

Spring25 CS598YP

19.2: Orca

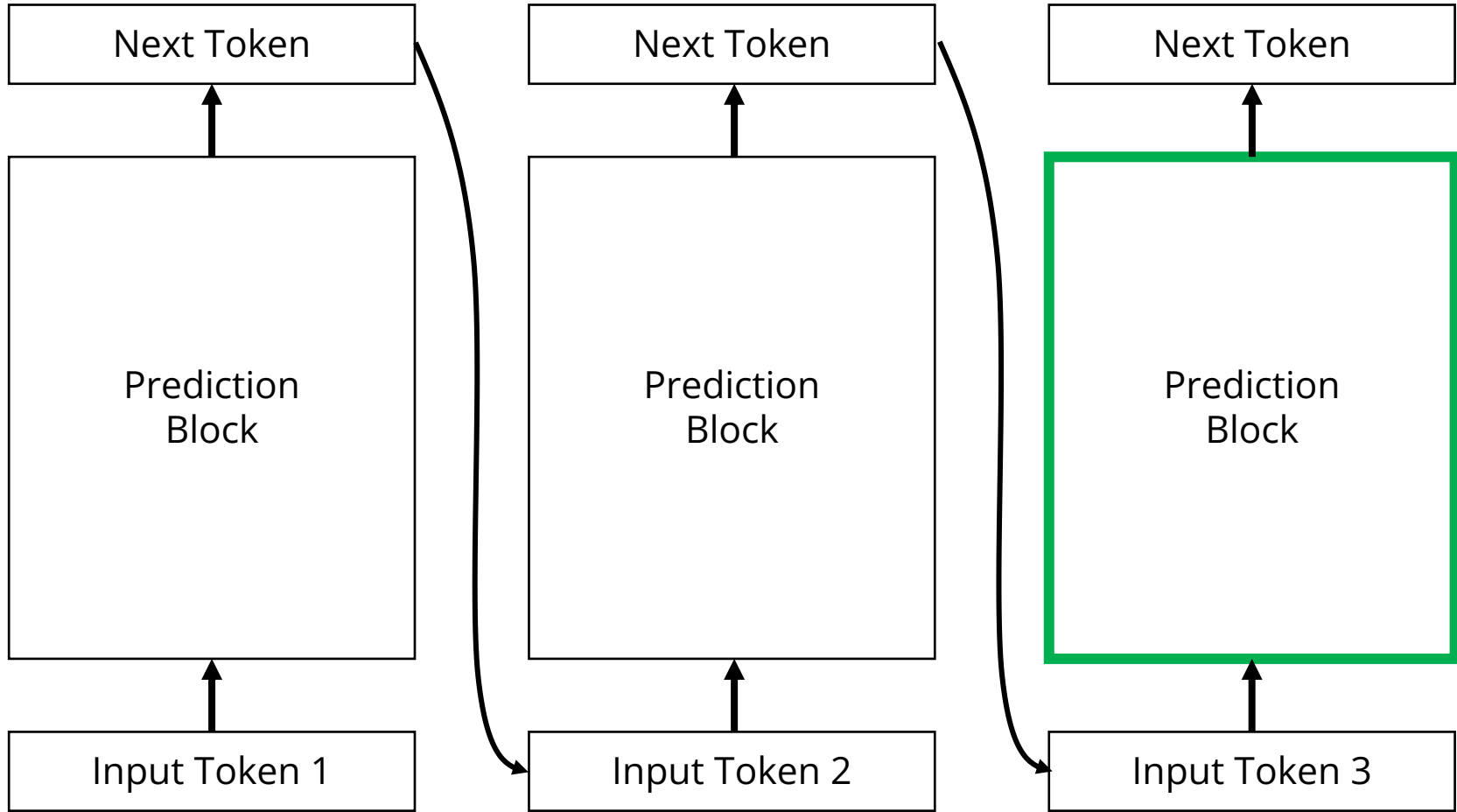
Yongjoo Park

University of Illinois Urbana-Champaign

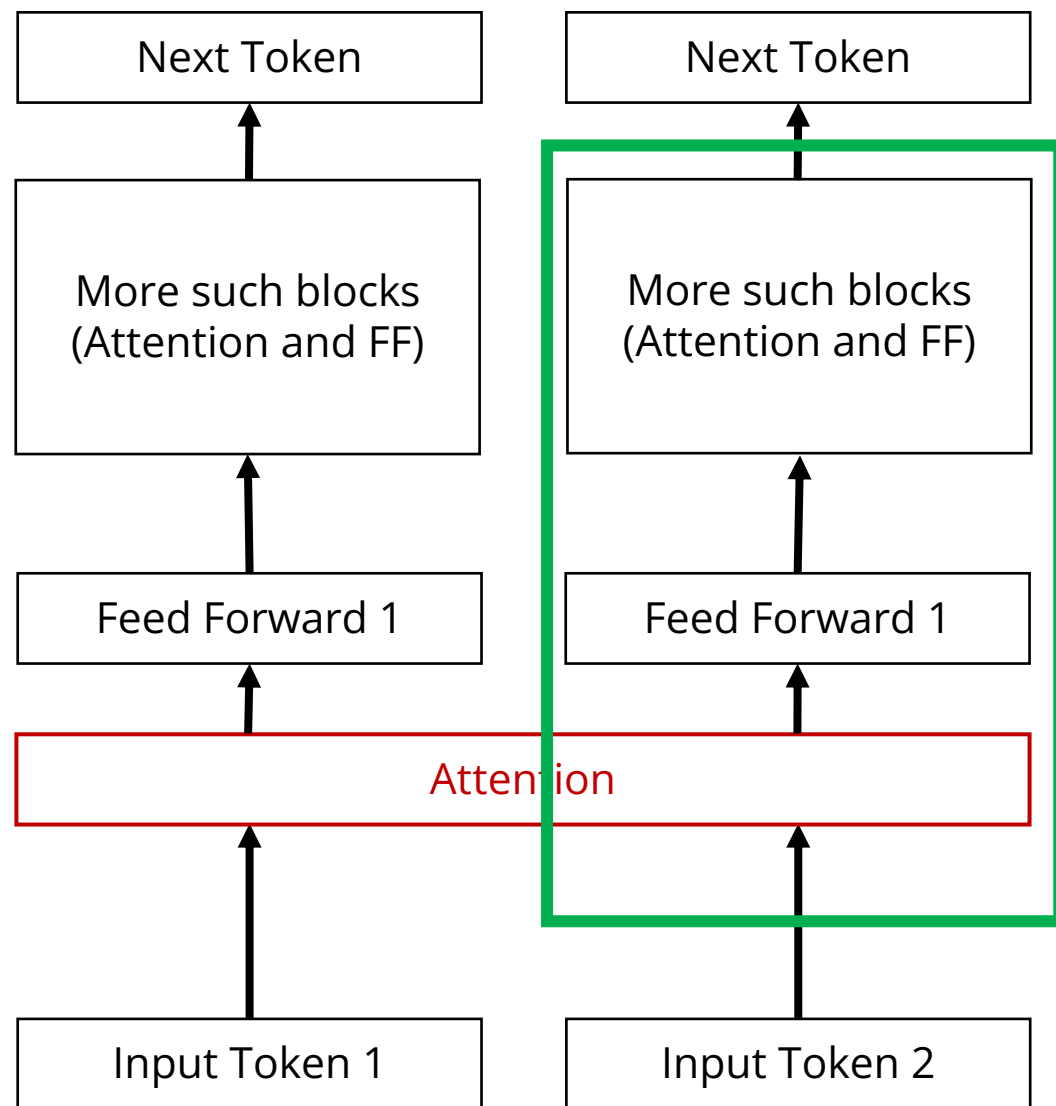
Outline

- ***Attention*** inside Transformer
- Static batching vs ***continuous batching***
- ***Selective batching*** for continuous batching

