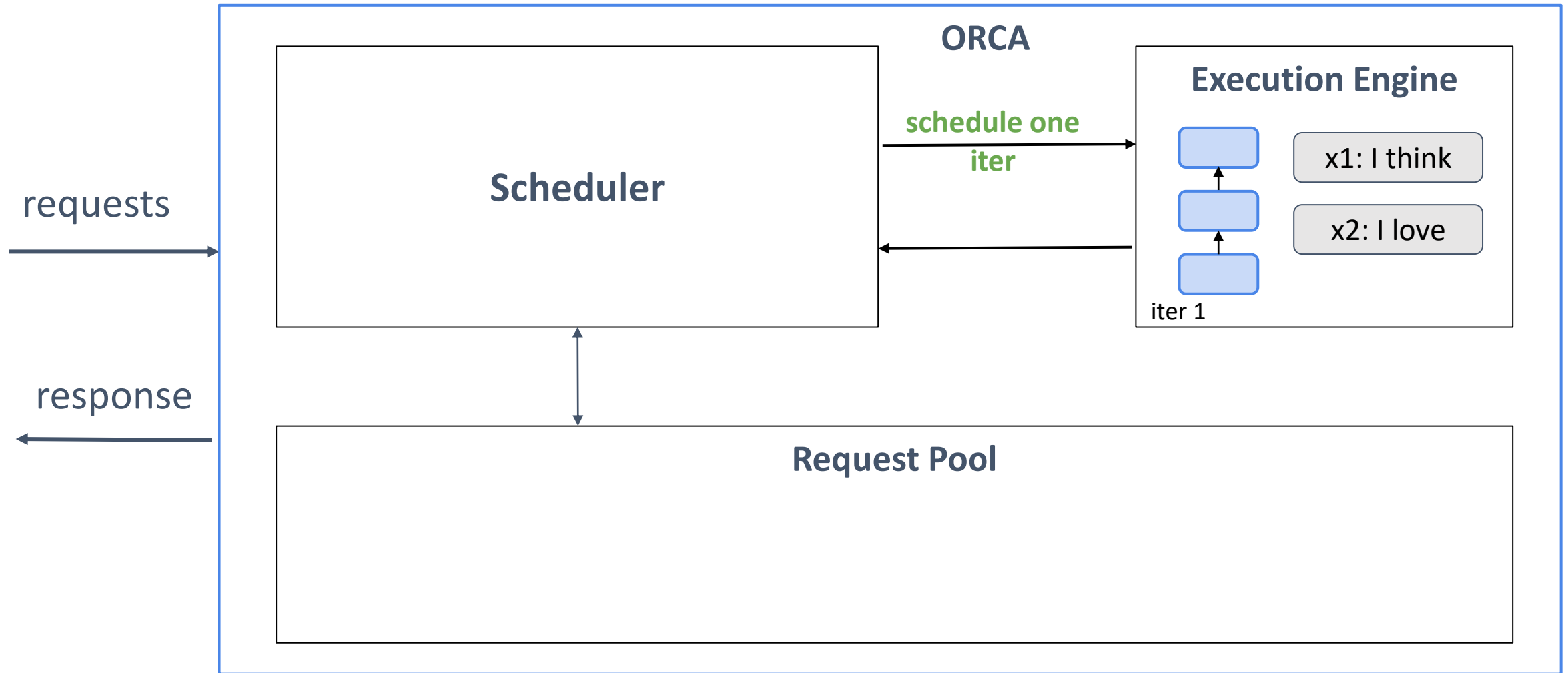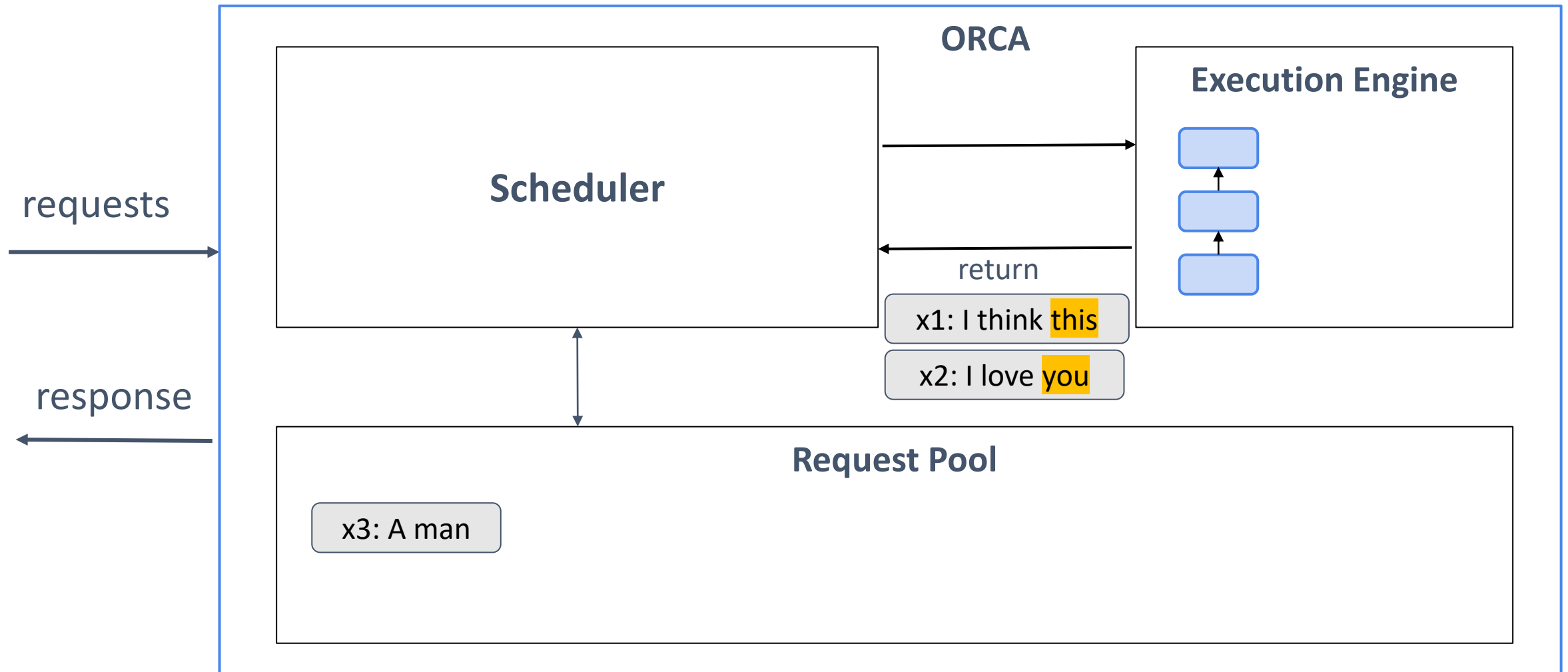