



CS 498: Machine Learning System Spring 2025

Minjia Zhang

The Grainger College of Engineering

DL Inference

- LLM Inference Basic
- LLM Serving System
- Continuous Batching

Learning objectives

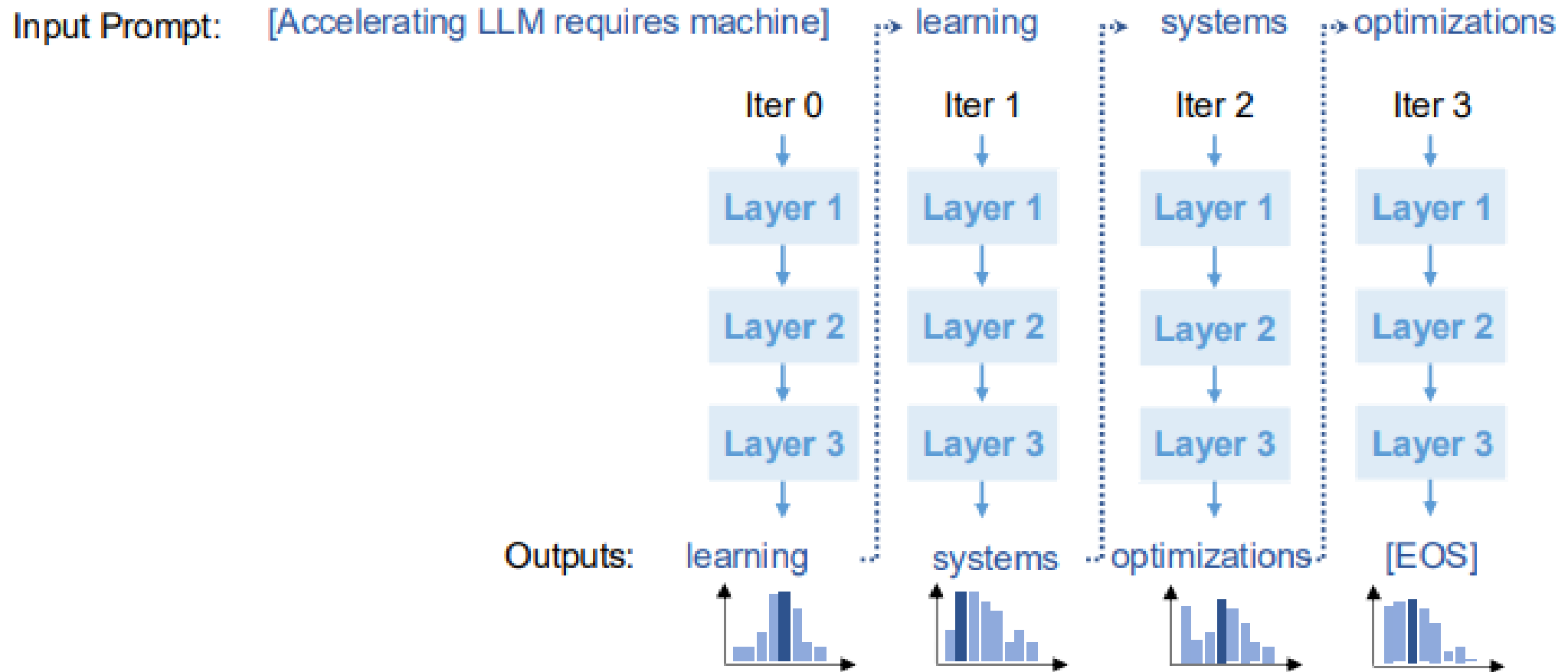
- Understand LLM inference basics
- Able to explain why sequence batching leads to inefficiencies
- Analyze how continuous batching and iteration-level scheduling improve GPU utilization

LLMs are Slow and Expensive to Serve

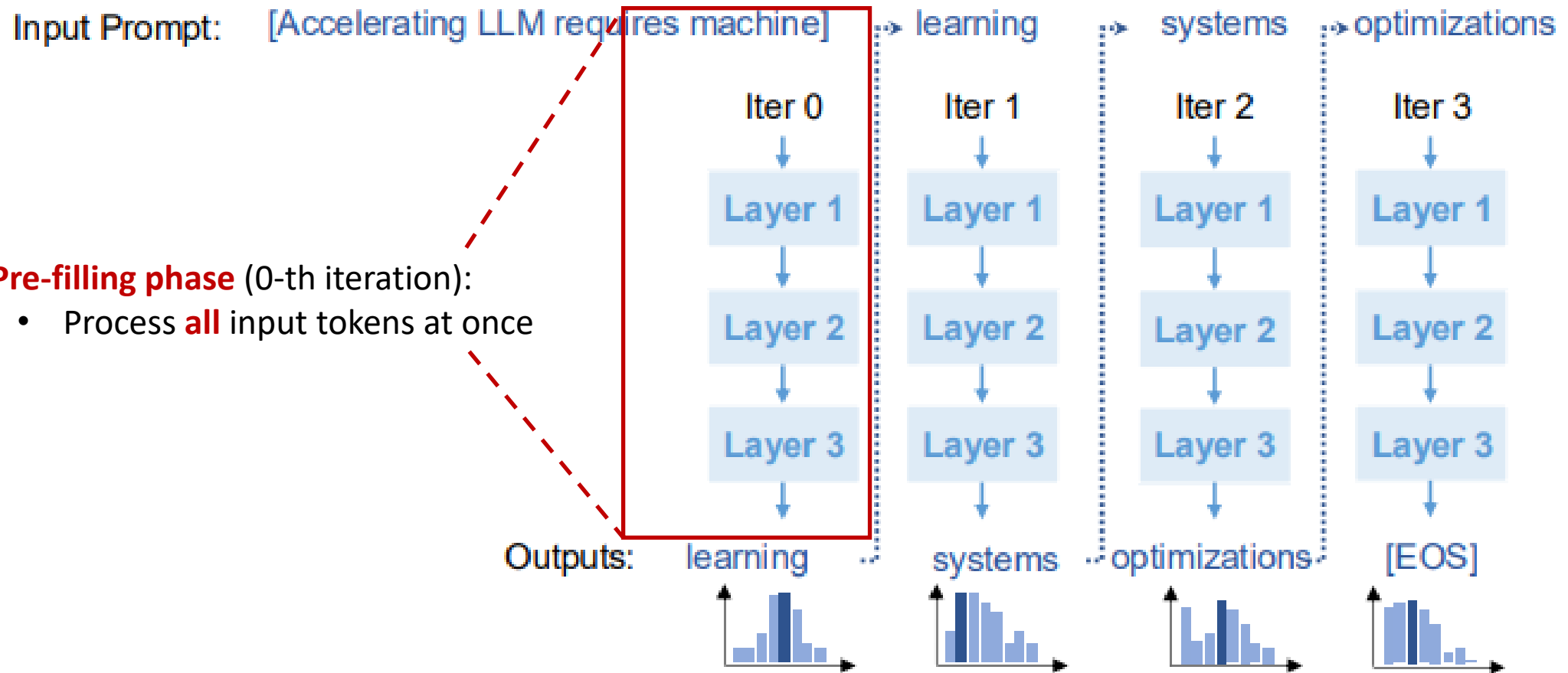


- **At least ten** A100-40GB GPUs to serve 175B GPT-3 in half-precision
- Generating 256 tokens takes **~20 seconds**
- Cannot process requests in parallel
 - Per-request key-value cache takes **3GB GPU memory**

Generative LLM Inference: Autoregressive Decoding

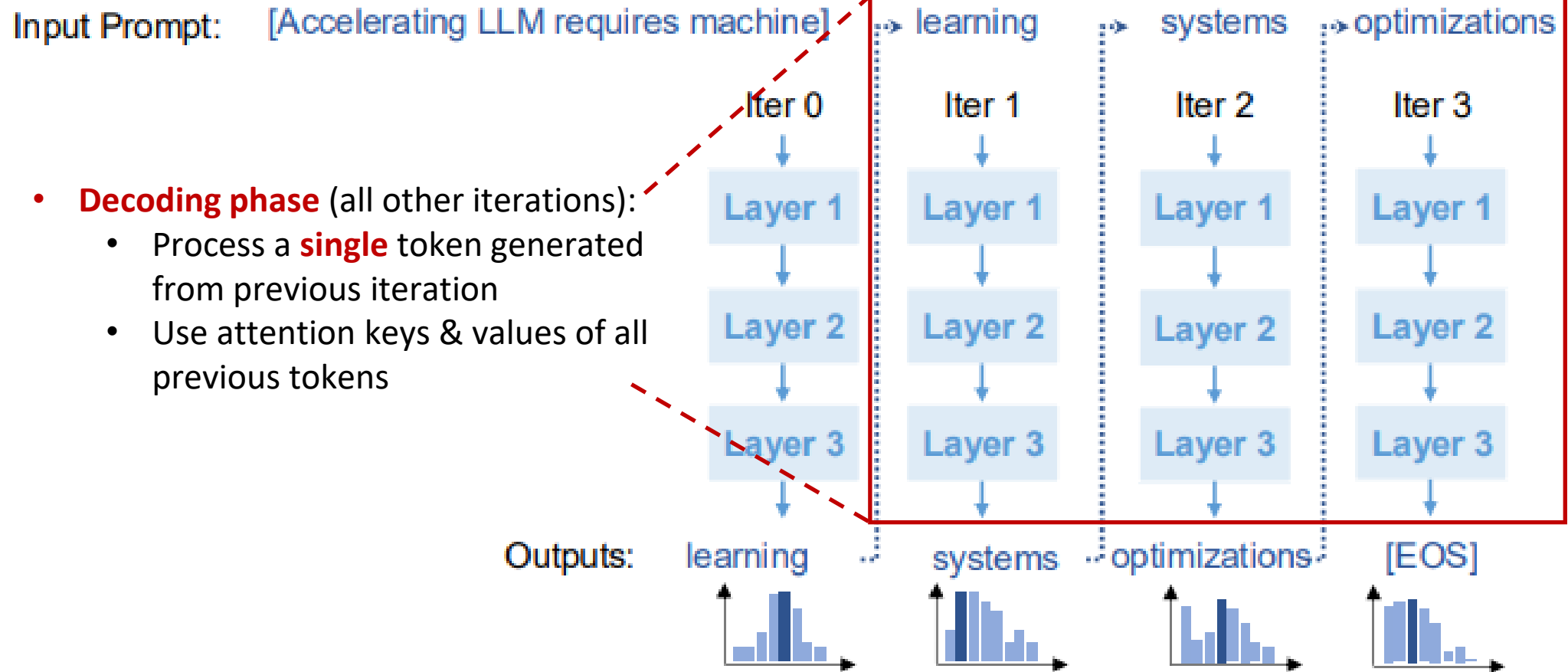


Generative LLM Inference: Autoregressive Decoding

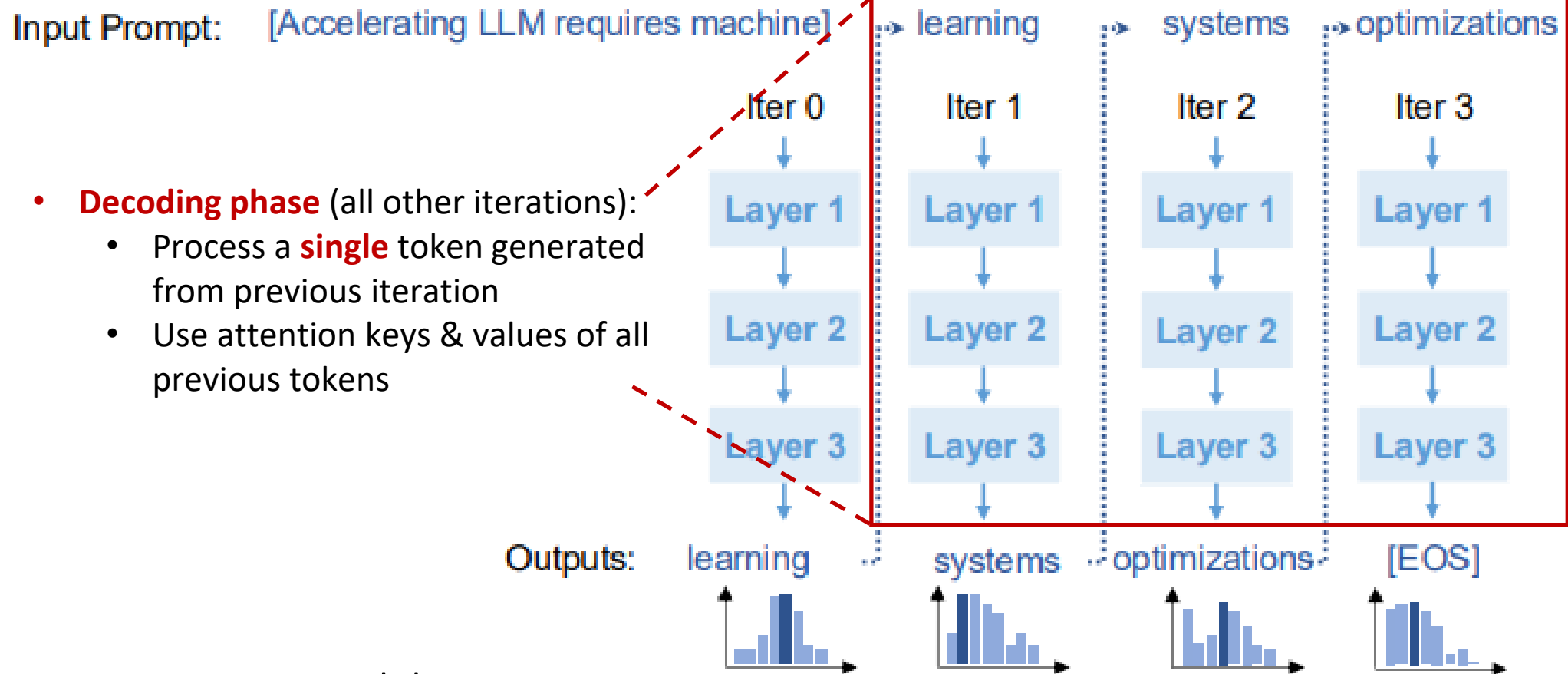


- **Pre-filling phase** (0-th iteration):
 - Process **all** input tokens at once