Decoding-only task

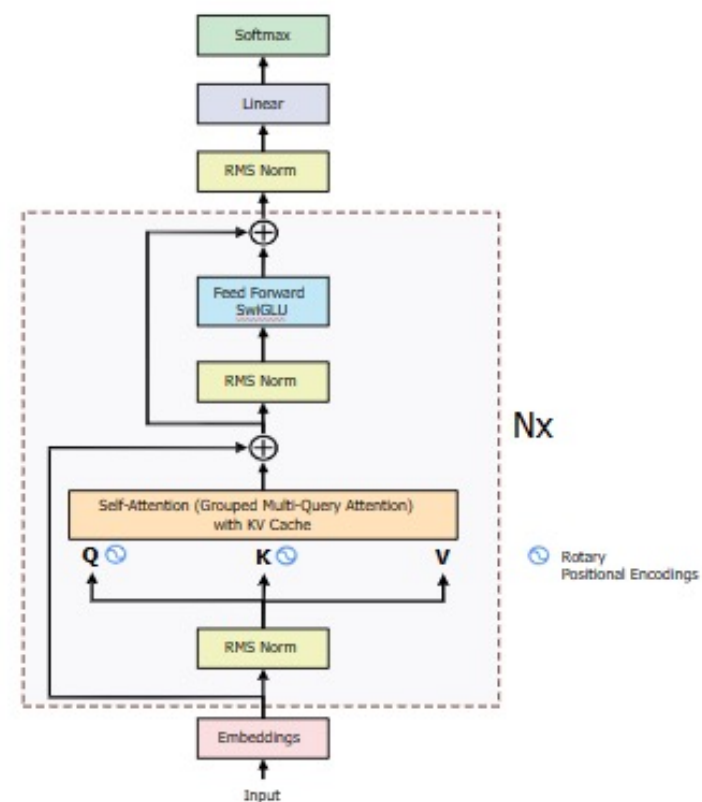


Attention: Captures dependency

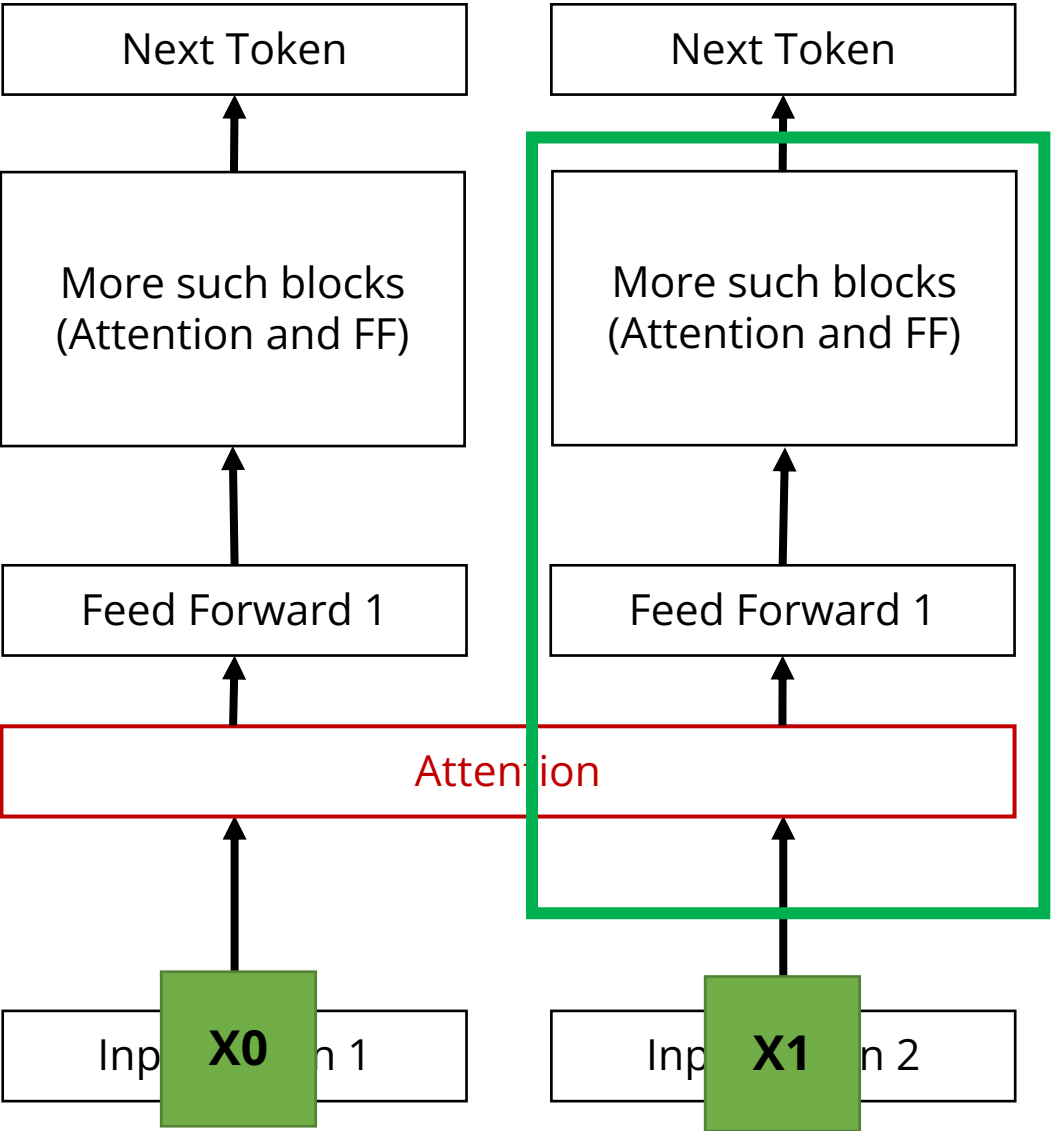


prediction block

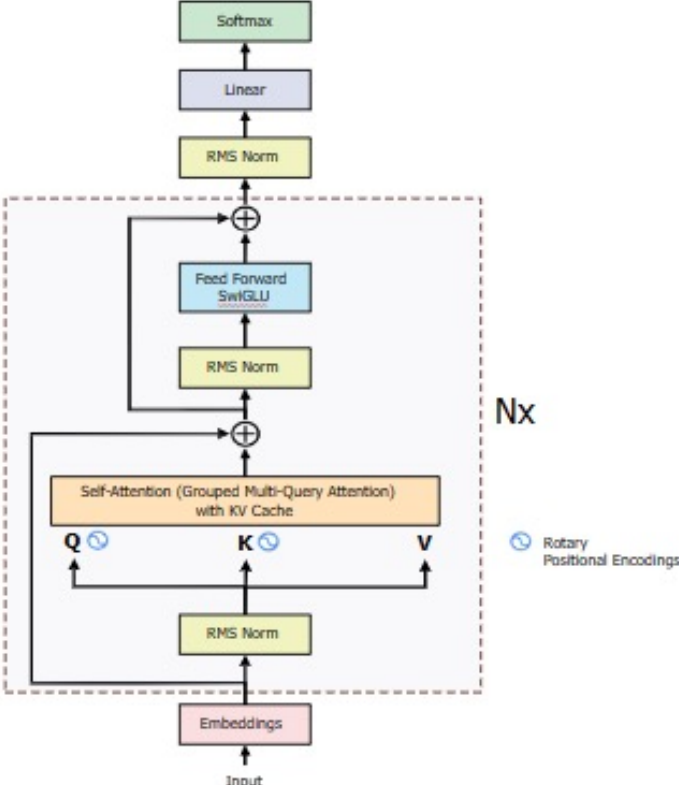