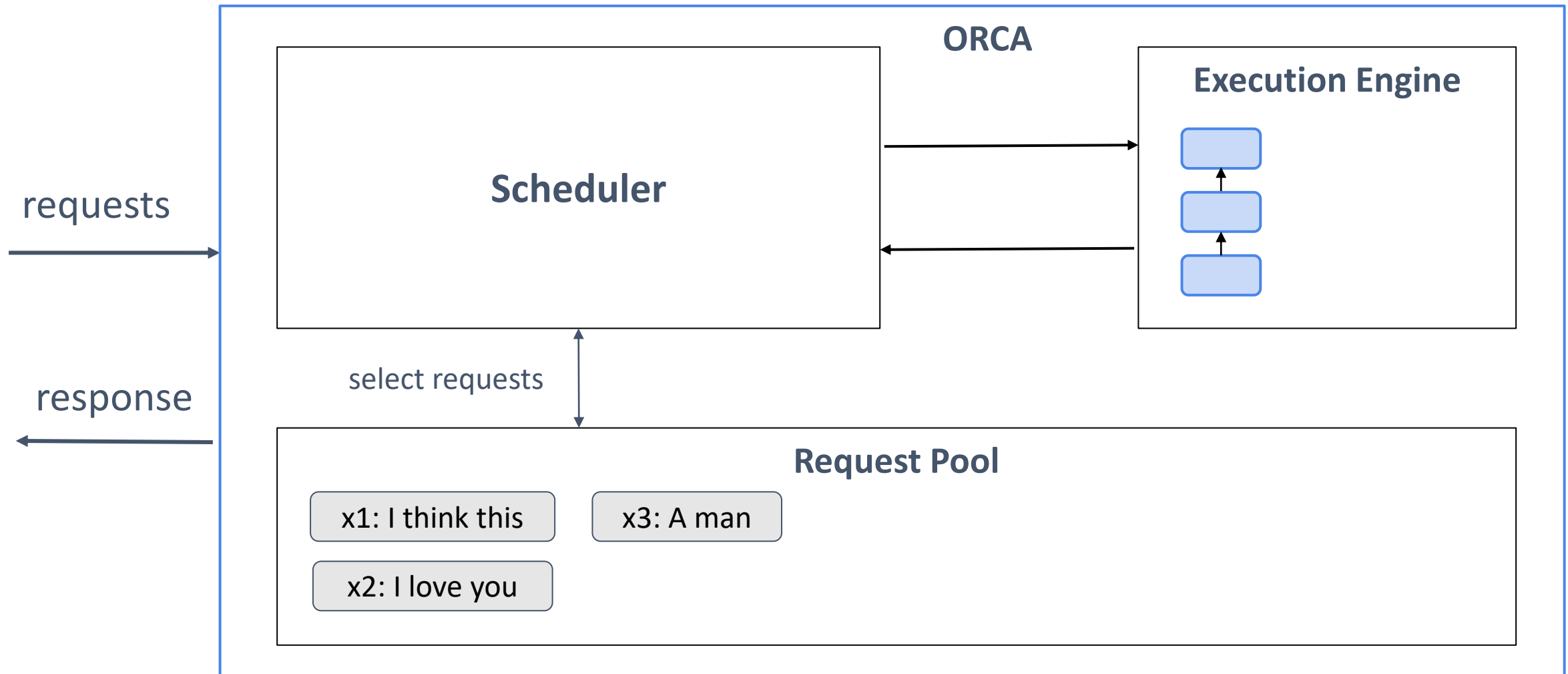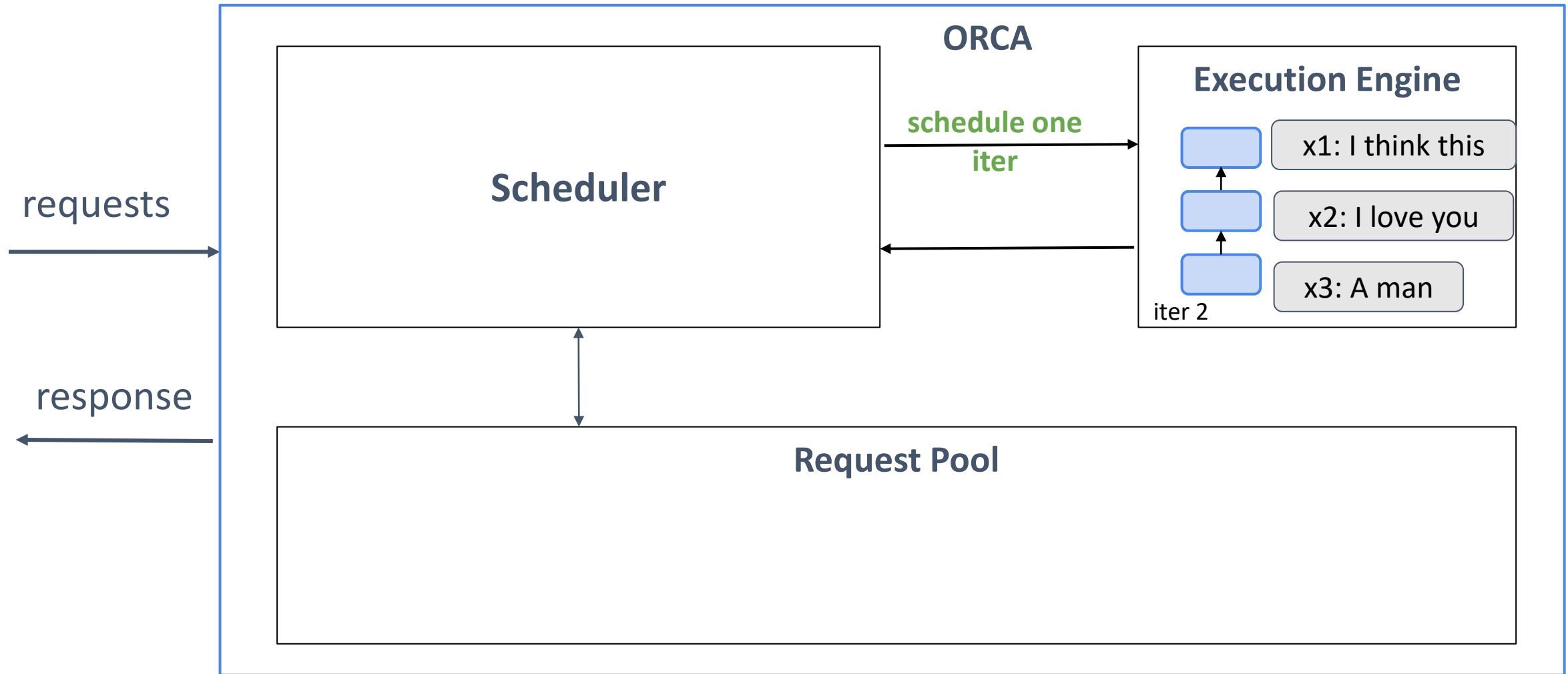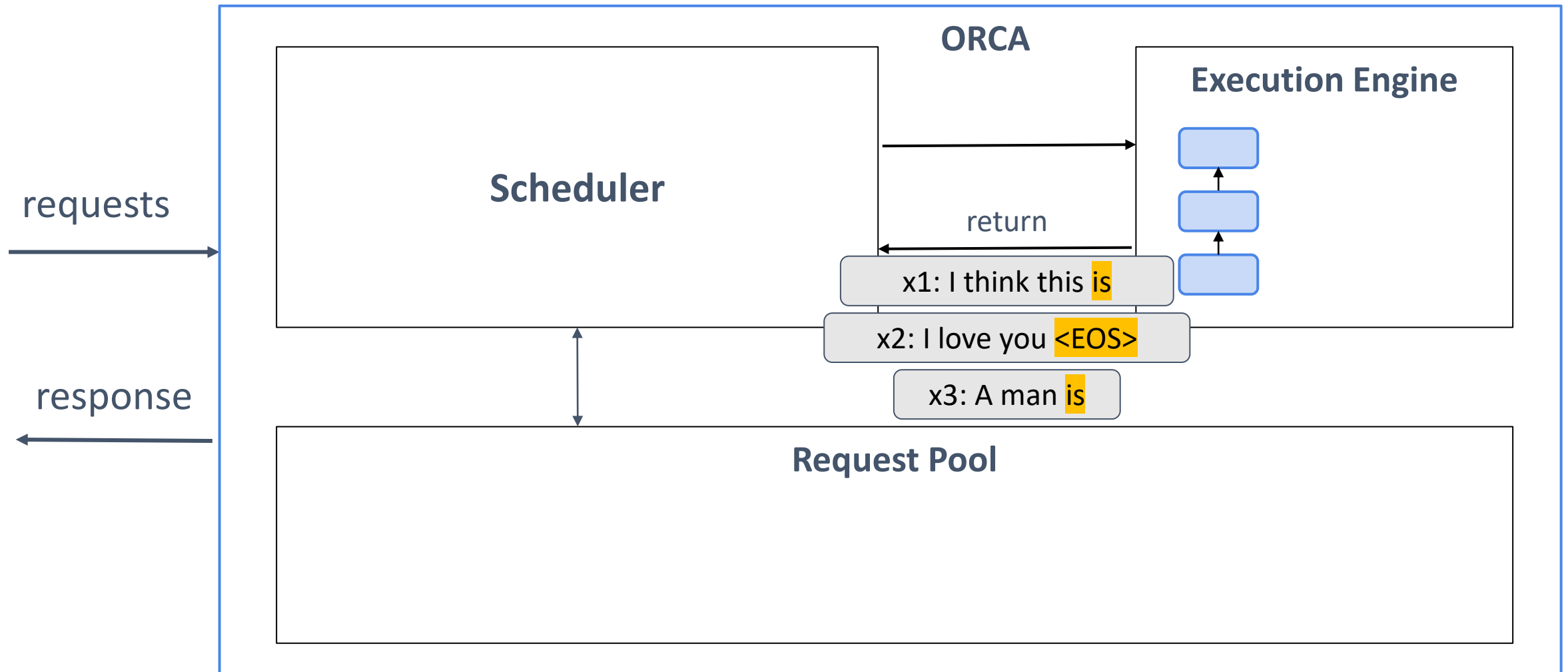