Generative LLM Inference: Autoregressive Decoding



Generative LLM Inference: Autoregressive Decoding

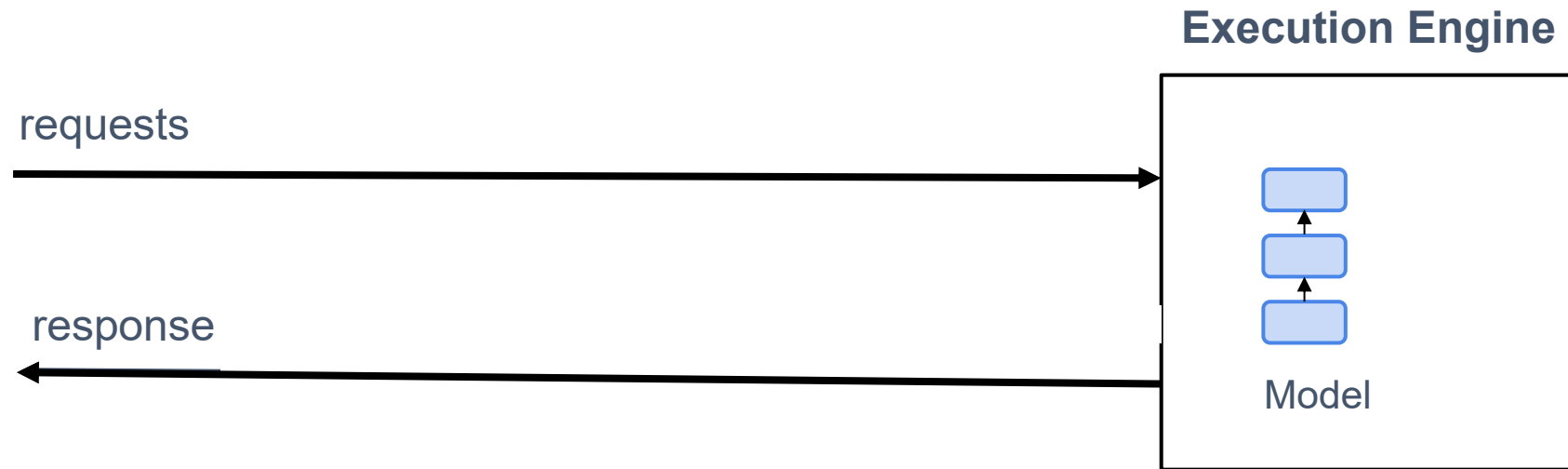


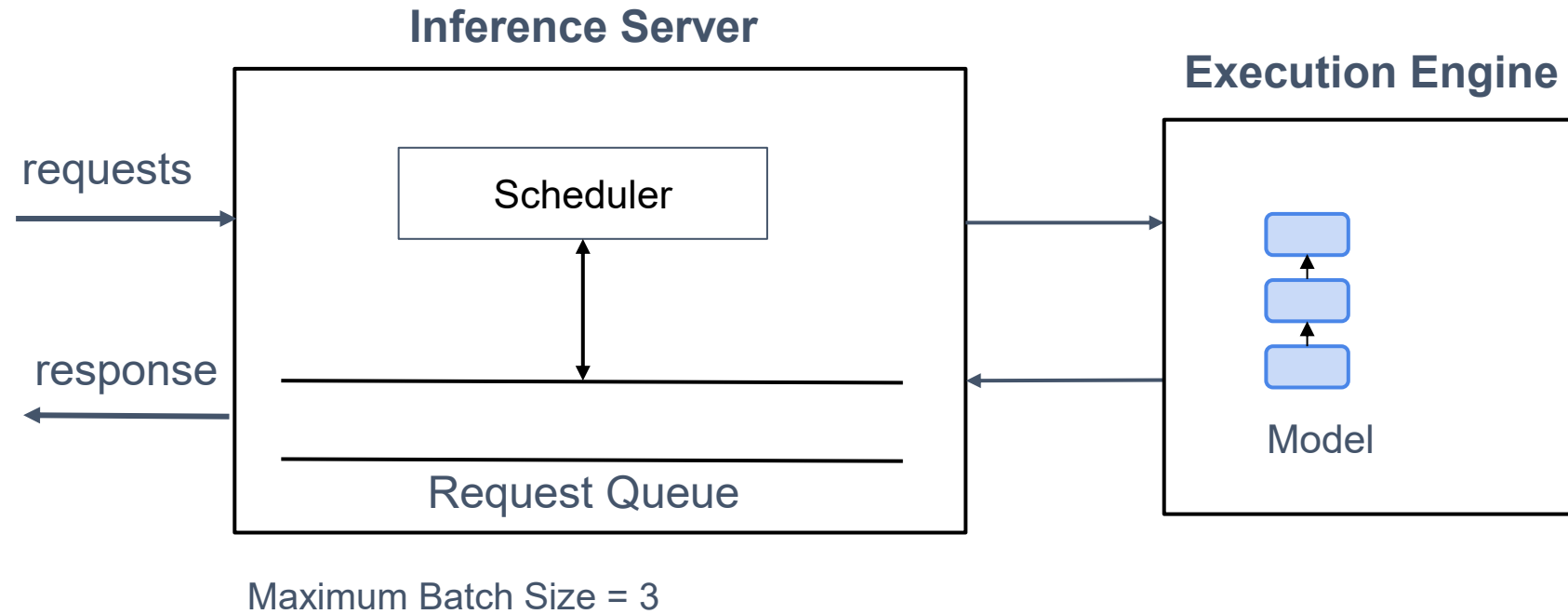
Repeat until the sequence

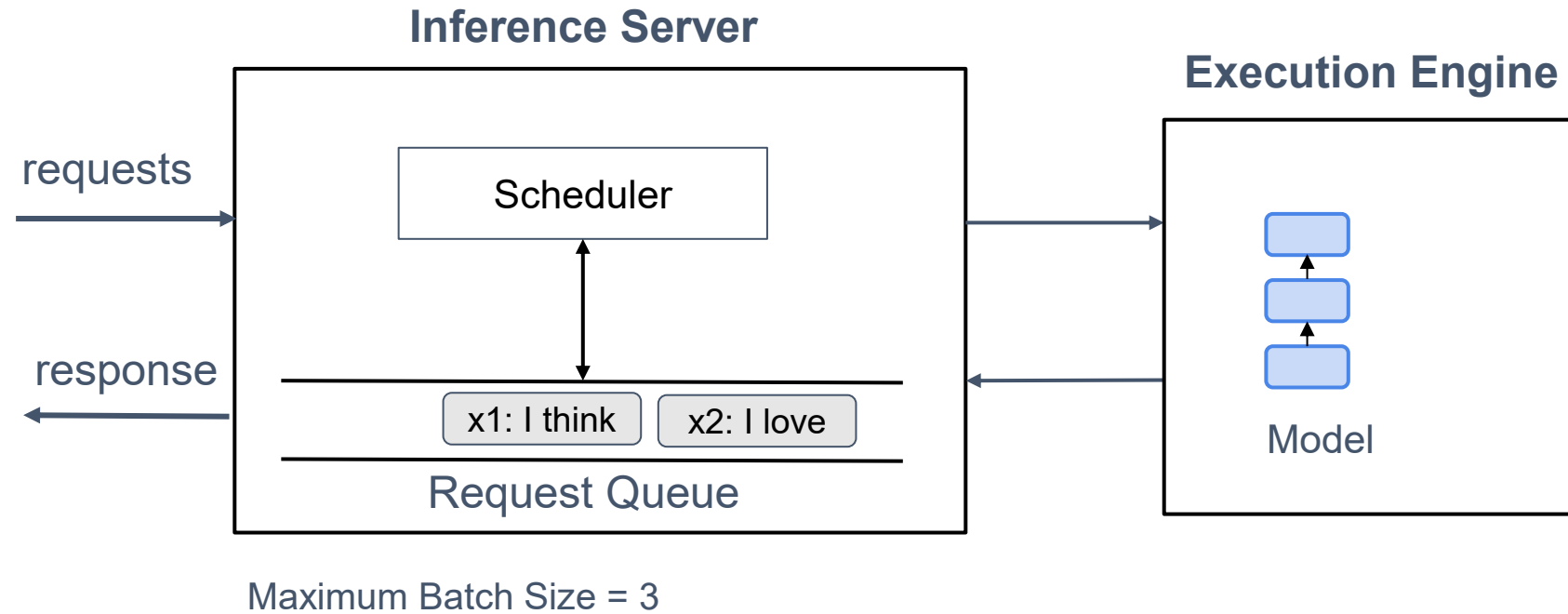
- Reaches the pre-defined maximum length (e.g., 2048 tokens)
- Generates the stop tokens (e.g., “<|end of sequence|>”)

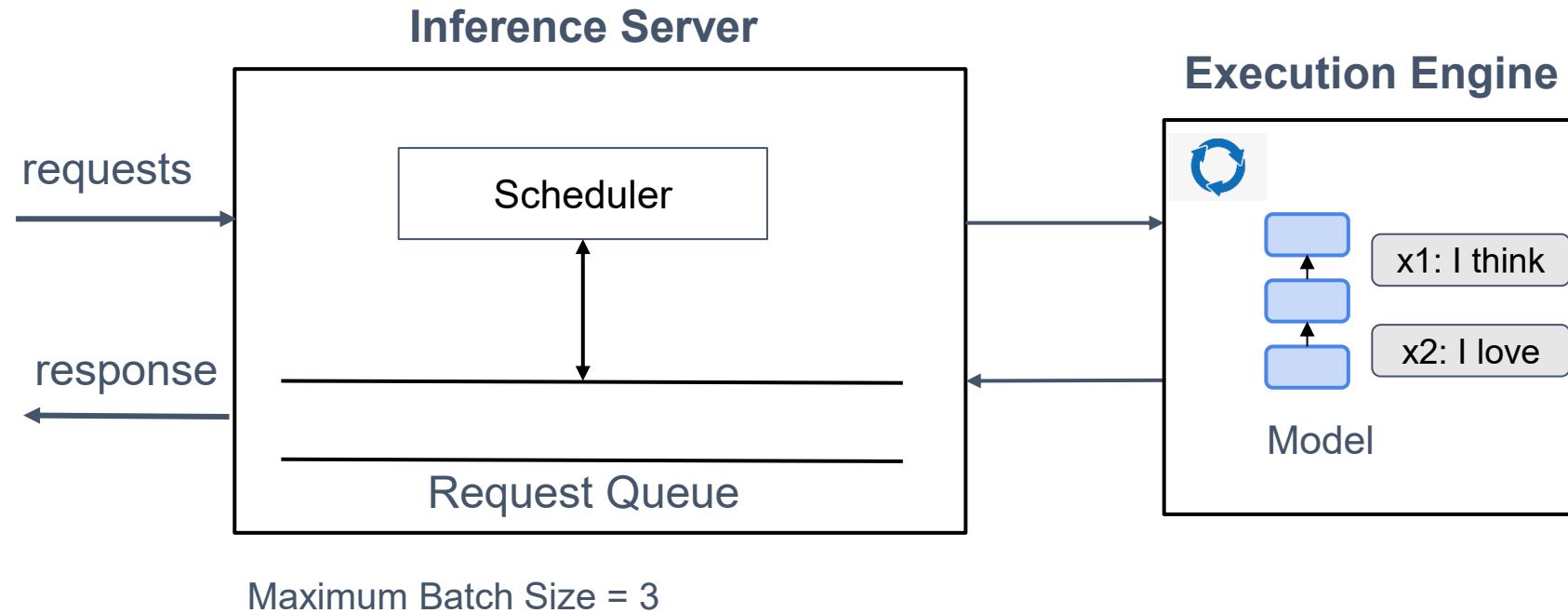
- **Time to First Token (TTFT):** Measures how quickly users begin to see the model output token after submitting a query.
 - Critical for real-time interactions
 - Driven by prompt processing time and the generation of the first token
- **Time per Output Token (TPOT):** Time taken to generate each output token.
 - Impacts user perception of speed (e.g., 100ms/token = 10 tokens/second)
- **E2E Latency = TTFT + (TPOT * the number of generated tokens)**
 - Total time to generate the complete response
- **Throughput:** Number of tokens generated per second across all requests by the inference server

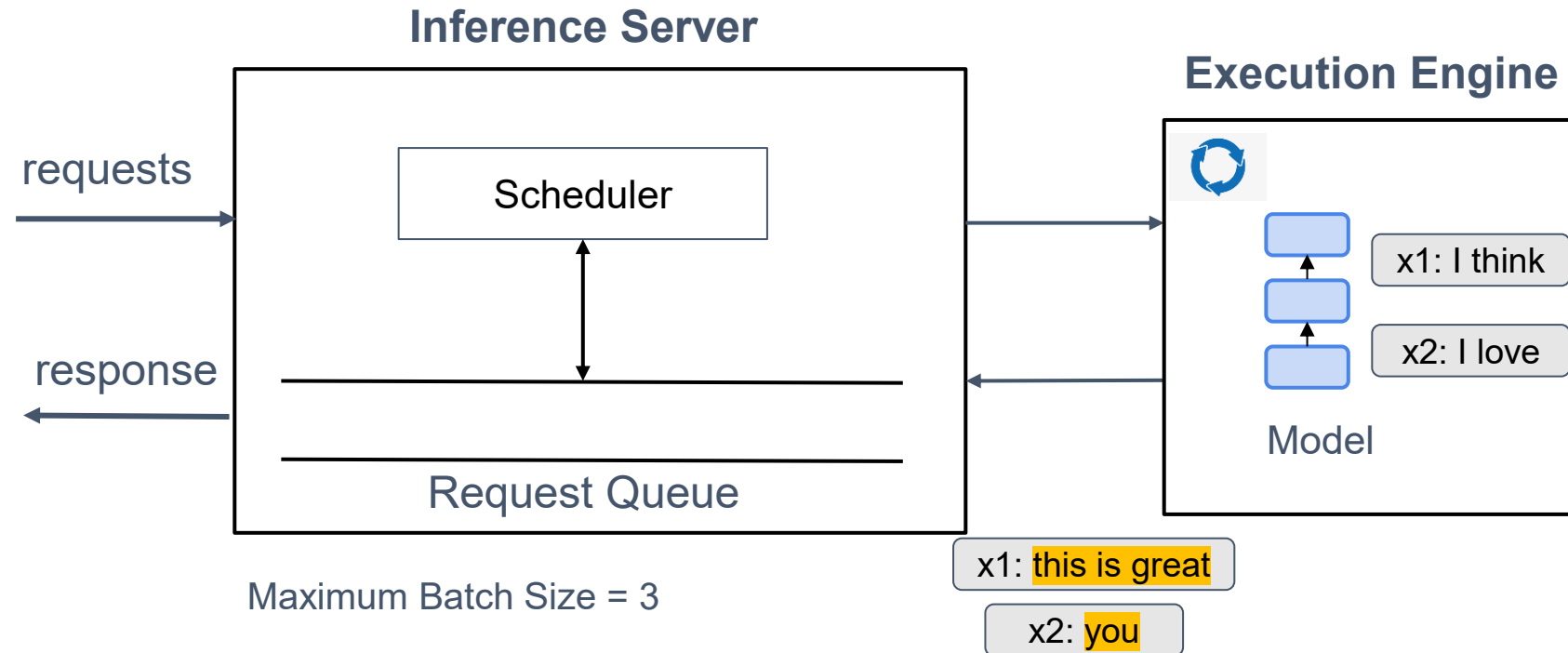
- **Goal:** Minimize TTFT, maximize throughput, and reduce TPOT
- **Throughput vs. TPOP Tradeoff:** Processing multiple queries concurrently increases throughput extends TPOT for each user.

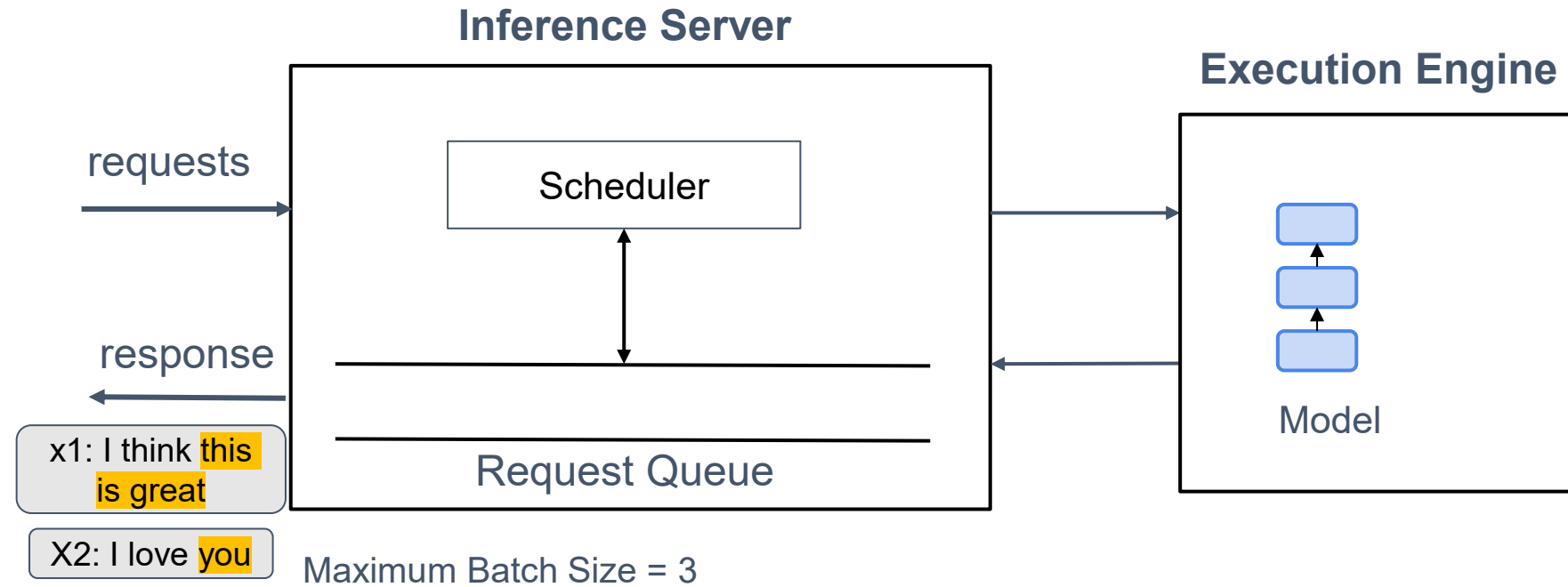




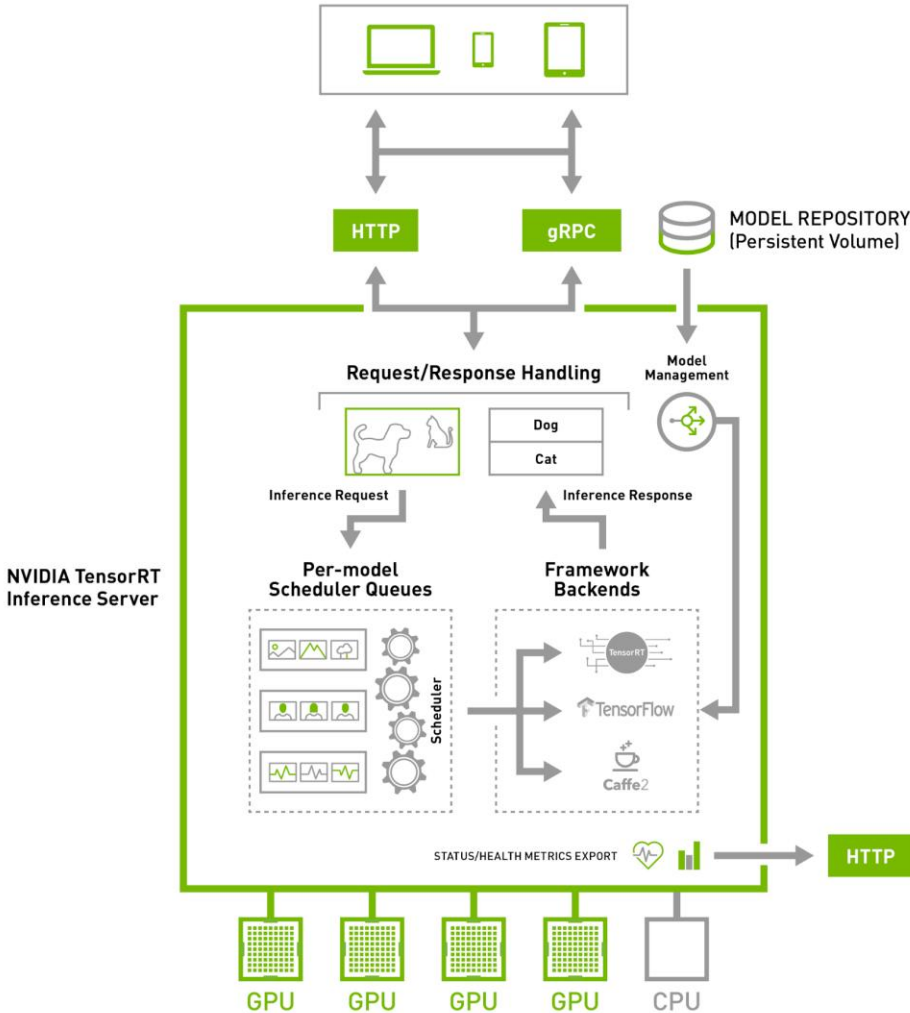








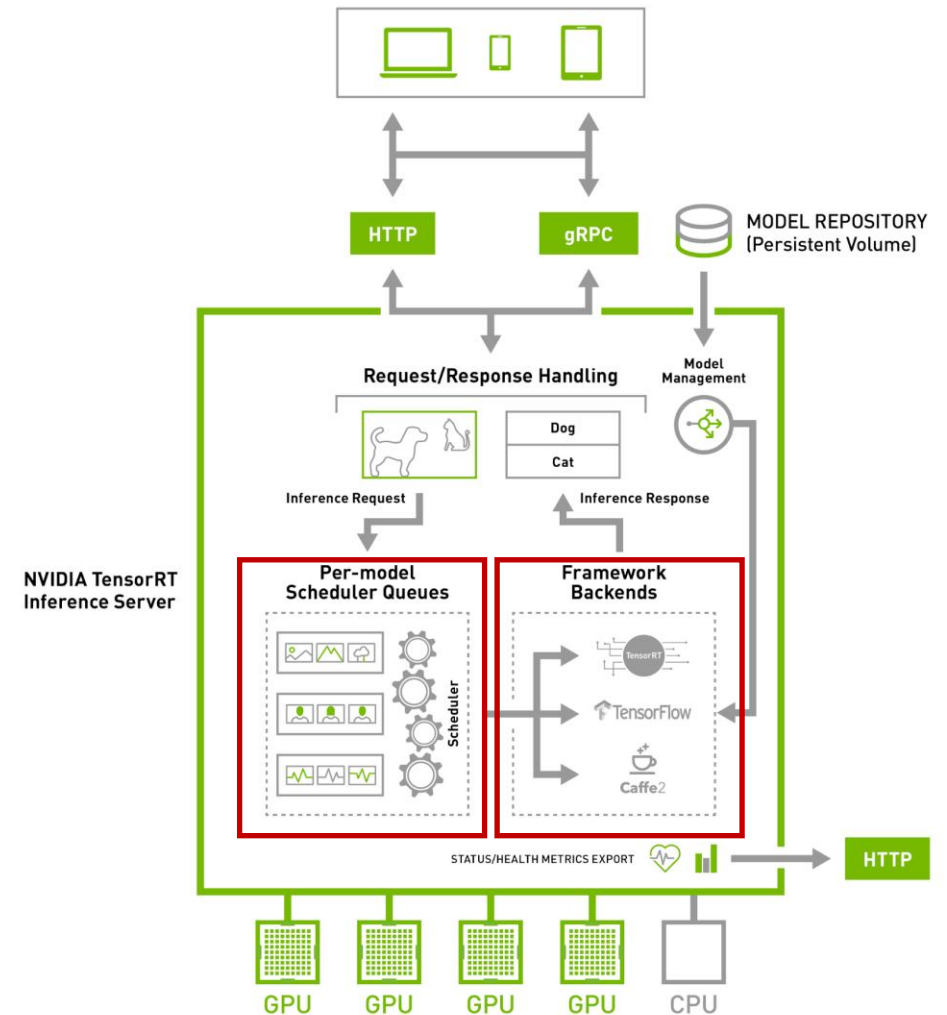
Example: TensorRT Inference Server



Example: TensorRT Inference Server



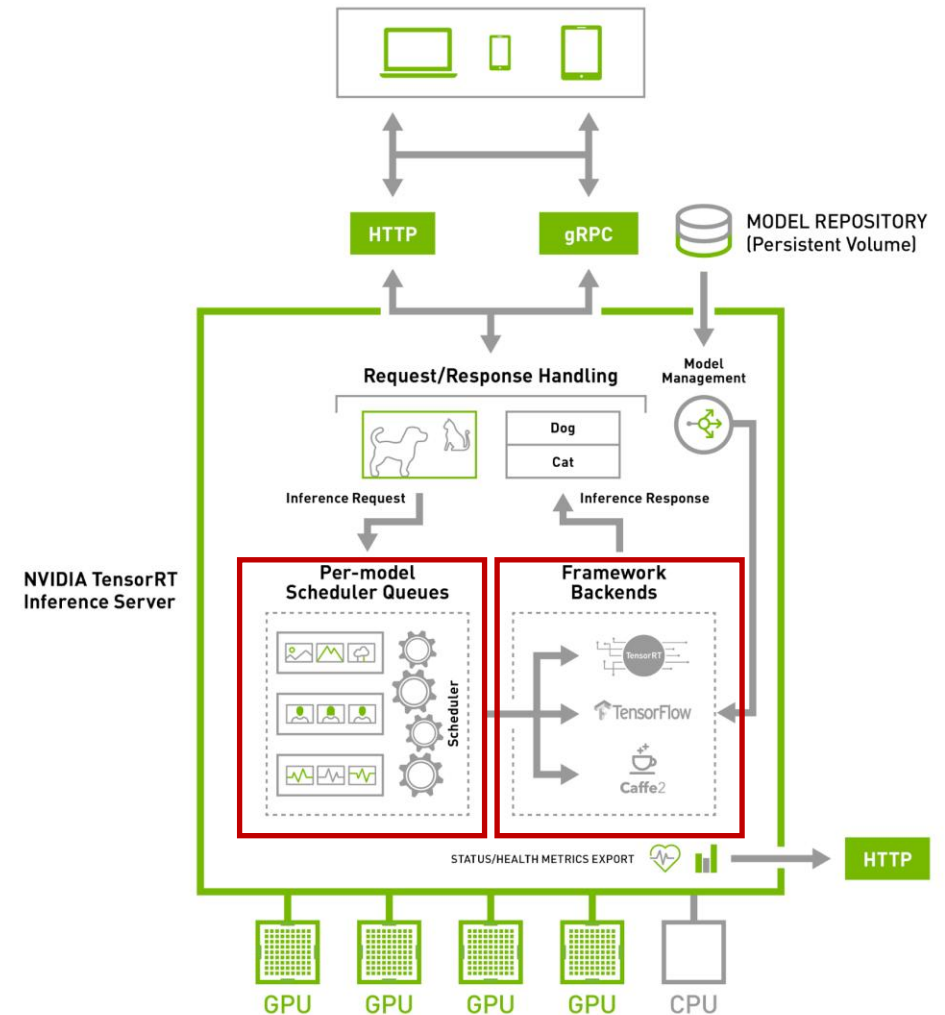
- Separates implementation of serving layer and execution layer



Example: TensorRT Inference Server



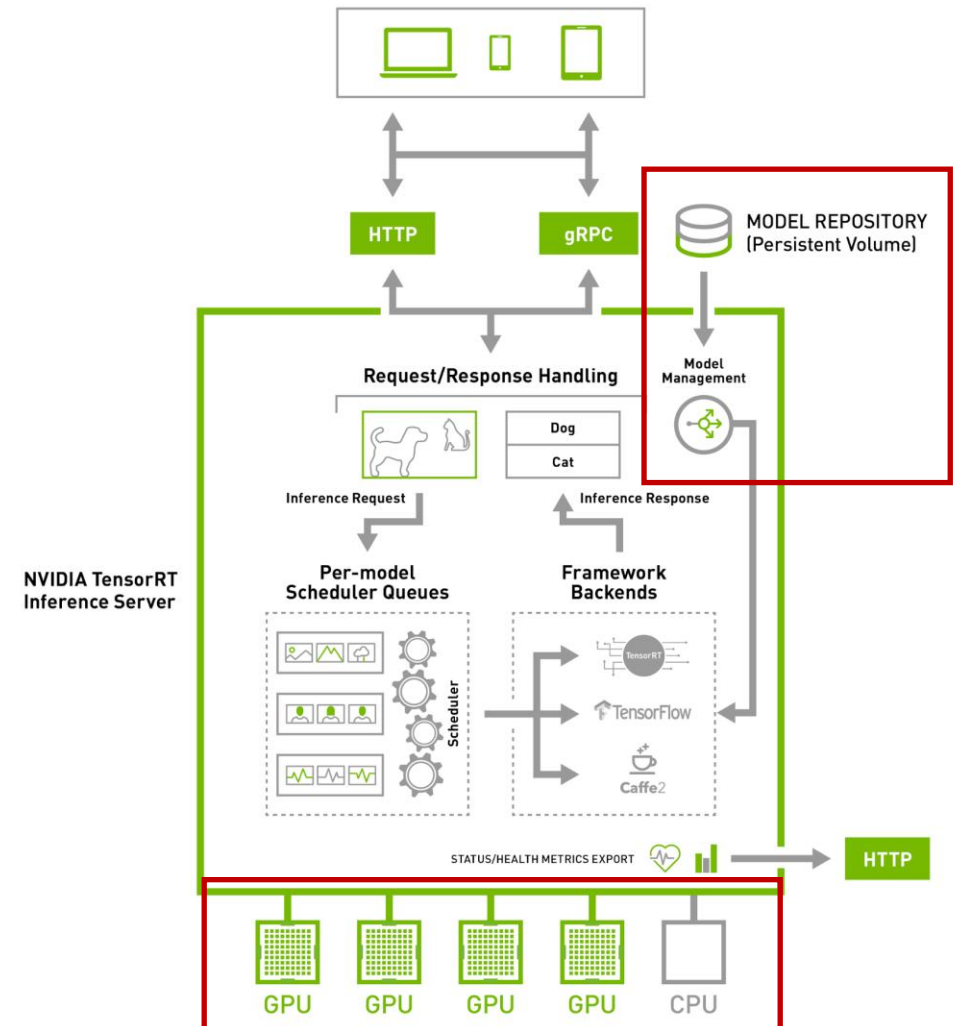
- Separates implementation of serving layer and execution layer
- Implements scheduling and batching algorithms
 - Sequence Batching
 - Continuous Batching



Example: TensorRT Inference Server



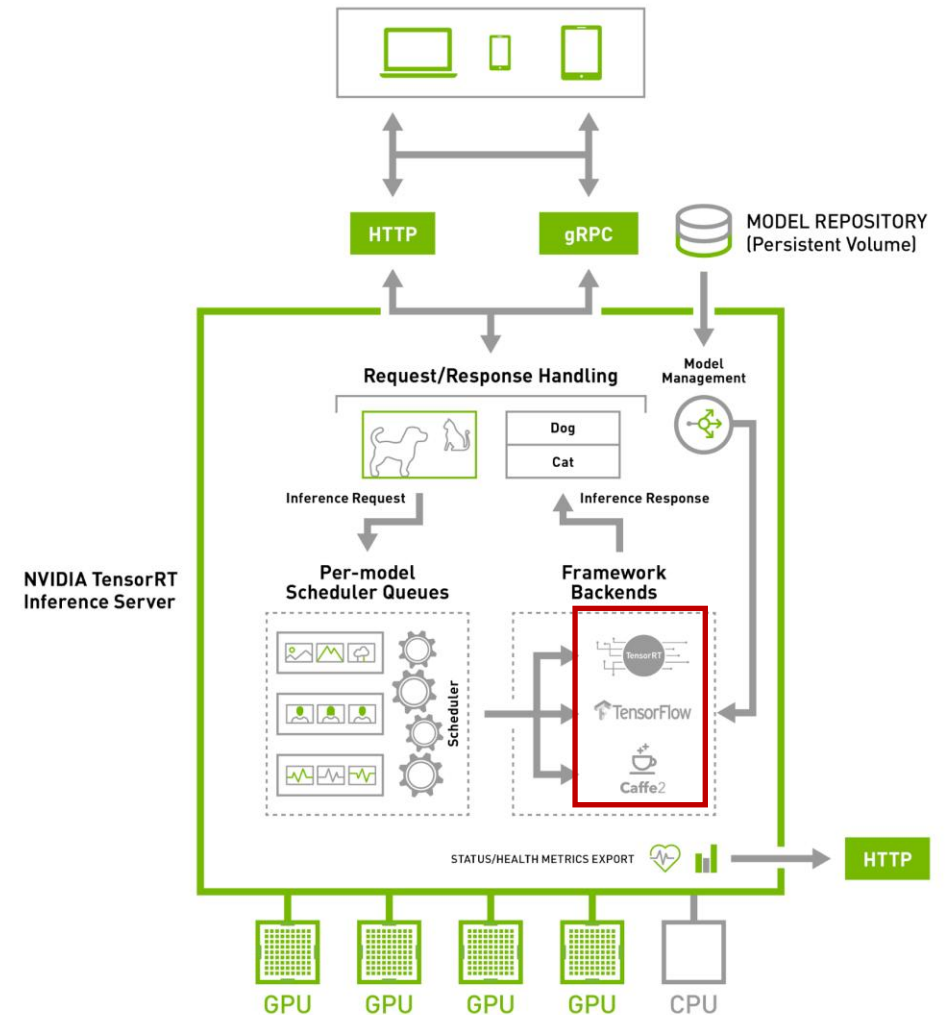
- Separates implementation of serving layer and execution layer
- Implements scheduling and batching algorithms
 - Sequence Batching
 - Continuous Batching
- Allows multiple models to concurrently execute



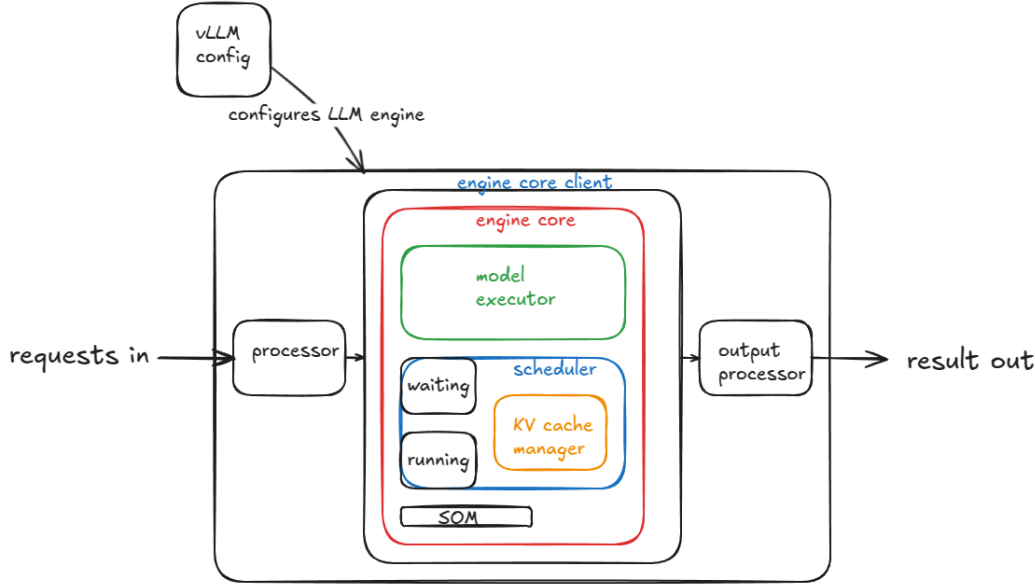
Example: TensorRT Inference Server



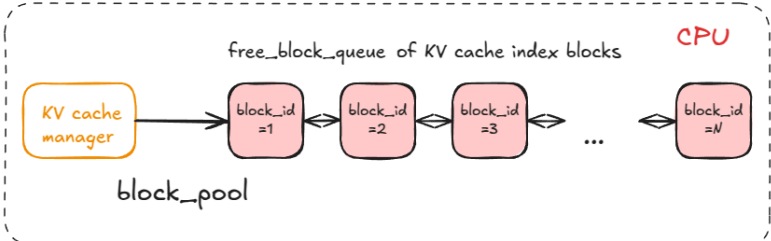
- Separates implementation of serving layer and execution layer
- Implements scheduling and batching algorithms
 - Sequence Batching
 - Continuous Batching
- Allows multiple models to concurrently execute
- Supports multiple frameworks
 - PyTorch
 - TensorFlow
 - ONNX
 - vLLM backend



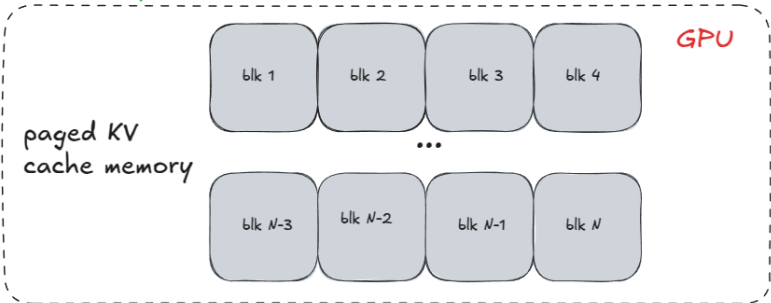
Example: vLLM Server



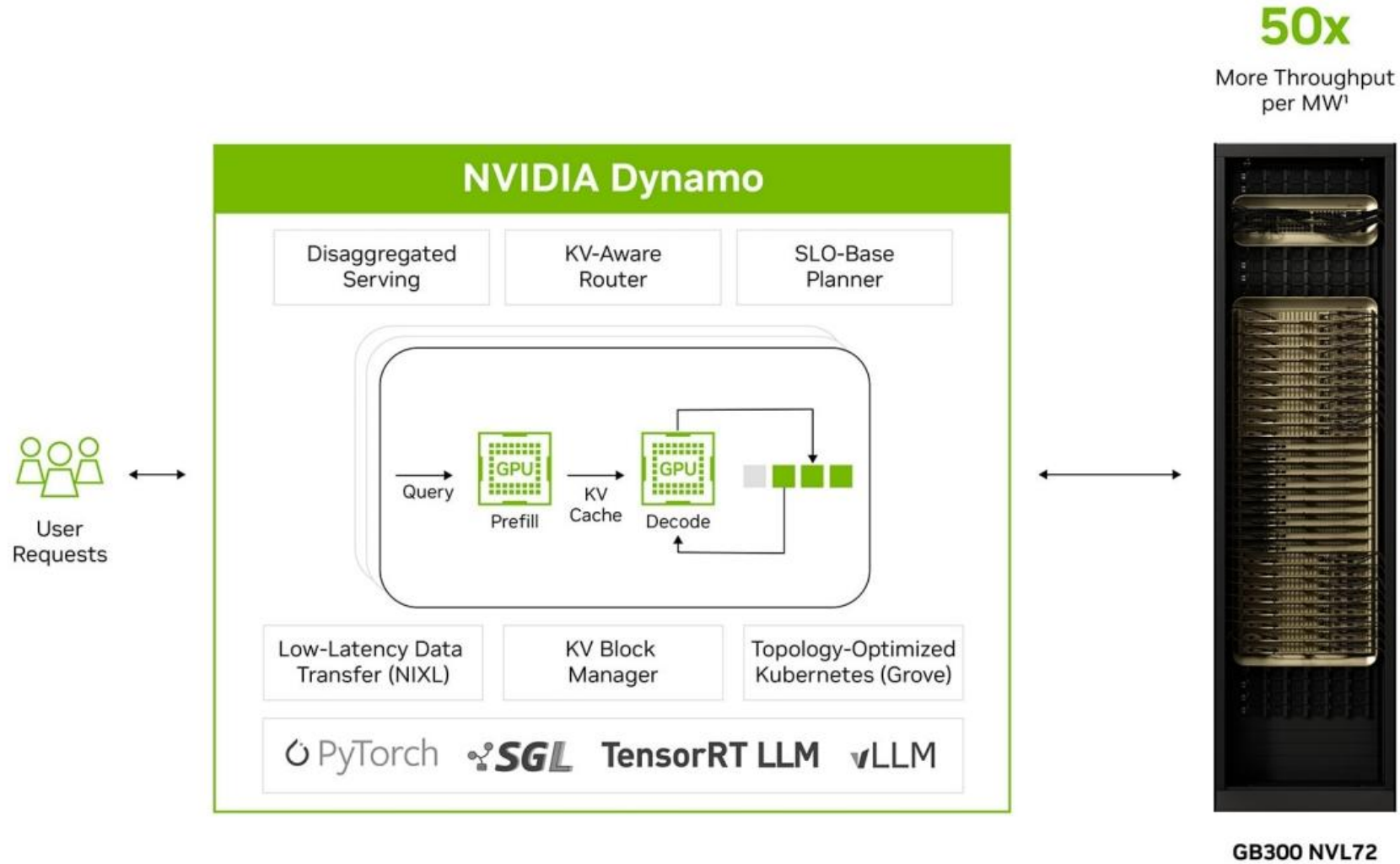
indexing structure



actual KV memory



Example: NVIDIA Dynamo



¹ - SemiAnalysis InferenceXv2 02/16/2026. NVIDIA GB300 NVL72 vs. H200. DeepSeek R1, TensorRT-LLM, 1k/1k

DL Inference

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

DL Inference

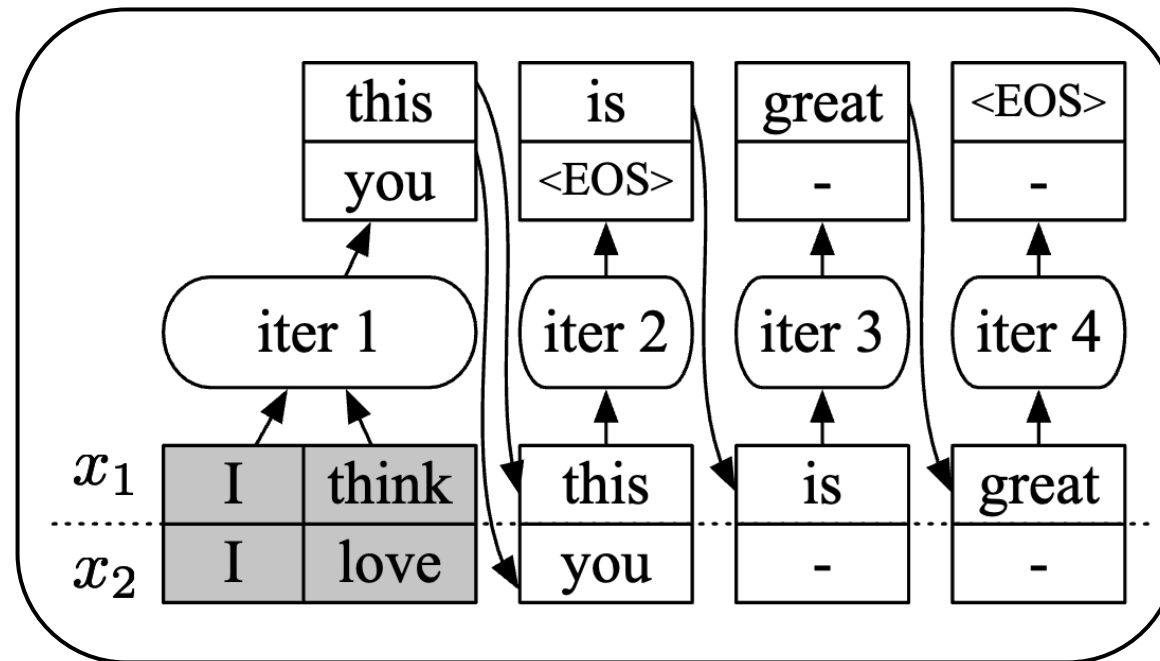
- LLM Serving System
- Continuous Batching
 - Sequence batching
 - Continuous batching