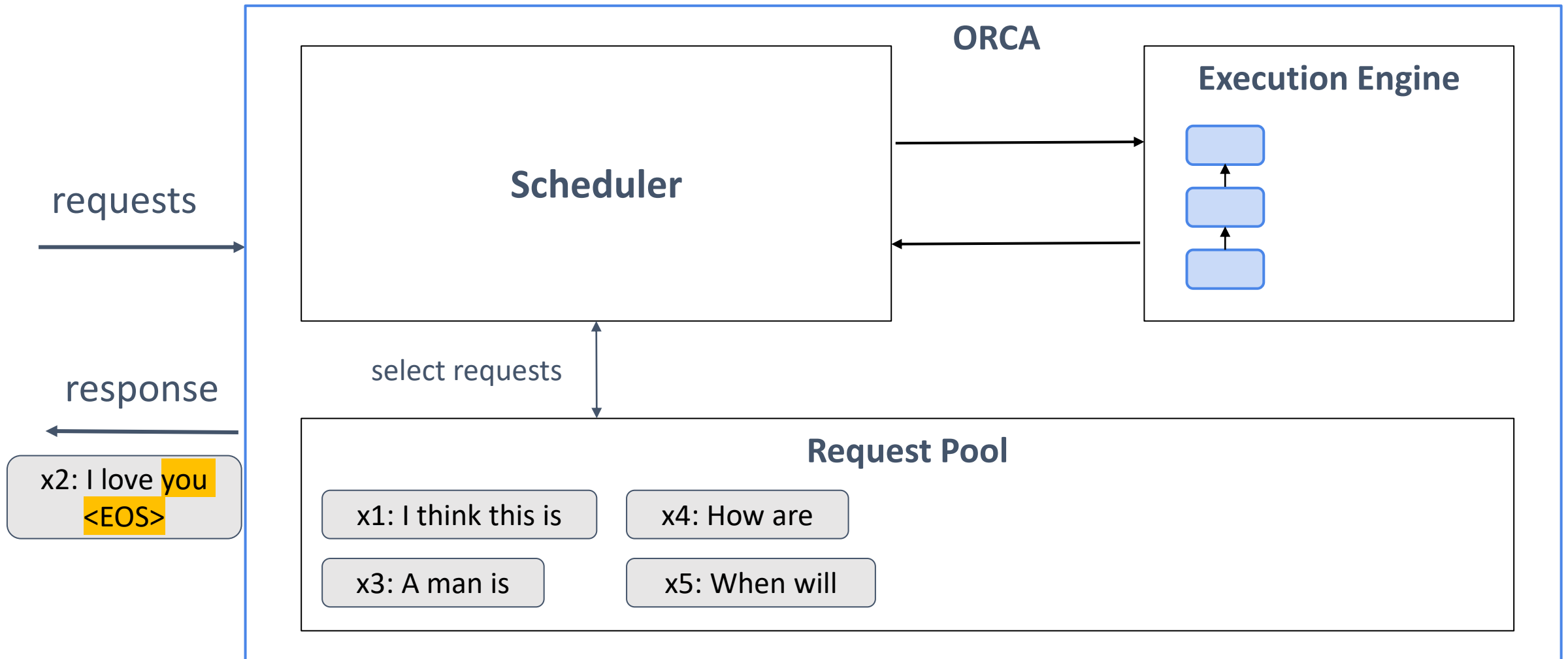
# Solution 1: Iteration Level Scheduling