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

Problem 1: Request Level Scheduling



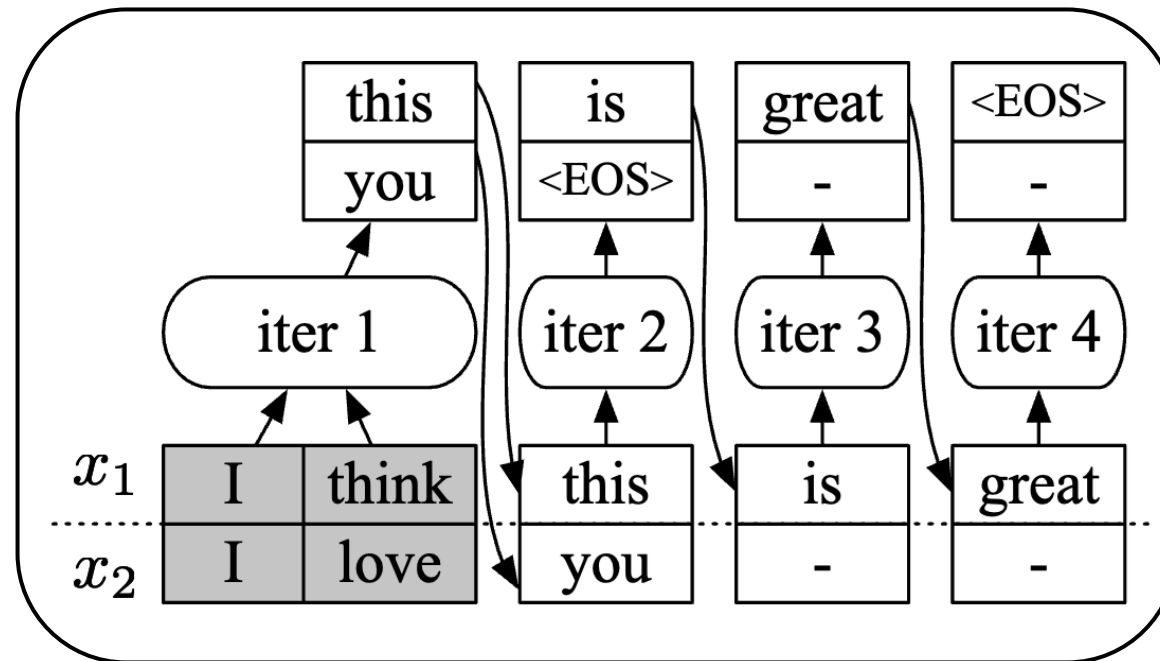
Execution Engine



Problem 1: Request Level Scheduling

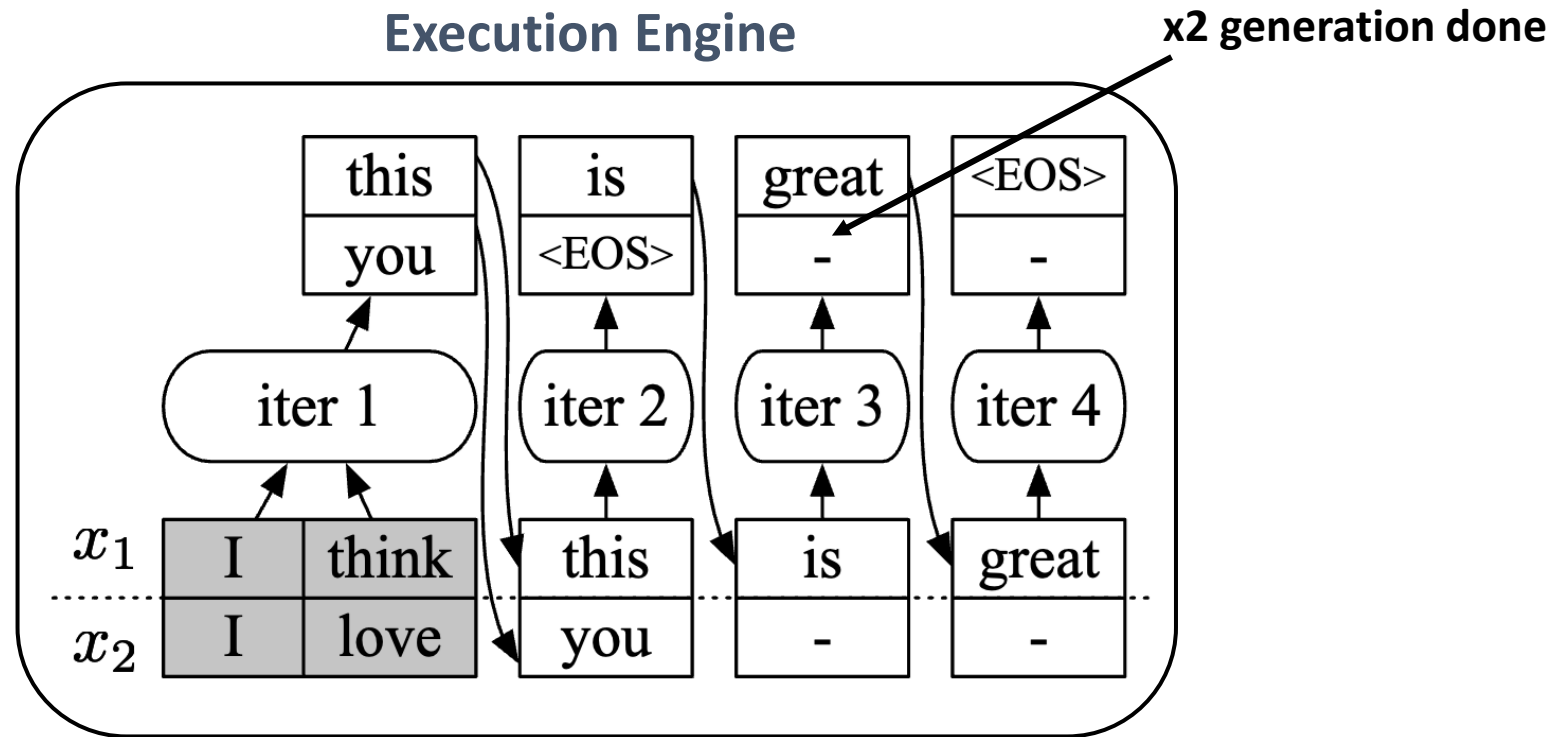


Execution Engine

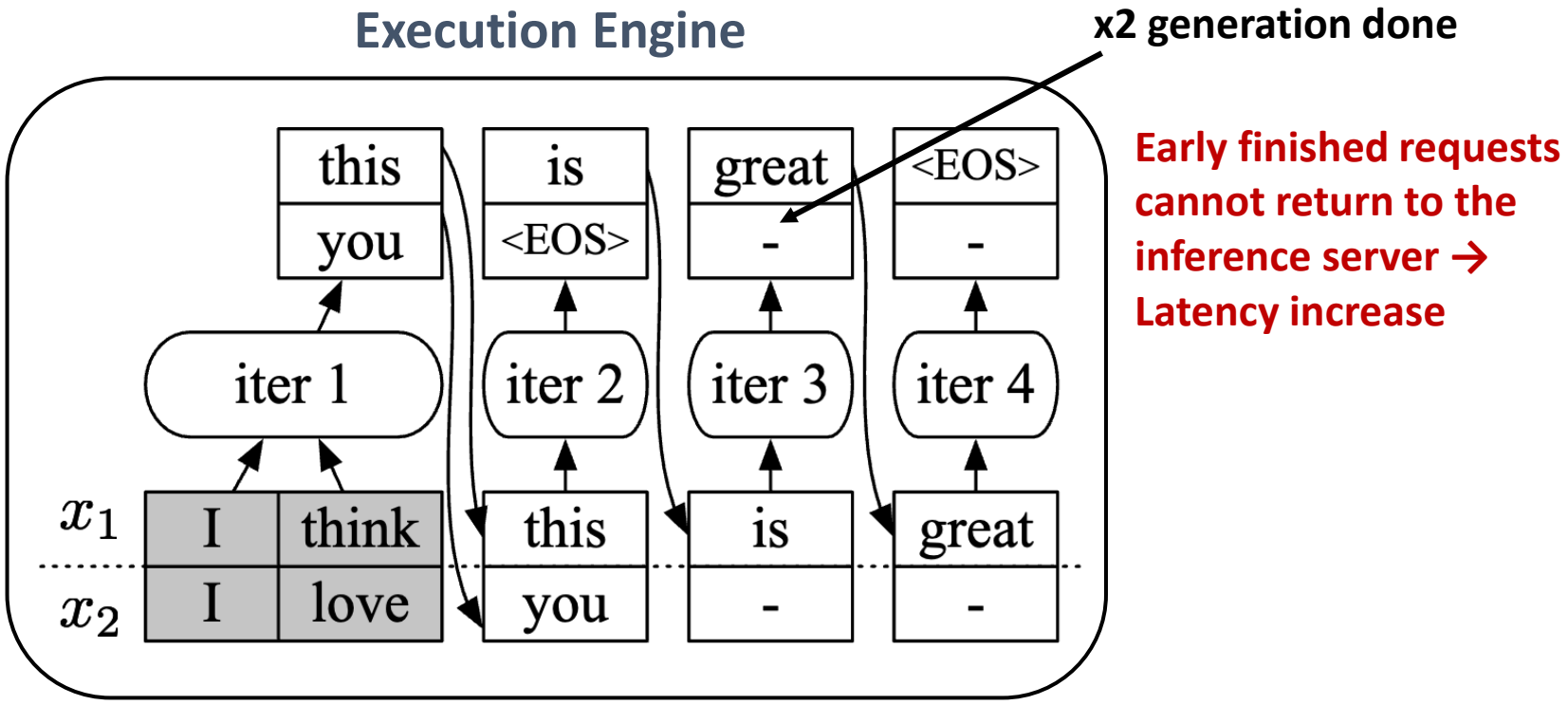


Do you see any problem?

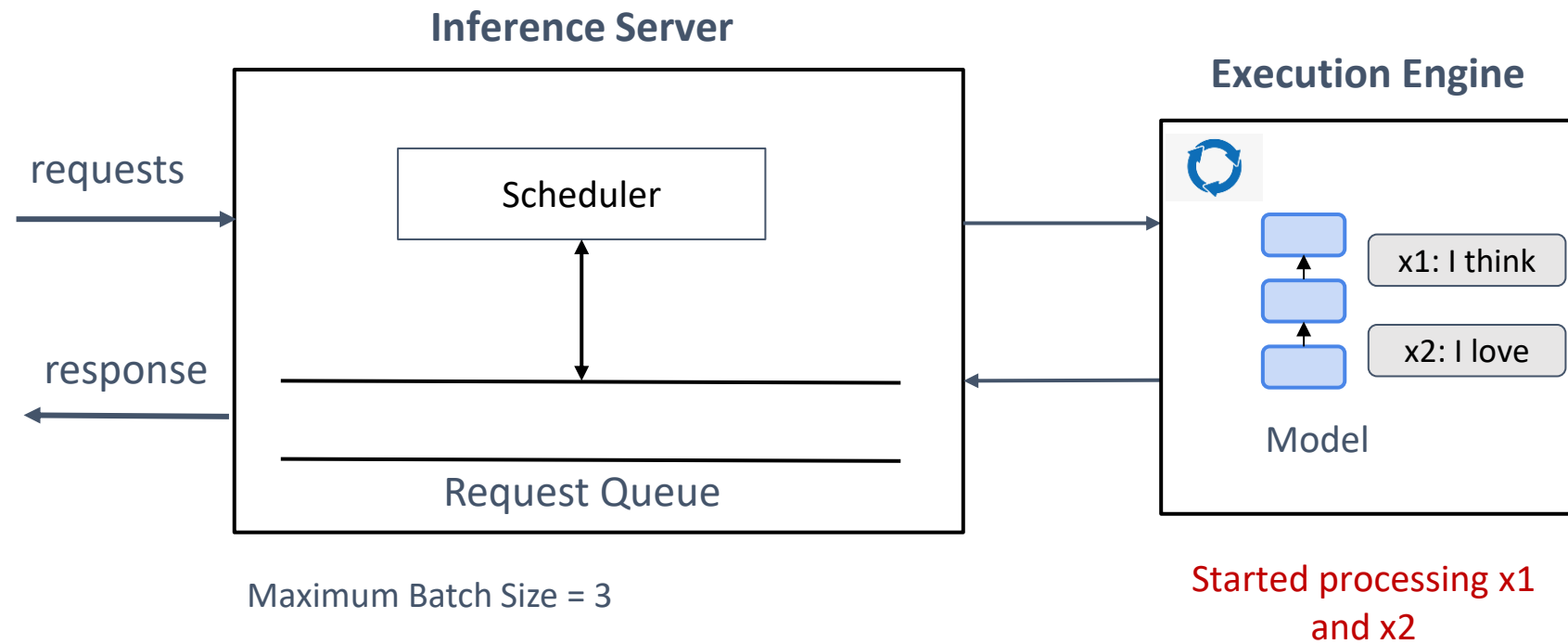
Problem 1: Request Level Scheduling



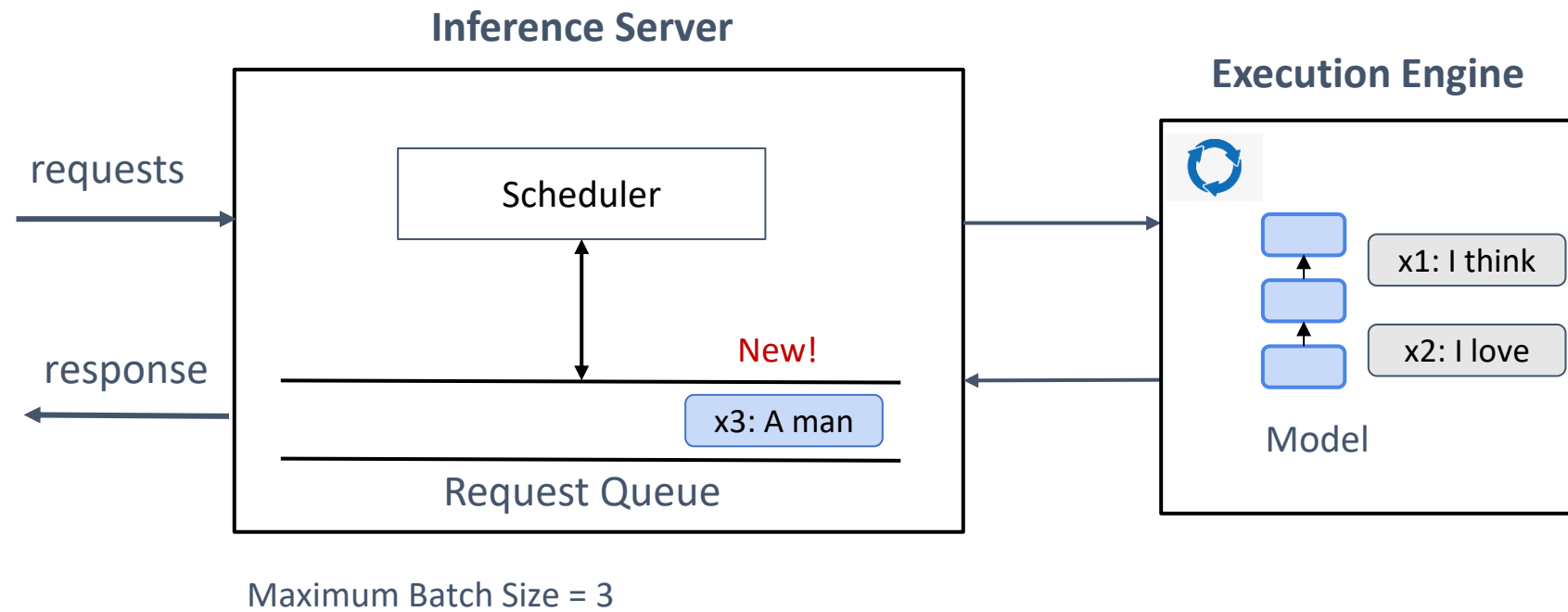
Problem 1: Request Level Scheduling



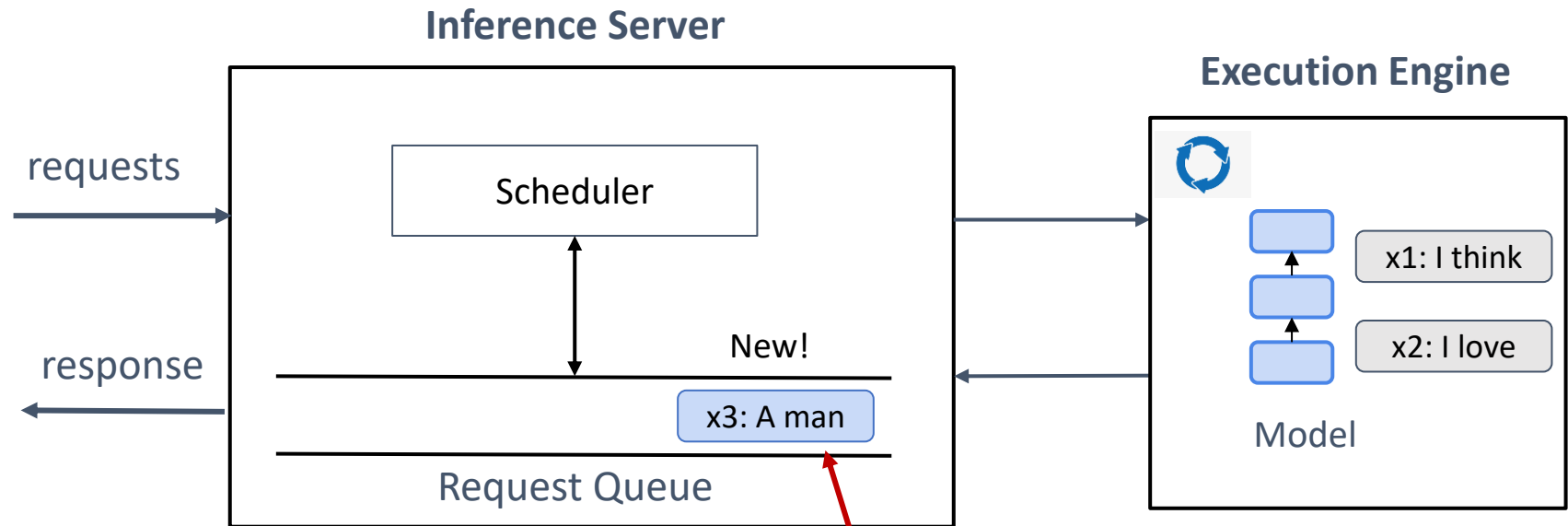
Problem 1: Request Level Scheduling



Problem 1: Request Level Scheduling



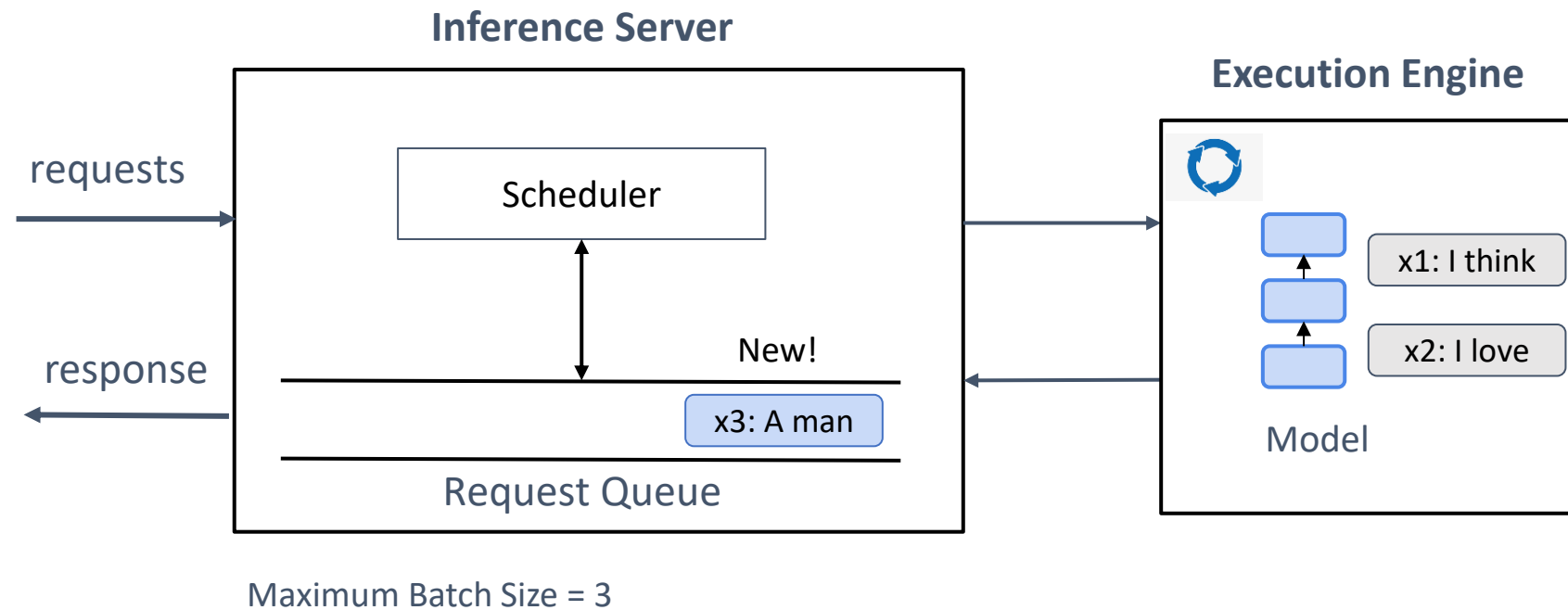
Problem 1: Request Level Scheduling



Maximum Batch Size = 3

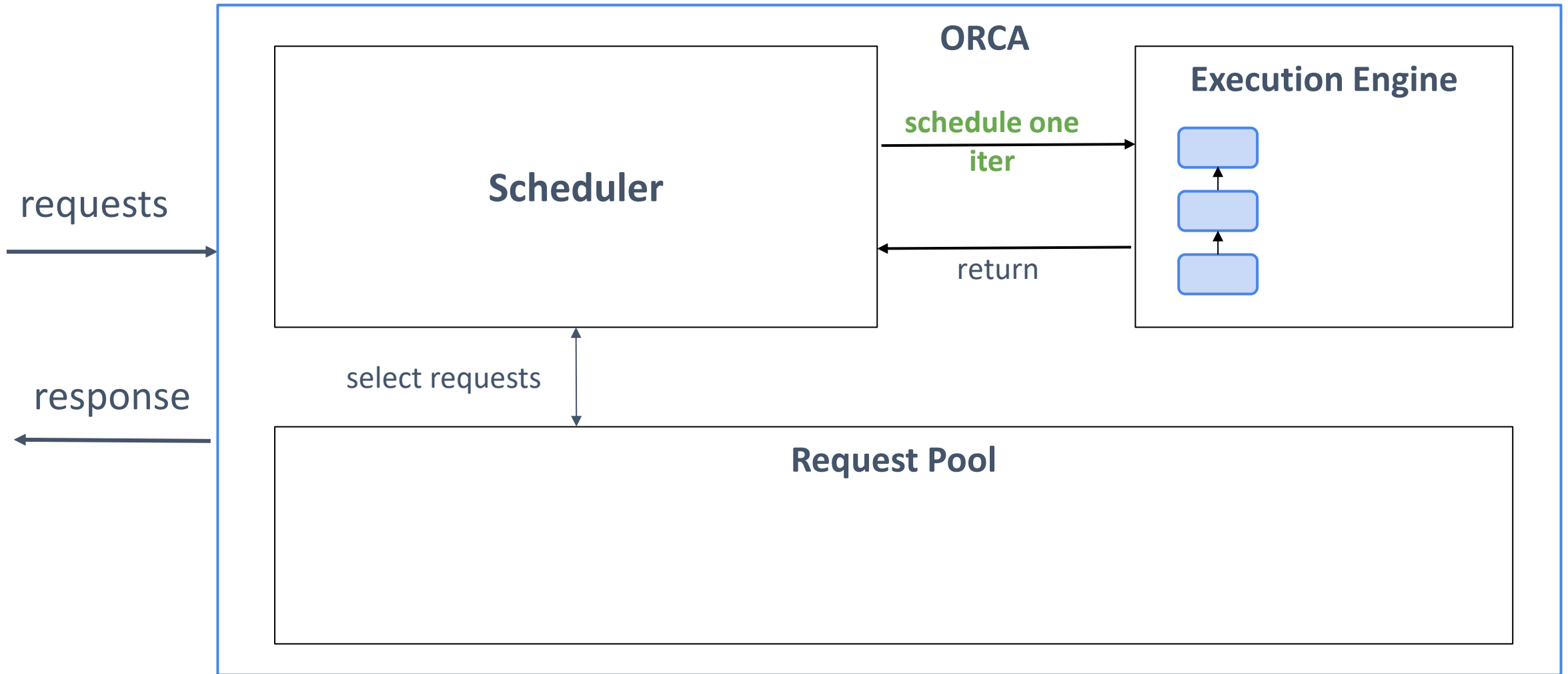
**Late join requests need to wait until engine finishes execution
→ Latency Increase**

Problem 1: Request Level Scheduling

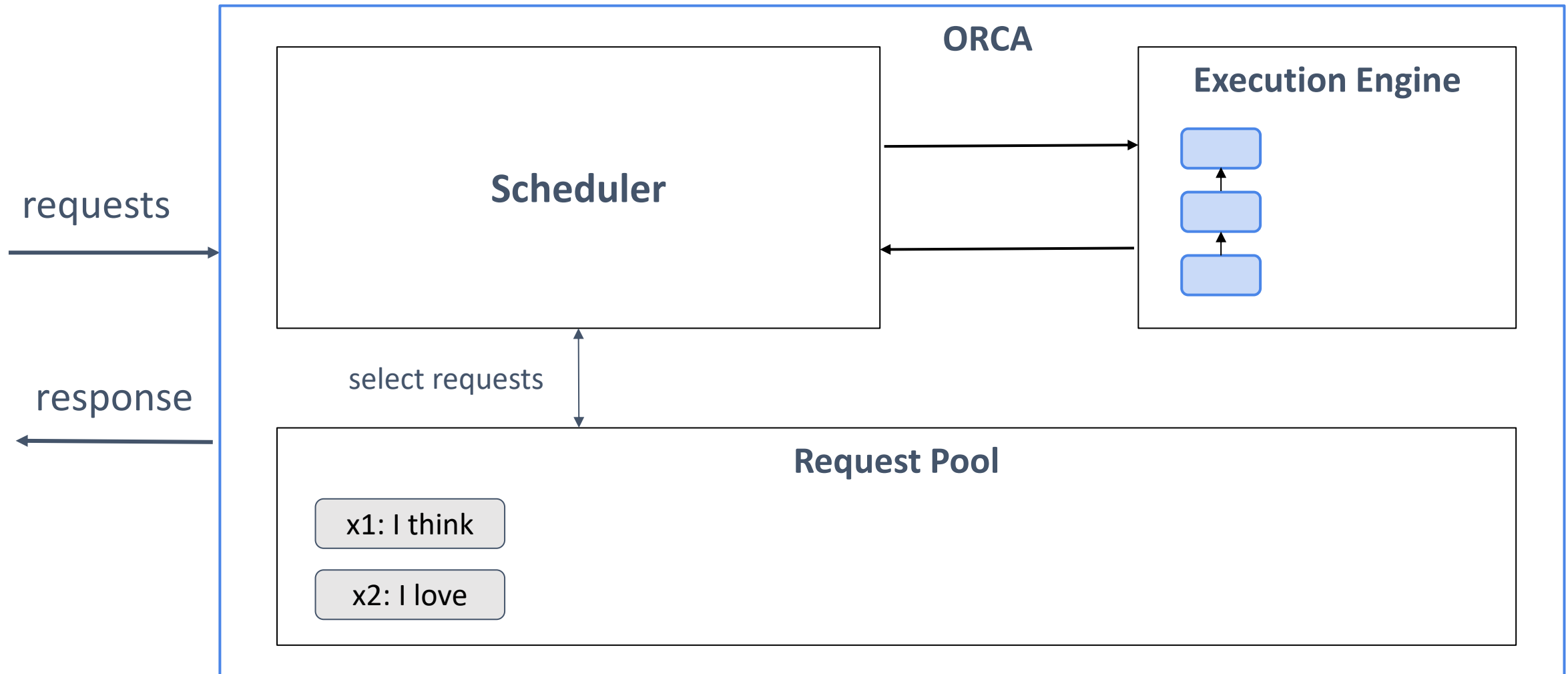


Question: How can we avoid redundant computation and ensure late-arriving requests to be processed more promptly?