# Solution 1: Iteration Level Scheduling

# Solution 1: Iteration Level Scheduling

**ORCA**

**Scheduler**

schedule one iter

**Execution Engine**

requests

x1: I think this

x2: I love you

x3: A man

iter 2

response

**Request Pool**

**ORCA**

**Scheduler**

**Execution Engine**

requests

response

x2: I love you <EOS>

select requests

**Request Pool**

x1: I think this is

x4: How are

x3: A man is

x5: When will

ORCA

**Execution Engine**

**Scheduler**

schedule one iter

x1: I think this is

...n is

...re

requests

response

**Iteration Level Scheduling handles early finished requests and late joining requests**

Request Pool

x5: When will

Let's assume Batch Size B = 1

**Input Dimension:** [L x H] (L=sequence length, H=hidden dim.)

**Attention Operation:**
1.  **$QK^T$** : [LxH] x [HxL] → [L x L]

2.  **P = softmax($QK^T$)** : [L x L]

3.  **O = PV** : [LxL] x [LxH] → [L x H]

Let's assume Batch Size B = 1

**Input Dimension:** [L x H] (L=sequence length, H=hidden dim.)

**Attention Operation:**

1. $\mathbf{QK^T}$ : [LxH] x [HxL] → [L x L]

2. $\mathbf{P = softmax(QK^T)}$ : [L x L]

3. $\mathbf{O = PV}$ : [LxL] x [LxH] → [L x H]

With Batch Size B, $\mathbf{QK^T}$ will be **[B x L x L]**
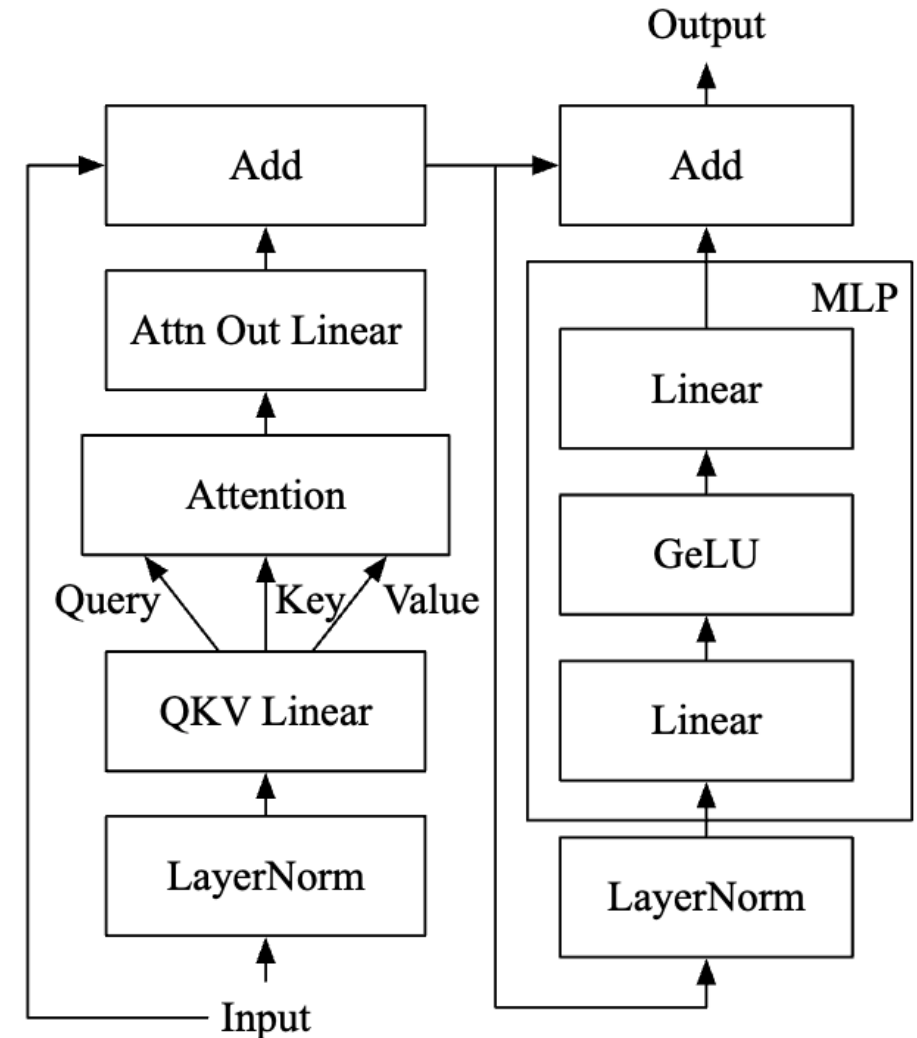
Let's assume Batch Size B = 1

**Input Dimension:** [L x H] (L=sequence length, H=hidden dim.)

**Attention Operation:**
1. $QK^T$ : [LxH] x [HxL] → [L x L]

2. $P = softmax(QK^T)$ : [L x L]

3. $O = PV$ : [LxL] x [LxH] → [L x H]

With Batch Size B, $QK^T$ will be [B x L x L]

> **With different sequence lengths, $QK^T$ cannot be easily computed**

Only **Attention operation** does not work with batching tensors with diff. $L_i$

Batch for other ops. (Layer Norm, GeLU, etc.)

Only **Attention operation** does not work with batching tensors with diff. $L_i$

Batch for other ops. (Layer Norm, GeLU, etc.)