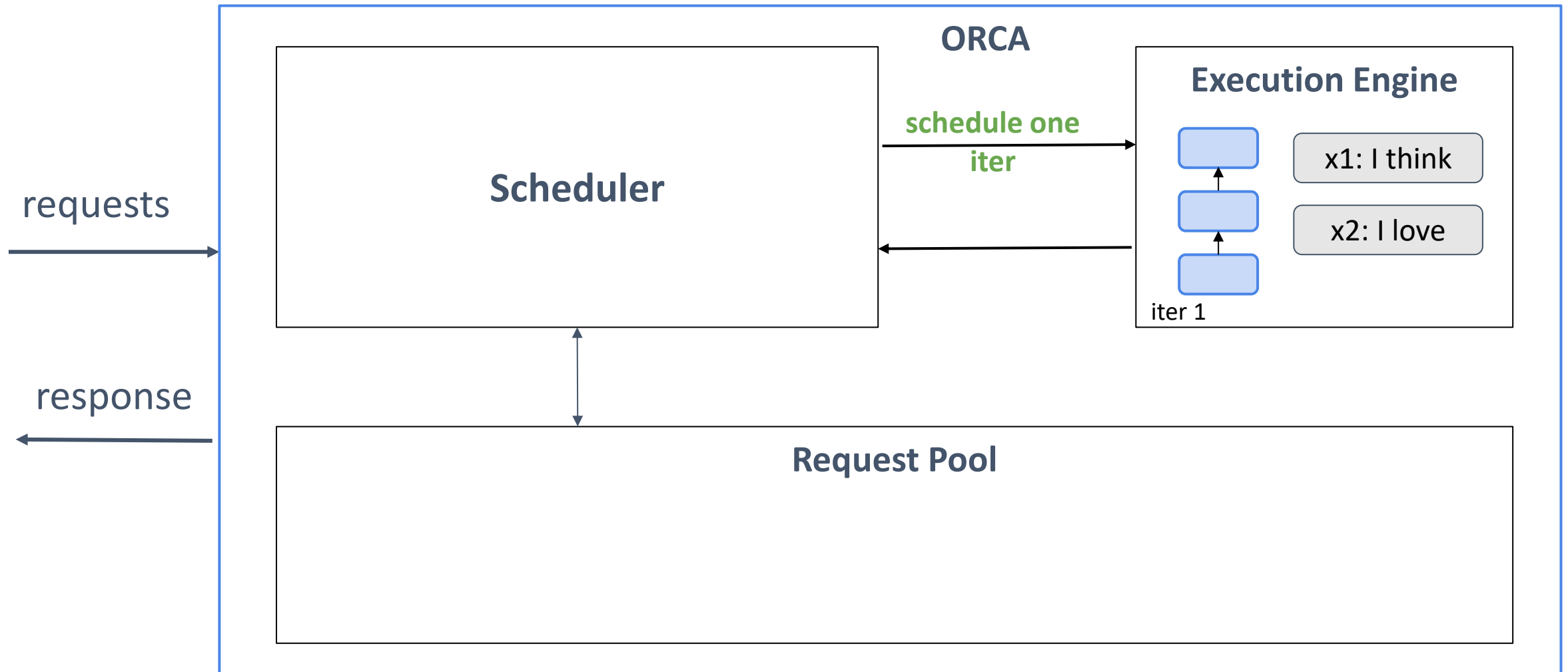
Solution 1: Iteration Level Scheduling



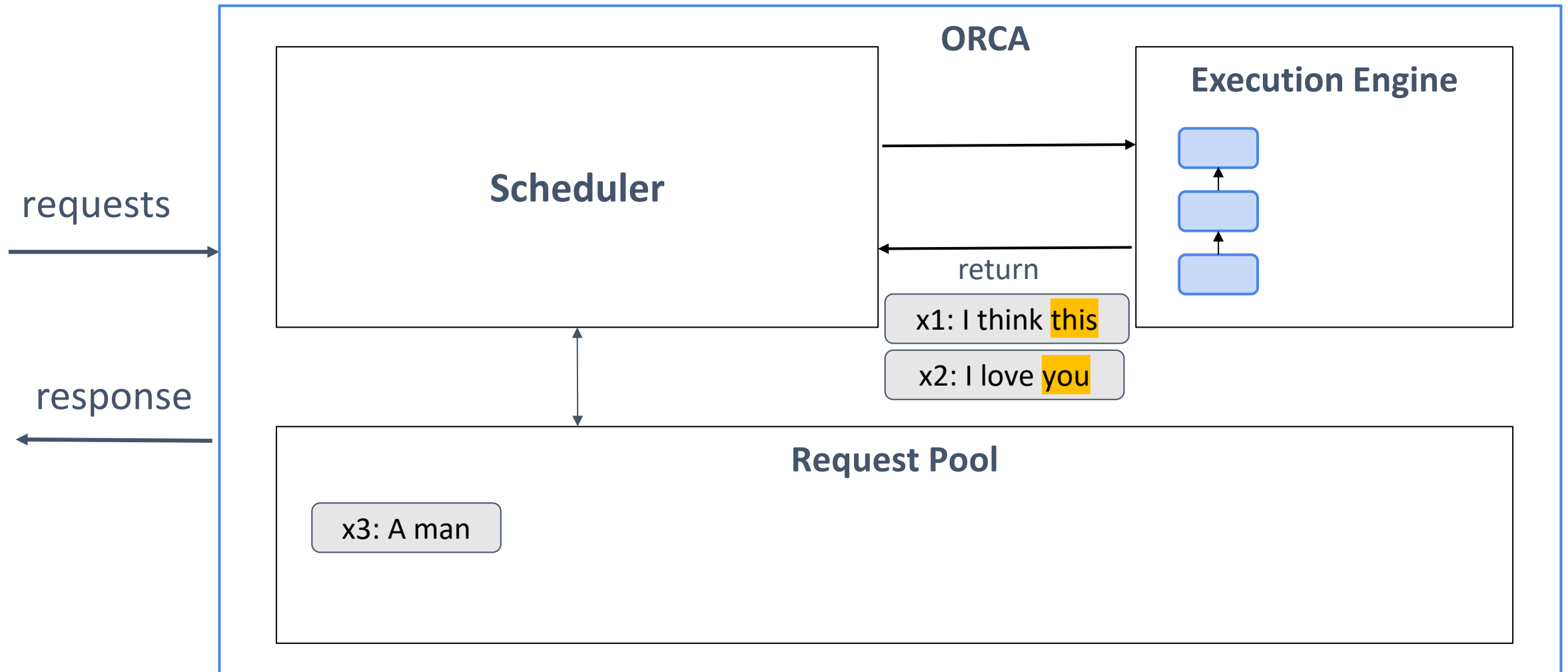
Solution 1: Iteration Level Scheduling



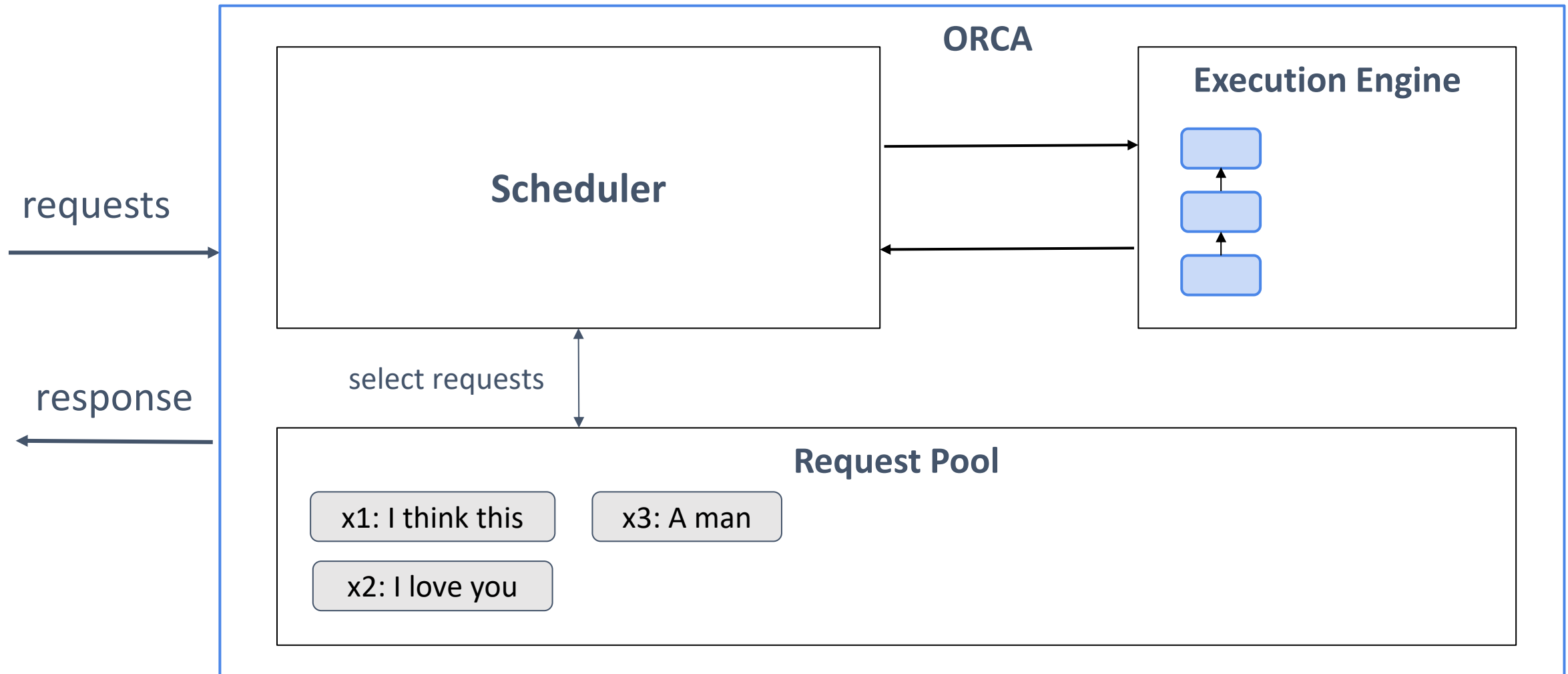
Solution 1: Iteration Level Scheduling



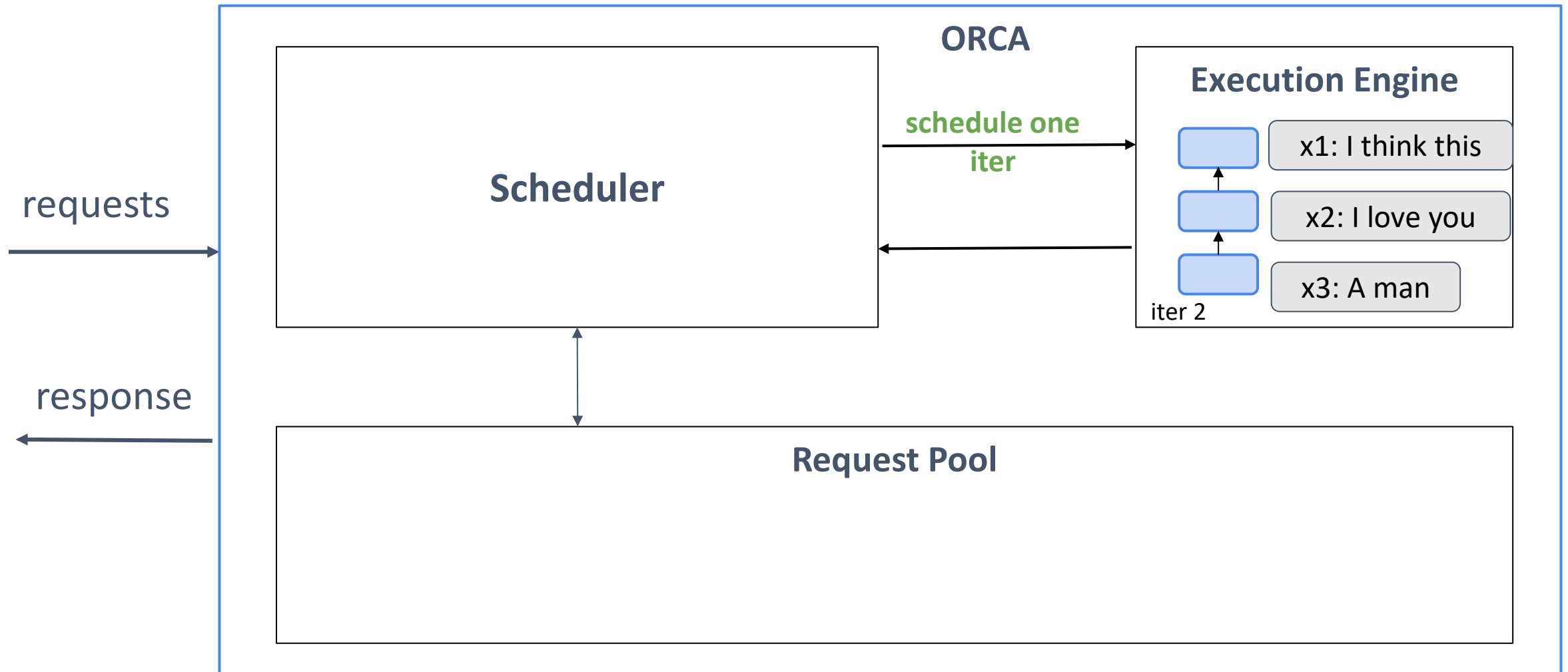
Solution 1: Iteration Level Scheduling



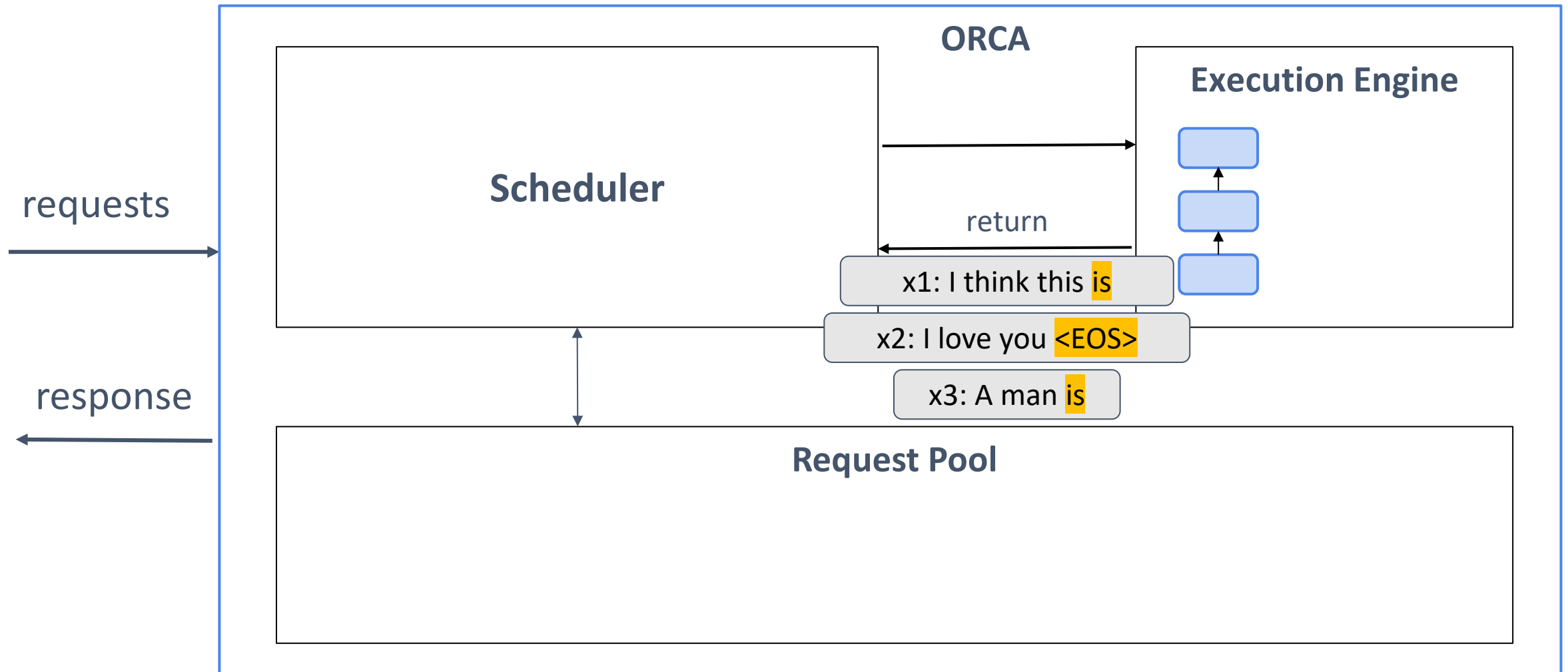
Solution 1: Iteration Level Scheduling



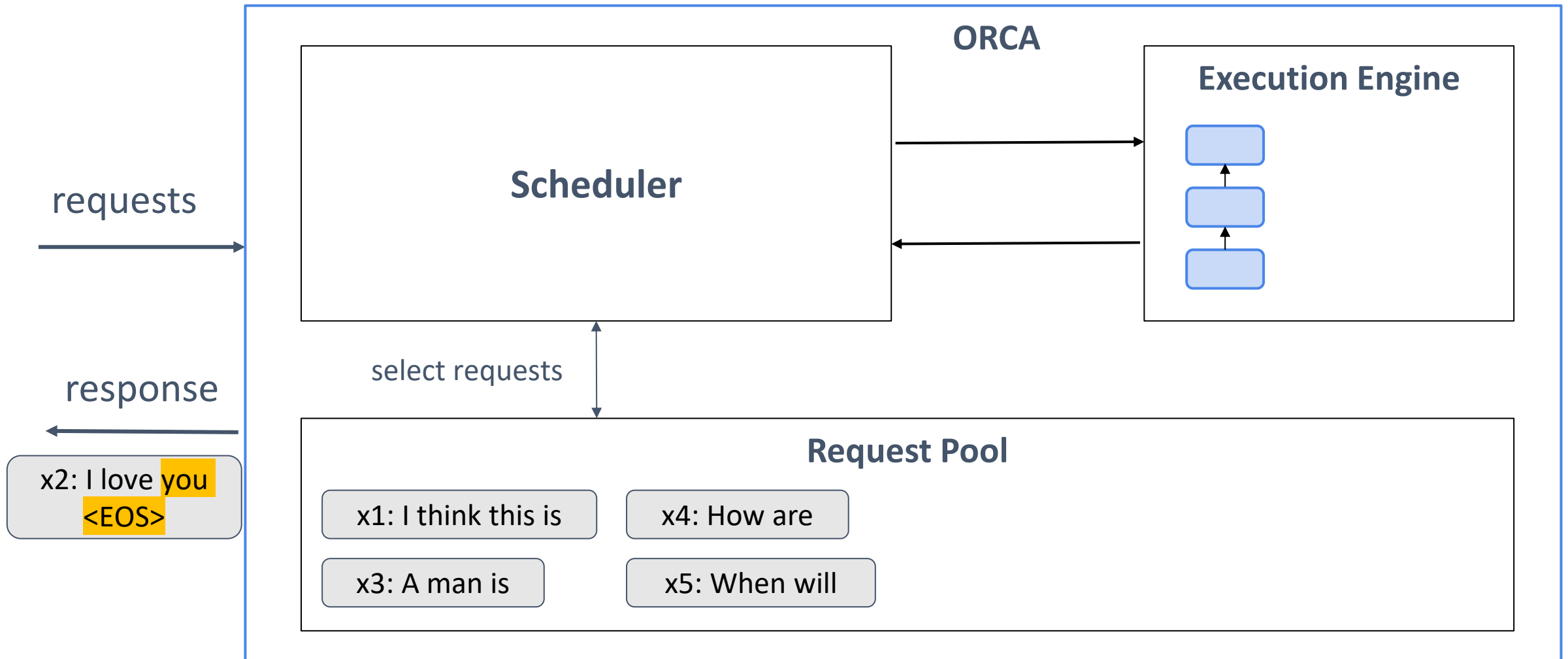
Solution 1: Iteration Level Scheduling



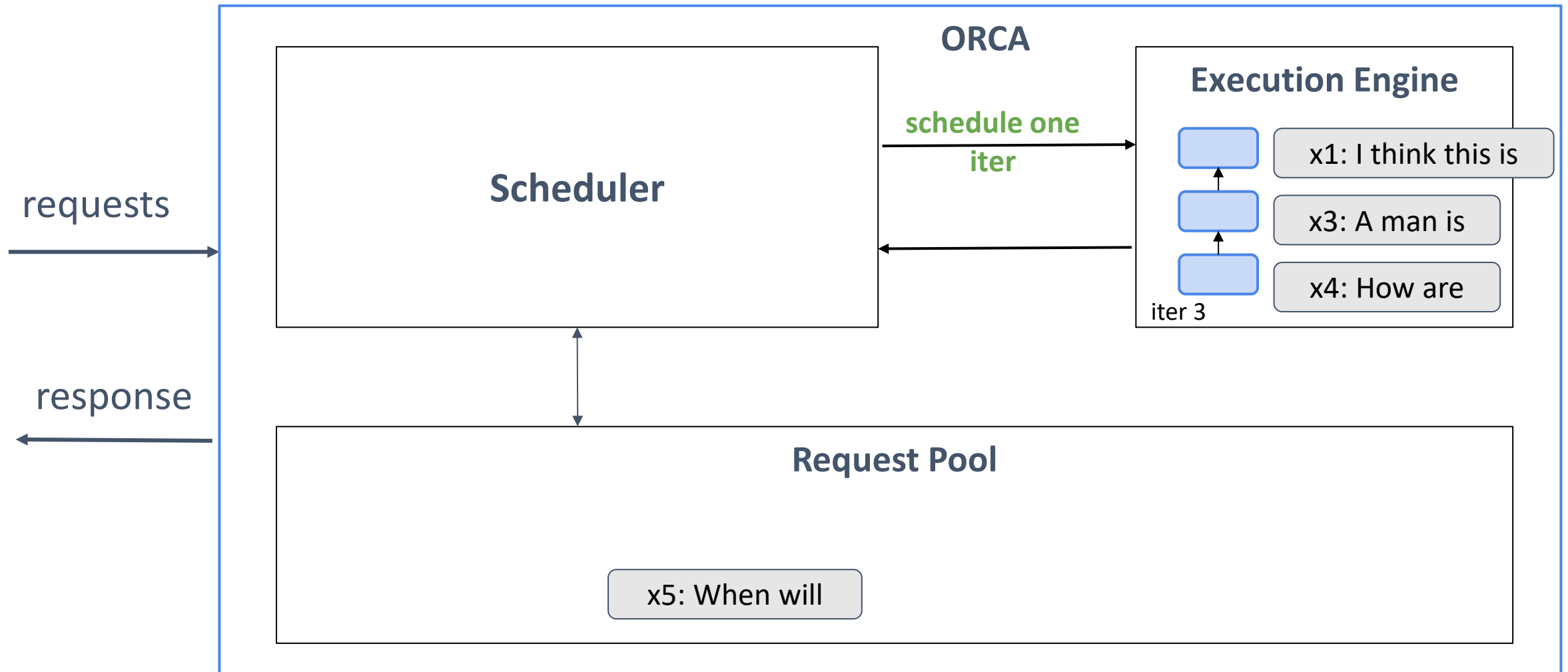
Solution 1: Iteration Level Scheduling



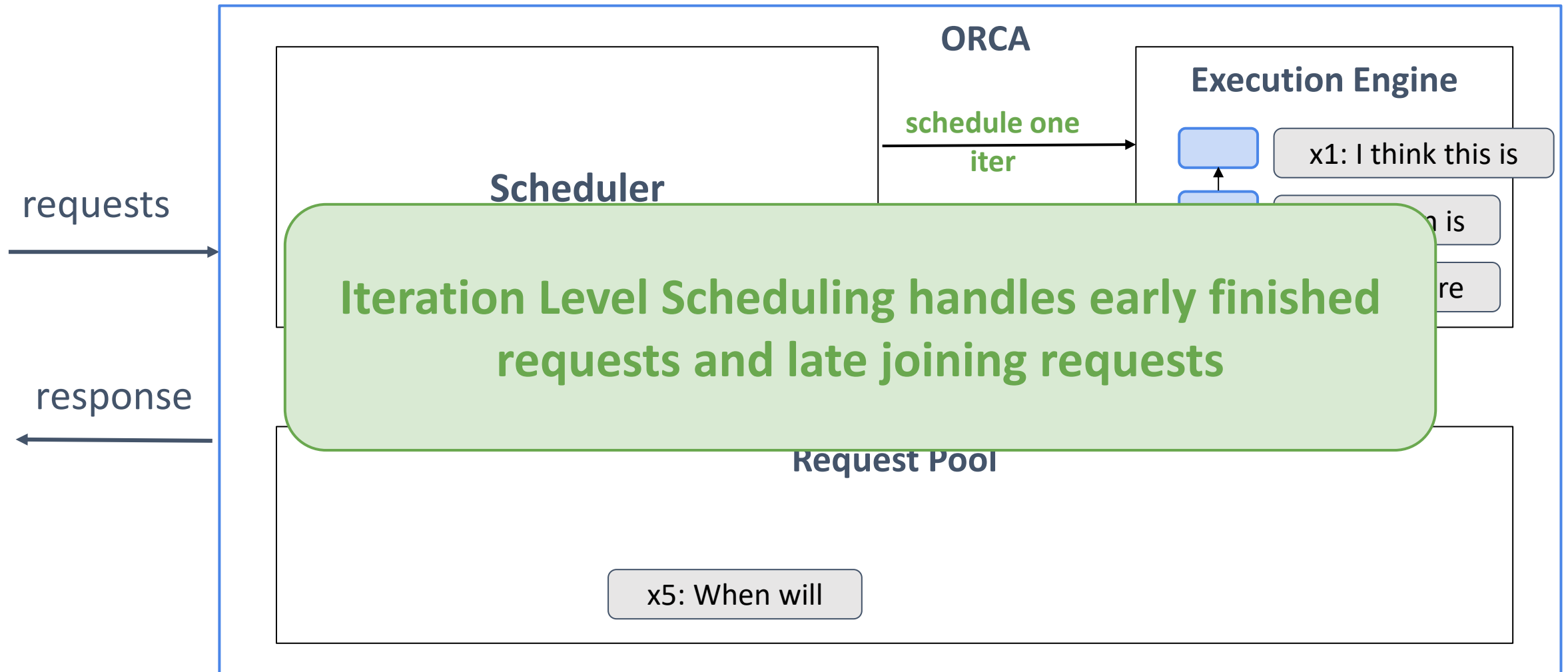
Solution 1: Iteration Level Scheduling



Solution 1: Iteration Level Scheduling



Solution 1: Iteration Level Scheduling



Problem 2: How to Batch Requests?