Coalesce $[L_i, H]$ tensor to $[\Sigma L_i, H]$ for batching

x1: [1,H]
x2: [1,H]
x3: [2,H]
x4: [3,H]

→ **[7,H] tensor**



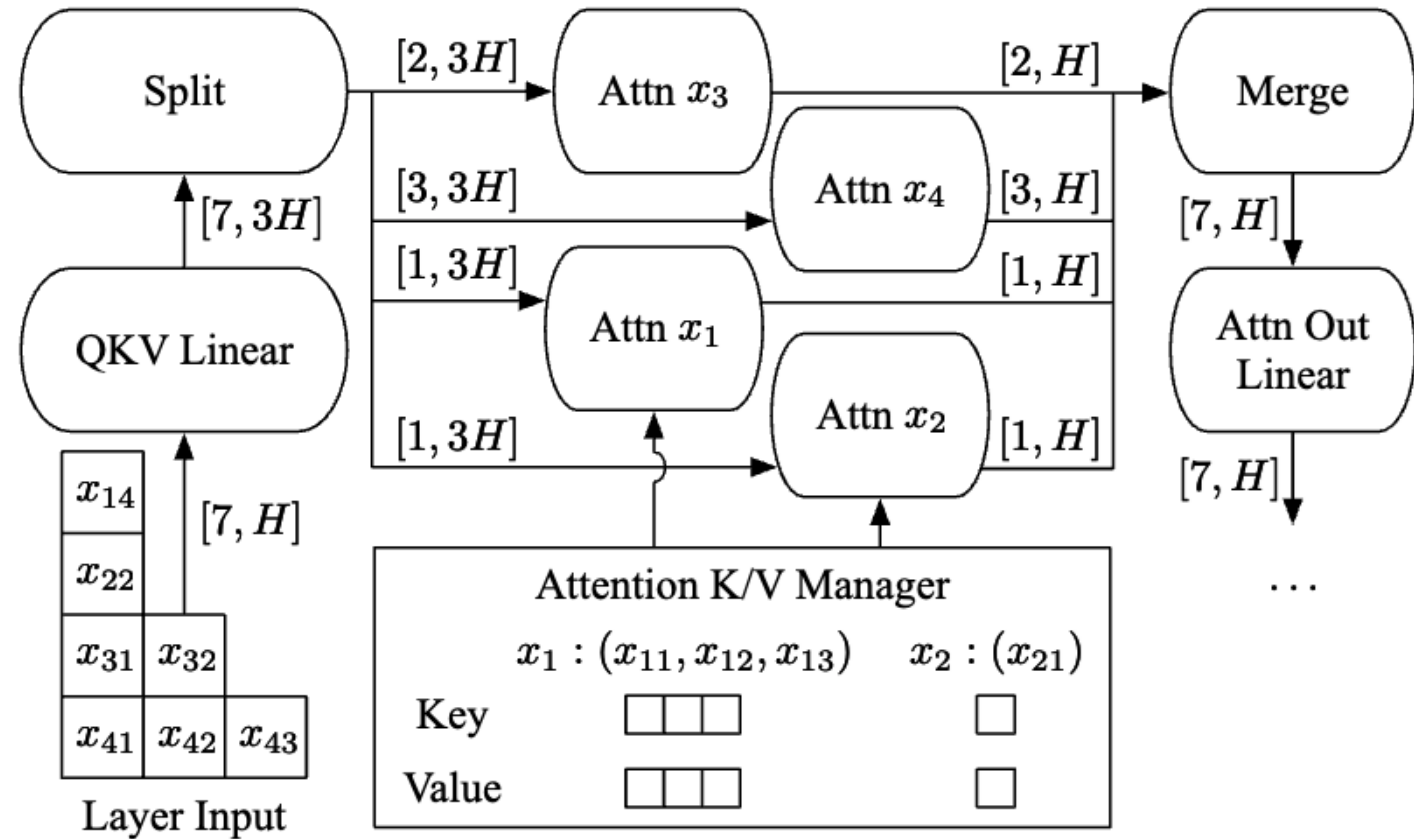Use flattened 2D tensor without batch dimension
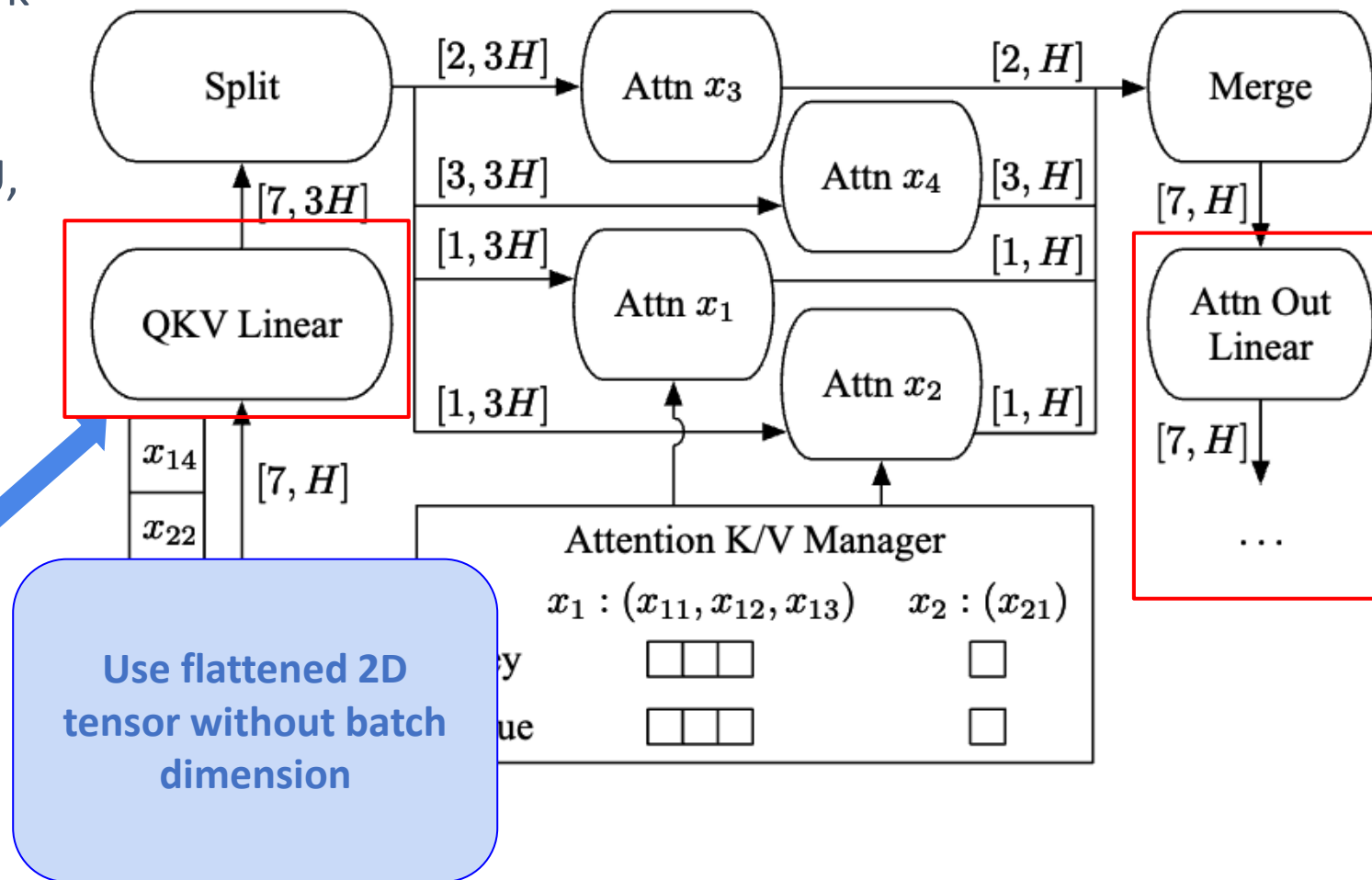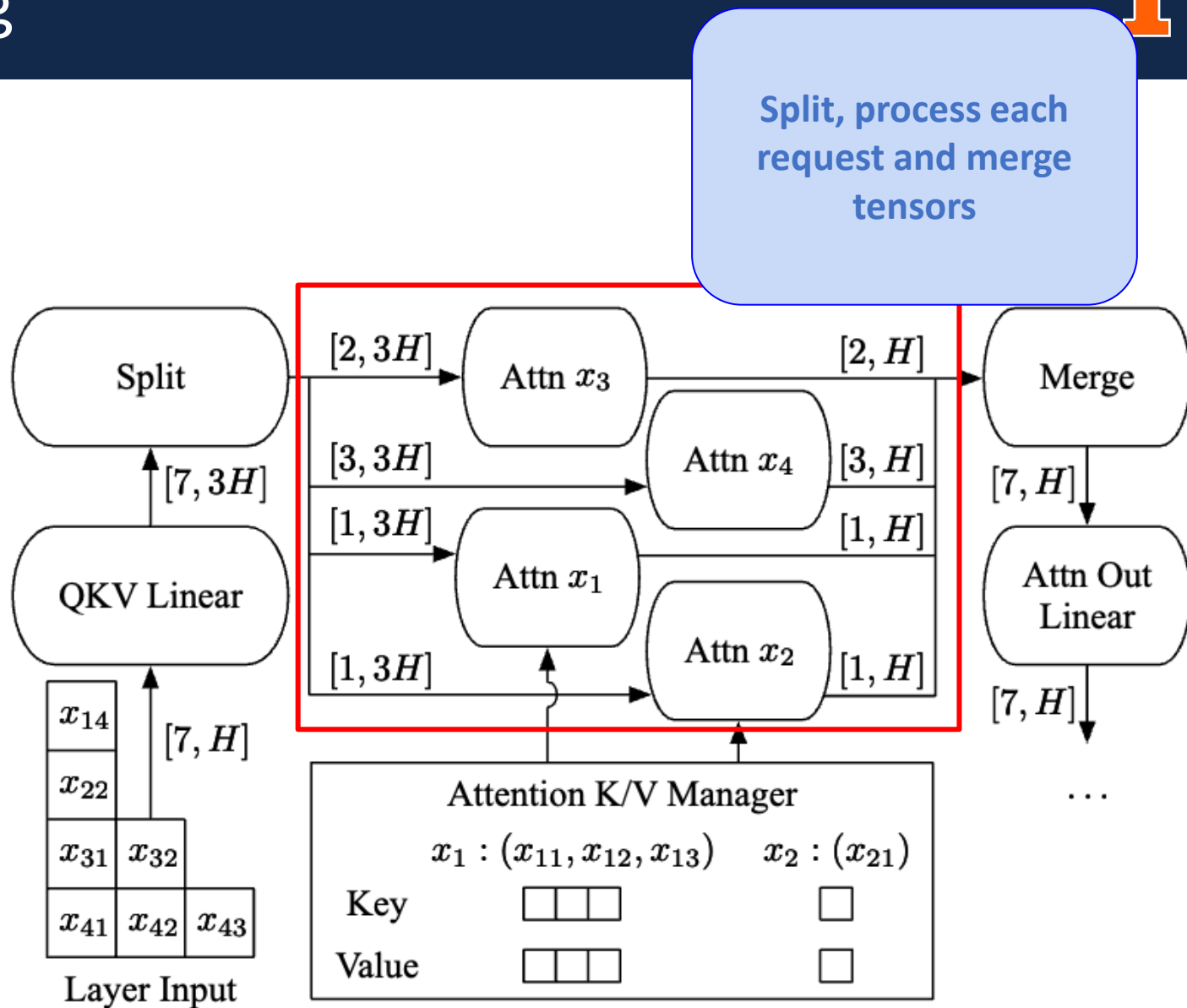
# Solution 2: Selective Batching