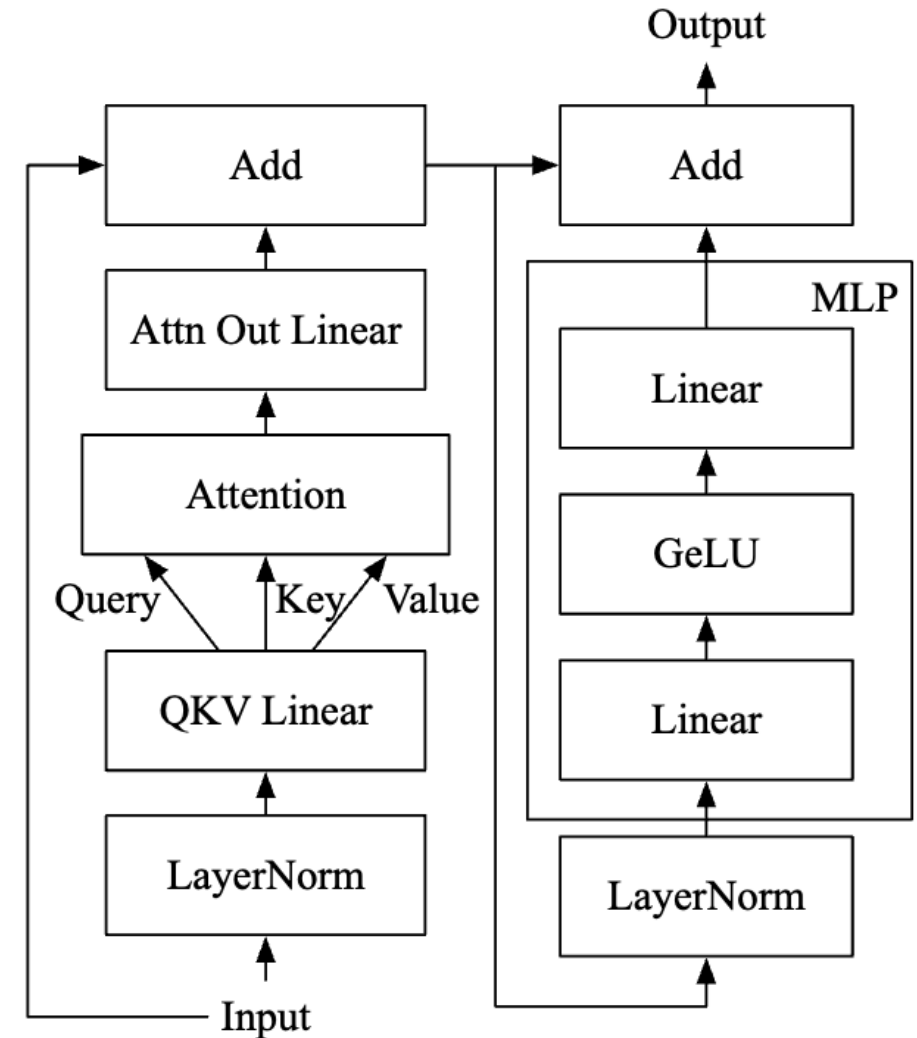


Let's assume Batch Size $B = 1$

Input Dimension: $[L \times H]$ (L=sequence length, H=hidden dim.)

Attention Operation:

1. $QK^T : [L \times H] \times [H \times L] \rightarrow [L \times L]$
2. $P = \text{softmax}(QK^T) : [L \times L]$
3. $O = PV : [L \times L] \times [L \times H] \rightarrow [L \times H]$



Problem 2: How to Batch Requests?



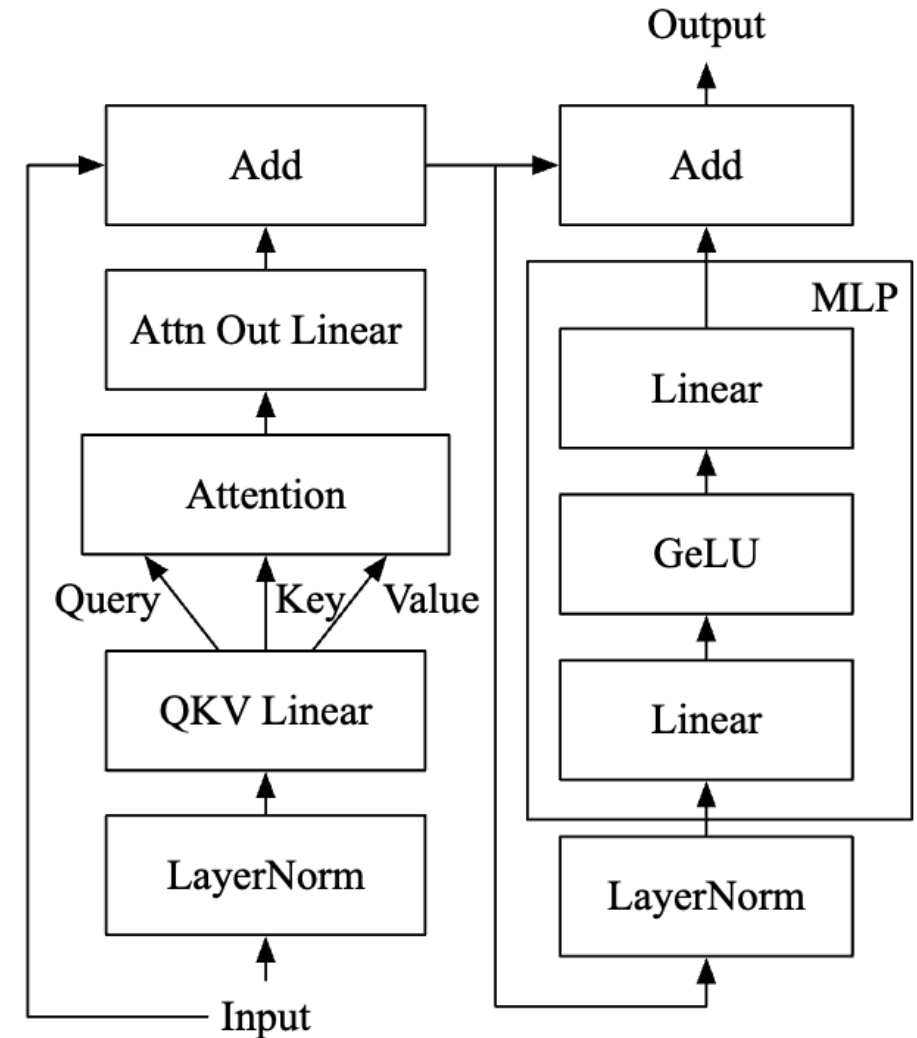
Let's assume Batch Size $B = 1$

Input Dimension: $[L \times H]$ (L=sequence length, H=hidden dim.)

Attention Operation:

1. $QK^T : [L \times H] \times [H \times L] \rightarrow [L \times L]$
2. $P = \text{softmax}(QK^T) : [L \times L]$
3. $O = PV : [L \times L] \times [L \times H] \rightarrow [L \times H]$

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



Problem 2: How to Batch Requests?



Let's assume Batch Size $B = 1$

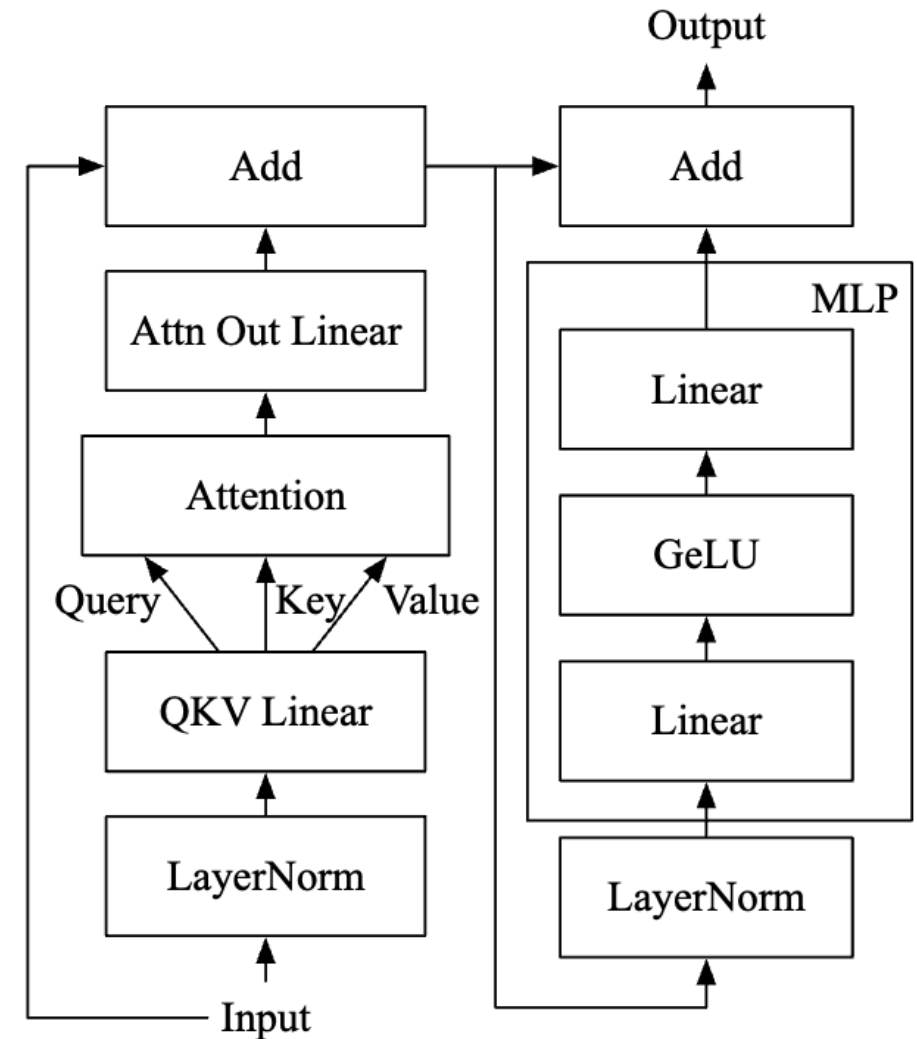
Input Dimension: $[L \times H]$ (L=sequence length, H=hidden dim.)

Attention Operation:

1. $QK^T : [L \times H] \times [H \times L] \rightarrow [L \times L]$
2. $P = \text{softmax}(QK^T) : [L \times L]$
3. $O = PV : [L \times L] \times [L \times H] \rightarrow [L \times H]$

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

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

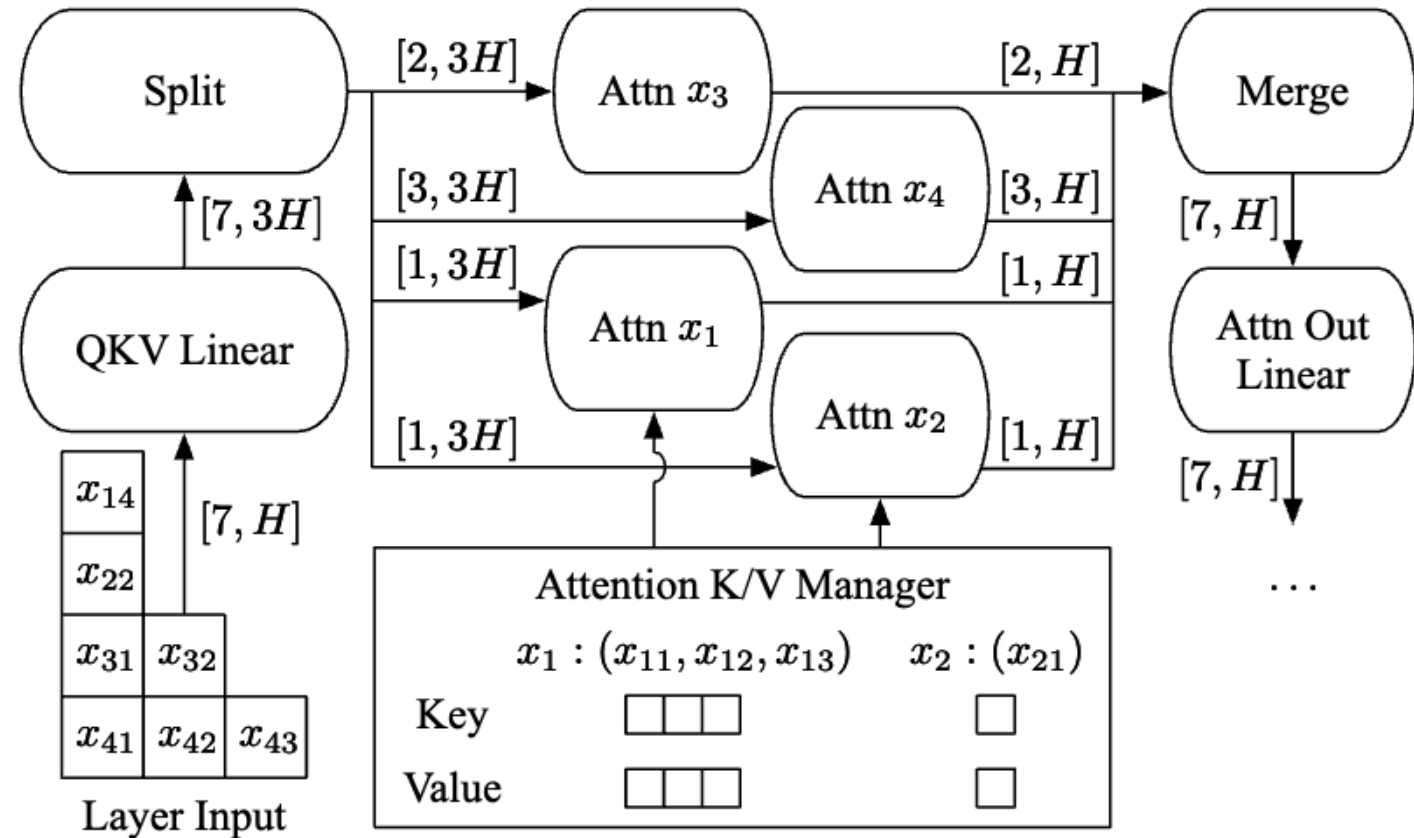


Solution 2: Selective Batching



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

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



Solution 2: Selective Batching



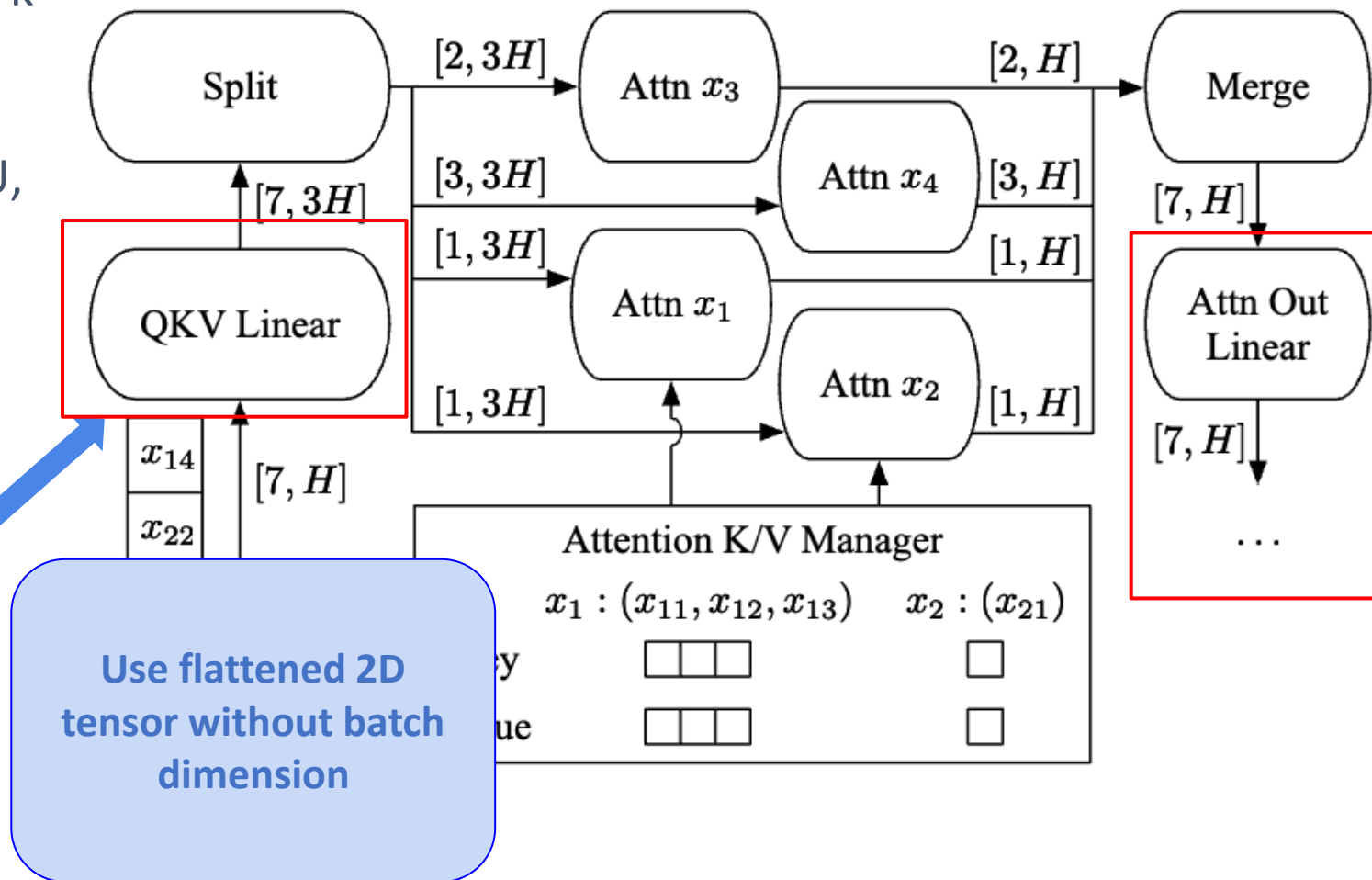
Only **Attention operation** does not work with batching tensors with diff. L_i

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

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

$x_1: [1, H]$
 $x_2: [1, H]$
 $x_3: [2, H]$
 $x_4: [3, H]$

**[7, H]
tensor**



Solution 2: Selective Batching



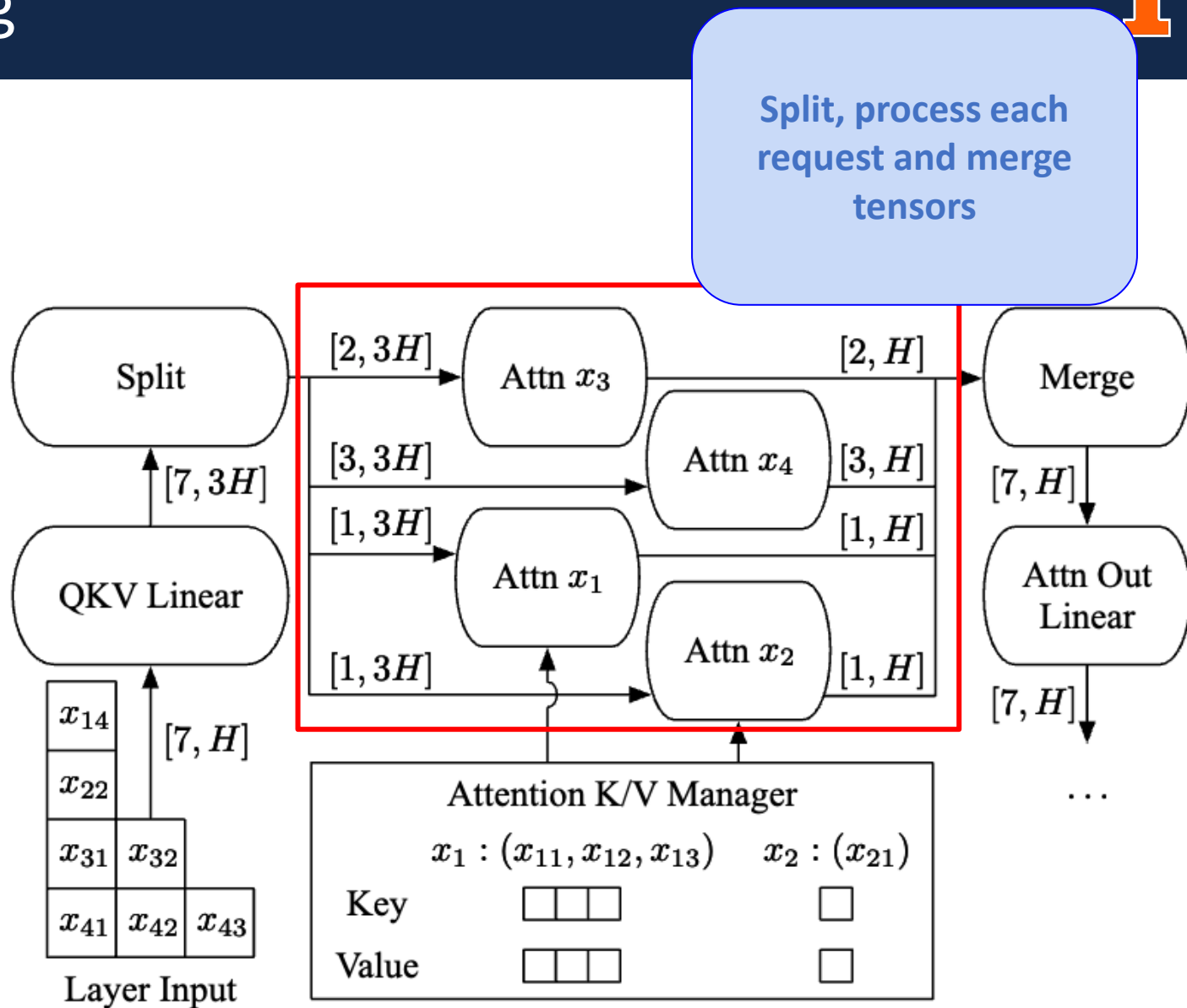
Only **Attention operation** does not work with batching tensors with diff. L_i

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

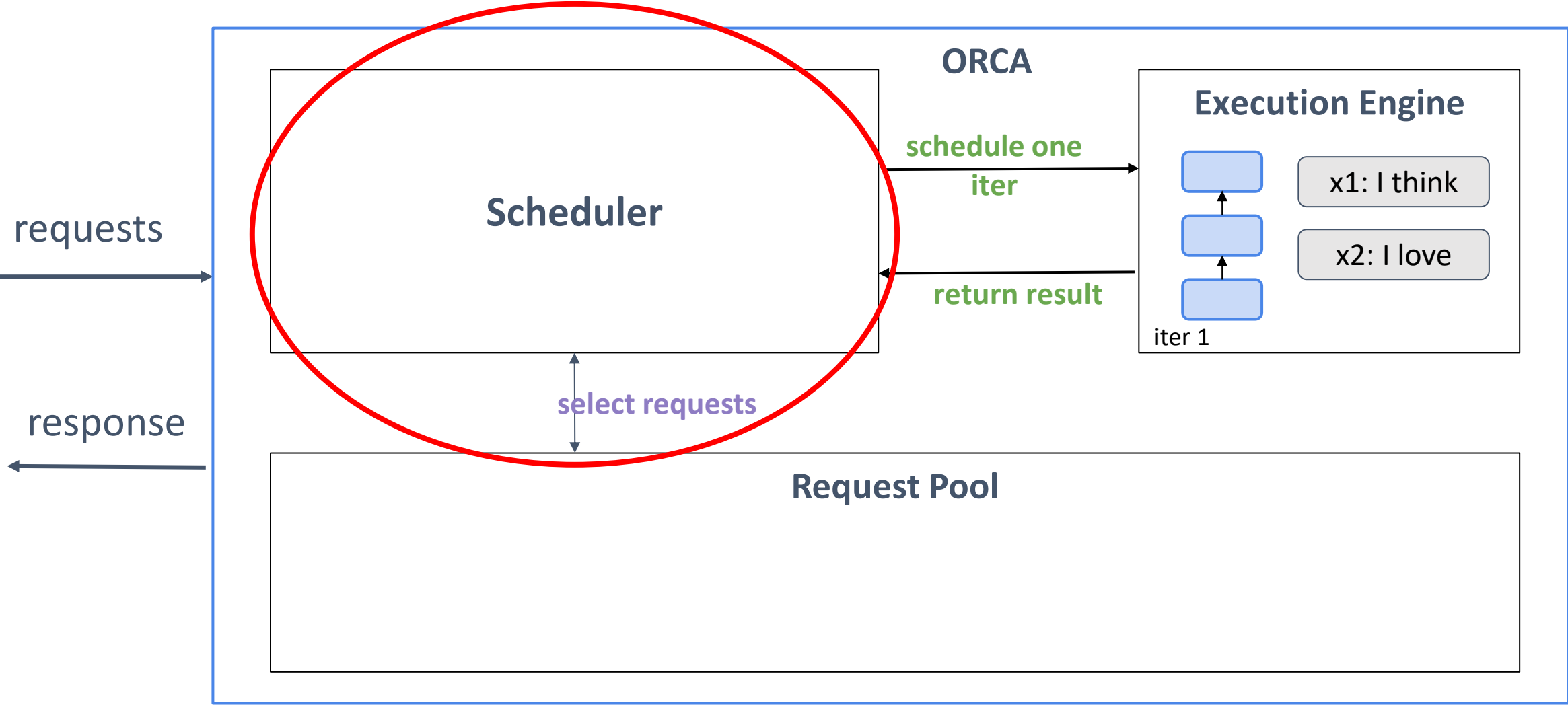
Coalesce $[L_i, H]$ tensor to $[\sum L_i, H]$ for batching

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

→ **[7,H] tensor**



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
- ...

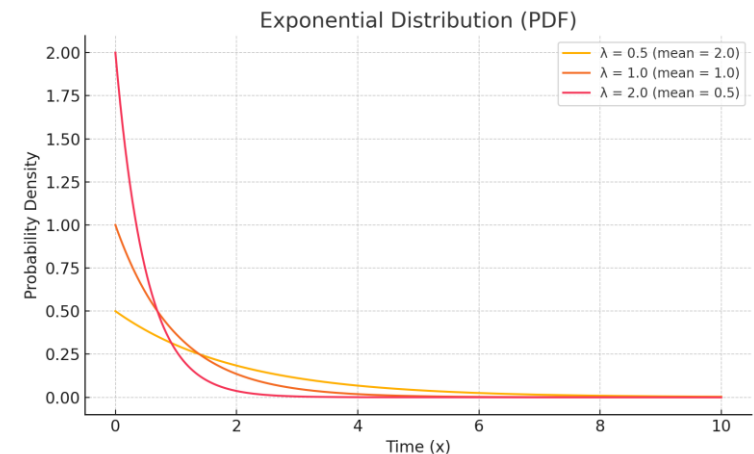
- Question: When would continuous batching provide more benefits than sequence batching?

- 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

- 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

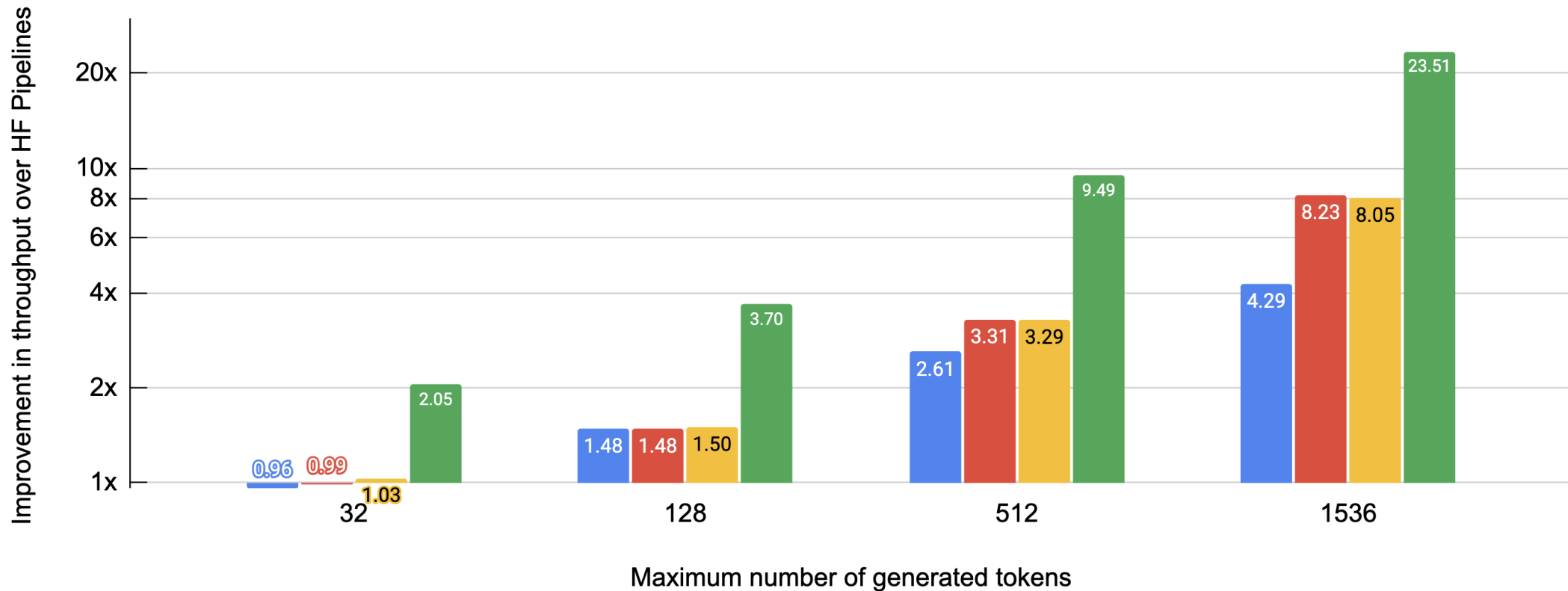


Throughput Improvement from Continuous Batching



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)



- vLLM uses PagedAttention – extra batch size space

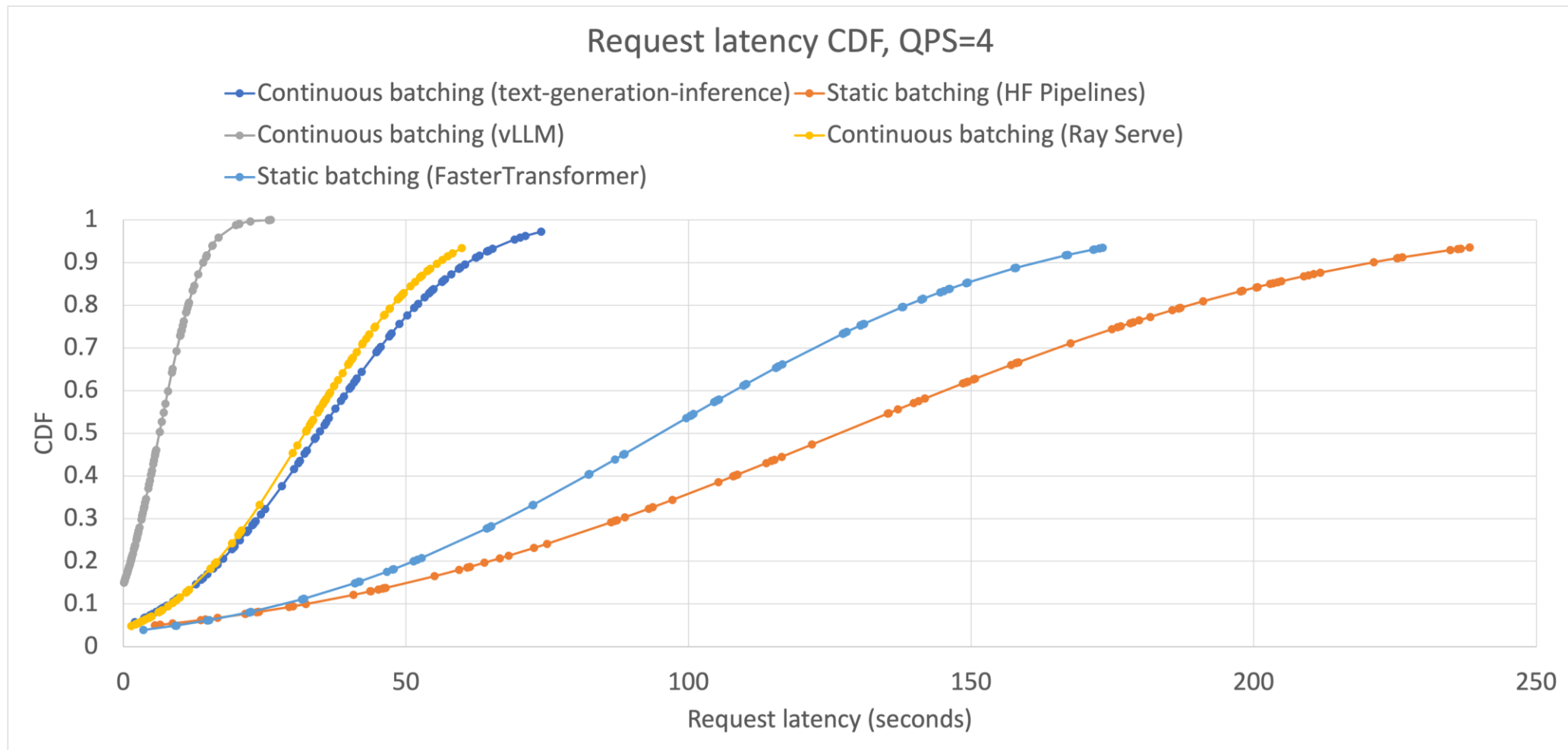


Contiguous Memory



Non-Contiguous Memory

E2E Latency Experiments: Results

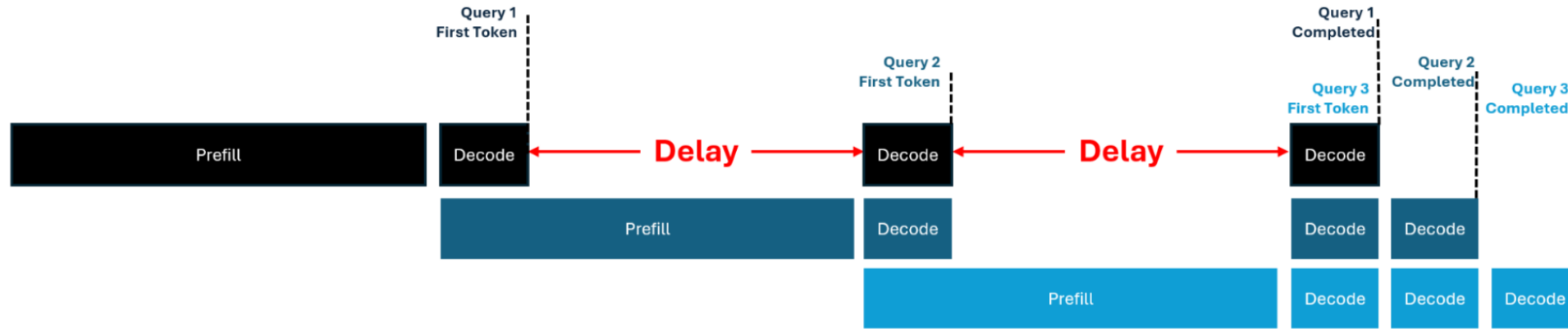


- Continuous batching handles early-finished and late-arrived requests more efficiently
- Fills GPU capacity after each token generation
- As variance in sequence length increases, continuous batching increases GPU utilization

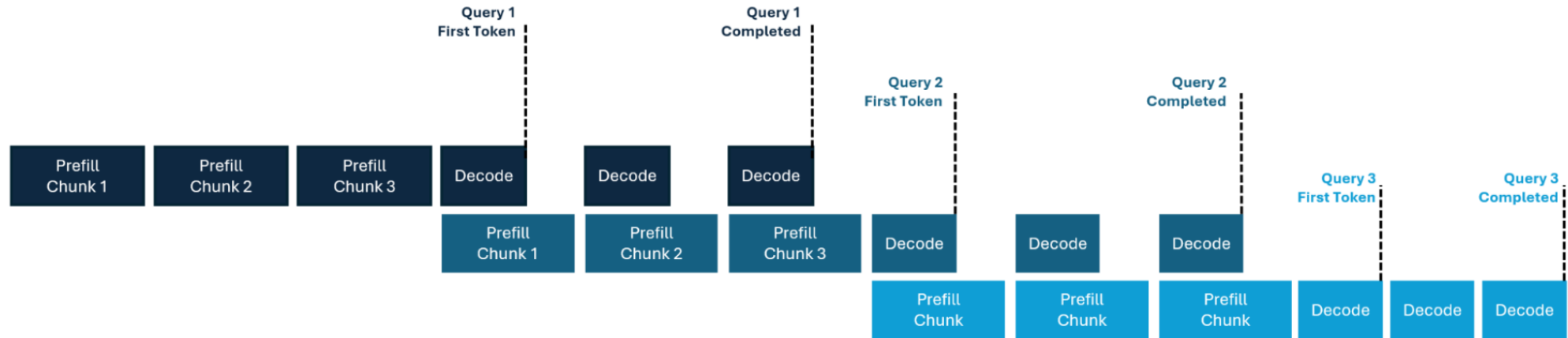
LLM Inference Scheduler: Chunked Prefill



W/O TensorRT-LLM Chunked Prefill



W/ TensorRT-LLM Chunked Prefill



■ Query#1 ■ Query#2 ■ Query#3

Visual for illustration purposes only. In production environments multiple Prefills are processed in parallel (not depicted). With TensorRT-LLM Inflight Batching queries are evicted once completed and new queries are added (not depicted).

Questions?

Sequence Batching (Static Batching)



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

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1				
S_2	S_2	S_2					
S_3	S_3	S_3	S_3				
S_4	S_4	S_4	S_4	S_4			

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1	S_1	END		
S_2	S_2	S_2	S_2	S_2	S_2	S_2	END
S_3	S_3	S_3	S_3	END			
S_4	S_4	S_4	S_4	S_4	S_4	END	

Legend:

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

Continuous Batching



Top: static batching
Bottom: continuous batching

Legend:

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

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1				
S_2	S_2	S_2					
S_3	S_3	S_3					
S_4	S_4	S_4					

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1	S_1	END		
S_2	S_2	S_2	S_2	S_2	S_2	S_2	END
S_3	S_3	S_3	S_3	END			
S_4	S_4	S_4	S_4	S_4	S_4	END	

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1				
S_2	S_2	S_2					
S_3	S_3	S_3					
S_4	S_4	S_4					

T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1	S_1	END	S_6	S_6
S_2	S_2	S_2	S_2	S_2	S_2	S_2	END
S_3	S_3	S_3	S_3	END	S_5	S_5	S_5
S_4	S_4	S_4	S_4	S_4	S_4	END	S_7

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