Only **Attention operation** does not work with batching tensors with diff. $L_i$

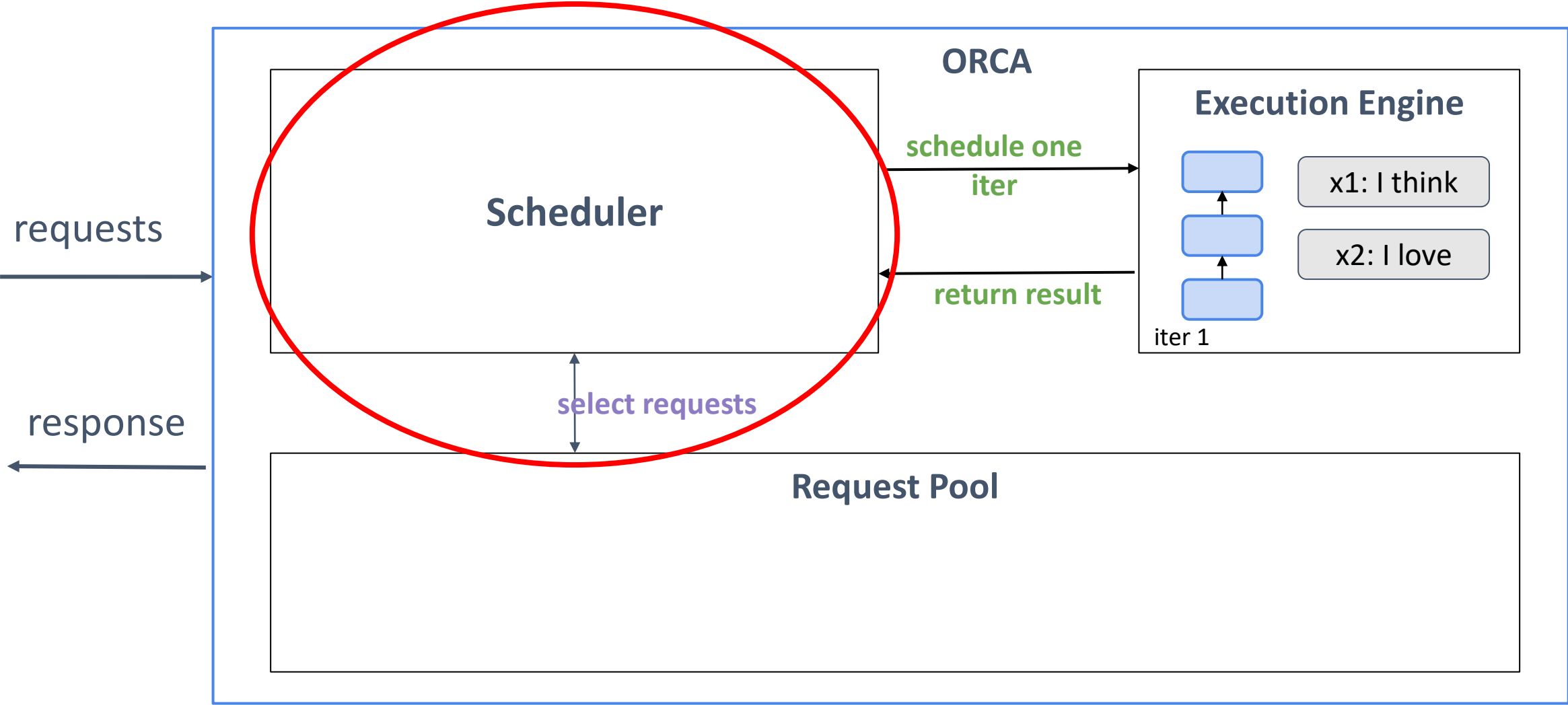Batch for other ops. (Layer Norm, GeLU, etc.)

Coalesce $[L_i, H]$ tensor to $[\Sigma L_i, H]$ for batching

x1: [1,H]
x2: [1,H]  ➡️  **[7,H] tensor**
x3: [2,H]
x4: [3,H]

# LLM Inference Scheduler

- Enforces iteration-level first-come-first-served (FCFS) property

- Maximum batch size → Throughput vs. Latency control knob

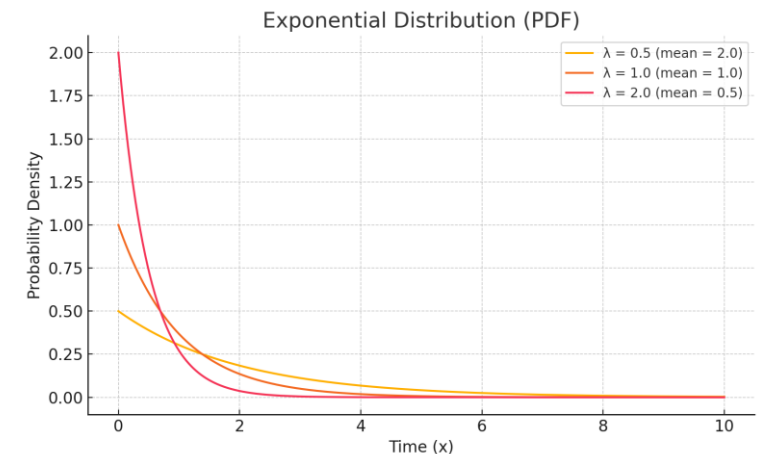- Reserves max_tokens memory slots per request

- ...

- Question: When would continuous batching provide more benefits than sequence batching?

- Hypothesis
  - Continuous batching performs better the more variance there is in sequence lengths

- Frameworks
- Setup – hardware/model
- Setup – data
- Results

- Static batching
  - HuggingFace
  - NVIDIA FasterTransformer

- Continuous batching
  - HuggingFace text-generation-inference (TGI)
  - Ray Serve
  - vLLM

https://www.anyscale.com/blog/continuous-batching-llm-inference

- NVIDIA A100-40GB

- Meta's OPT-13B
  - dtype = float16 → 26GB for parameters

- No tensor parallelism

- Hypothesis
  - Continuous batching performs better the more variance there is in sequence lengths

- How to test?
  - Generate 1000 prompts each with 512 input tokens
  - Generate predetermined output length for each prompt, following an exponential distribution
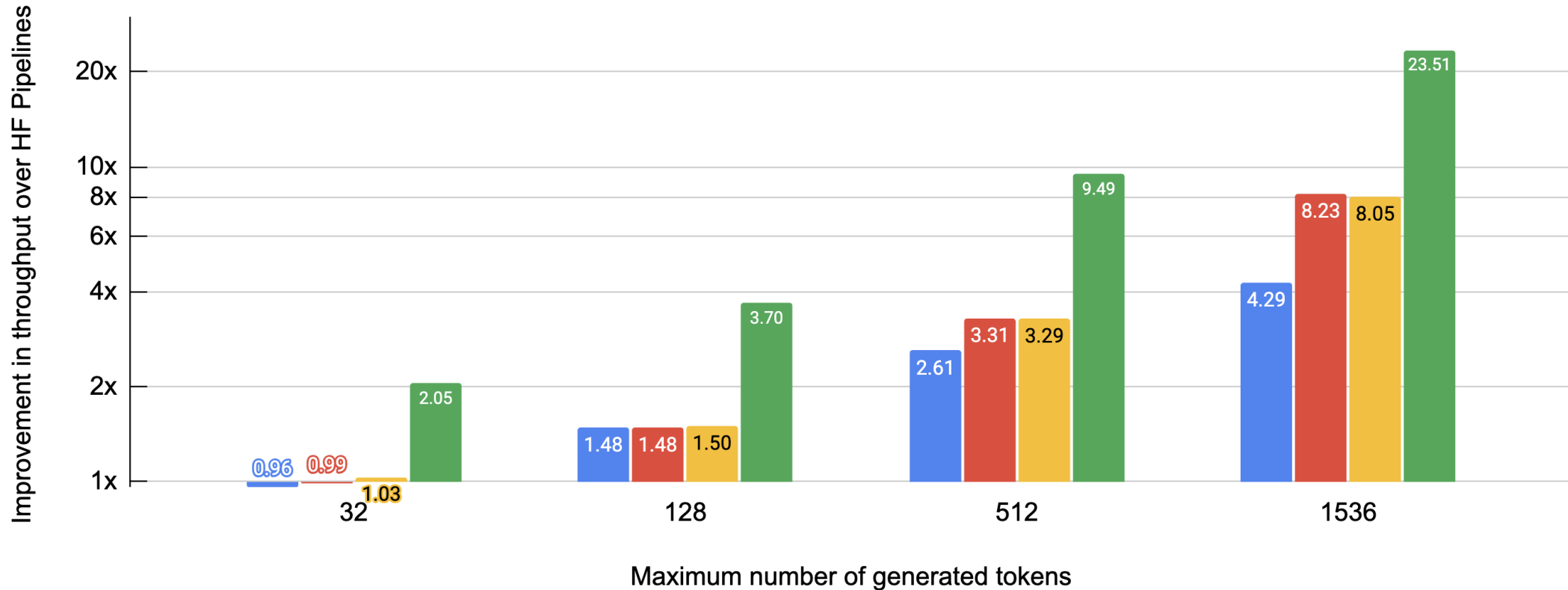  - Configure model to ignore EOS token



Exponential Distribution (PDF)

https://www.anyscale.com/blog/continuous-batching-llm-inference

# Throughput Improvement from Continuous Batching



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

https://www.anyscale.com/blog/continuous-batching-llm-inference

- vLLM uses PagedAttention – extra batch size space



Contiguous Memory

Non-Contiguous Memory

Request latency CDF, QPS=4

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

https://www.anyscale.com/blog/continuous-batching-llm-inference
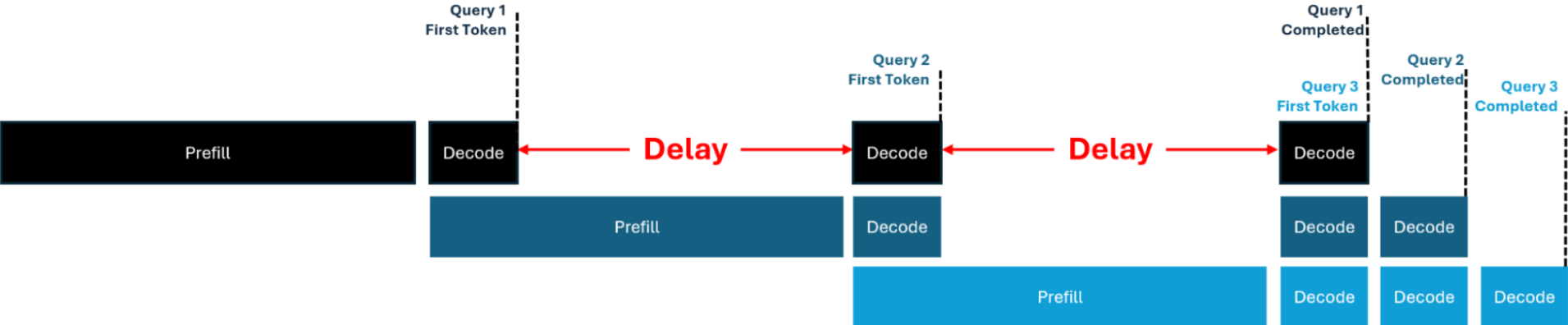
# Summary: Continuous Batching

- Continuous batching handles early-finished and late-arrived requests more efficiently

- Fills GPU capacity after each token generation

- As variance in sequence length increases, continuous batching increases GPU utilization
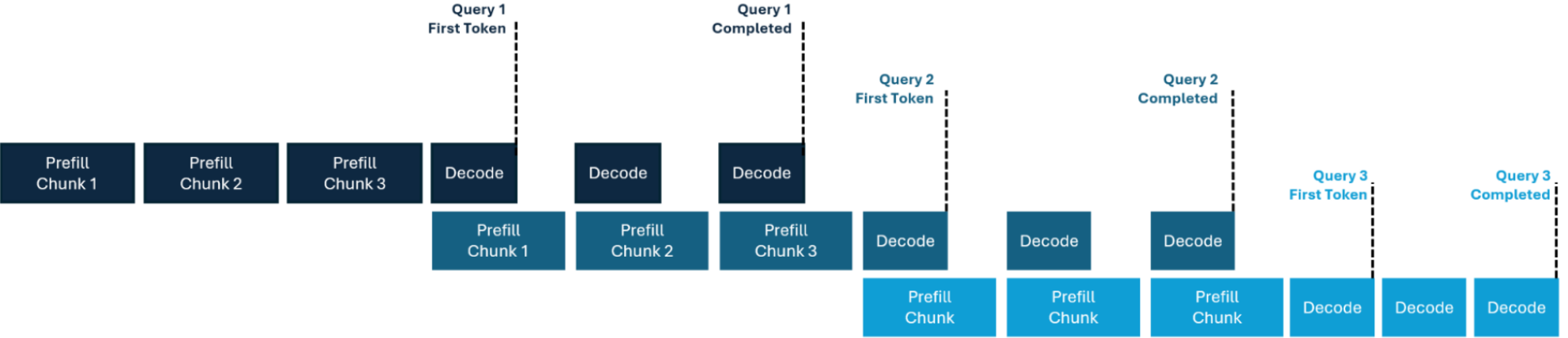
# Questions?

# Sequence Batching (Static Batching)

- Batching multiple sequences on GPU, aka "static batching"
- Problem: GPU utilization drops as sequences complete



Legend:
- Yellow: prompt token
- Blue: generated token
- Red: end-of-sequence token

Top: static batching
Bottom: continuous batching

Legend:
- Yellow: prompt token
- Blue: generated token
- Red: end-of-sequence token

# Throughput Experiments: Results

| Throughput (token/s) vs. variance in generated sequence lengths | Generation limit (higher limit implies higher variance in output sequence lengths) | | | |
|---|---|---|---|---|
| | max 32 tokens | max 128 tokens | max 512 tokens | max 1536 tokens |
| **Static batching** (HF Pipelines) | 2988 | 972 | 214 | 81 |
| **Static batching** (FasterTransformer) | 2869 | 1441 | 558 | 346 |
| **Continuous batching** (Ray Serve) | 3090 | 1460 | 703 | 650 |
| **Continuous batching** (text-generation-inference) | 2948 | 1442 | 707 | 665 |
| **Continuous batching** (vLLM) | 6121 | 3592 | 2029 | 1898 |

https://www.anyscale.com/blog/continuous-batching-llm-inference