Mut-mul to capture relationship
(w/ key, query, value vectors)



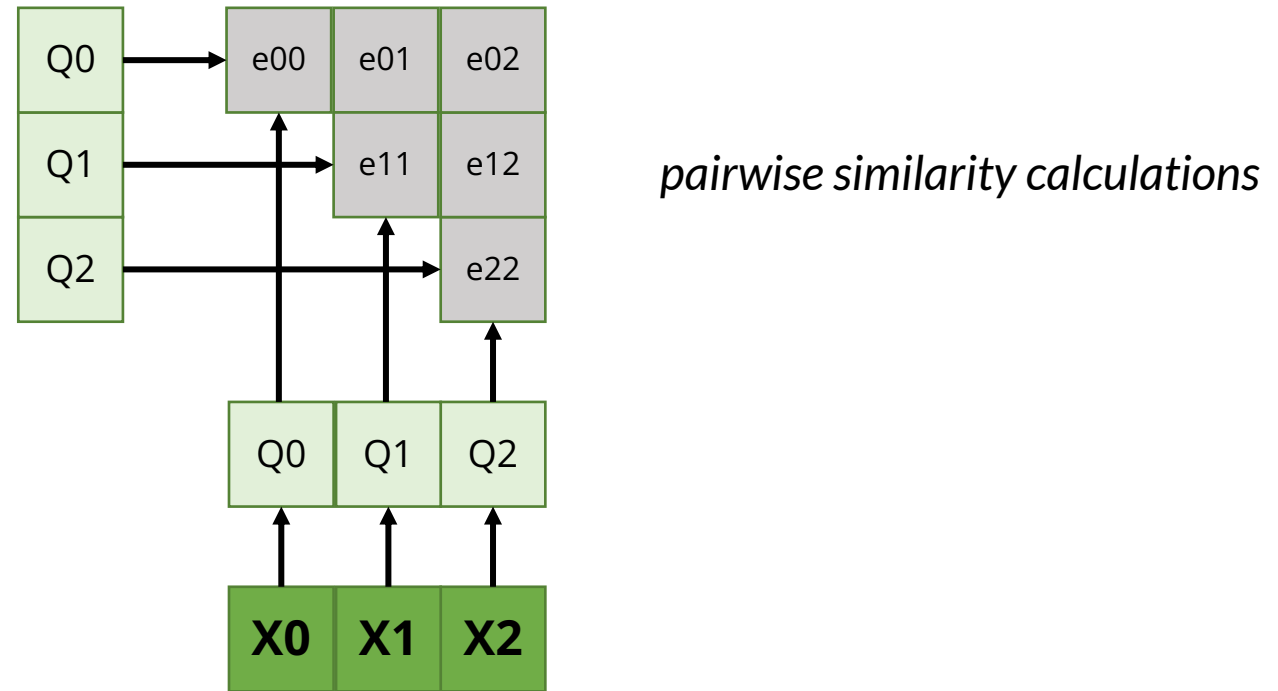
Attention: Captures dependency



Mut-mul to capture relationship
(w/ key, query, value vectors)



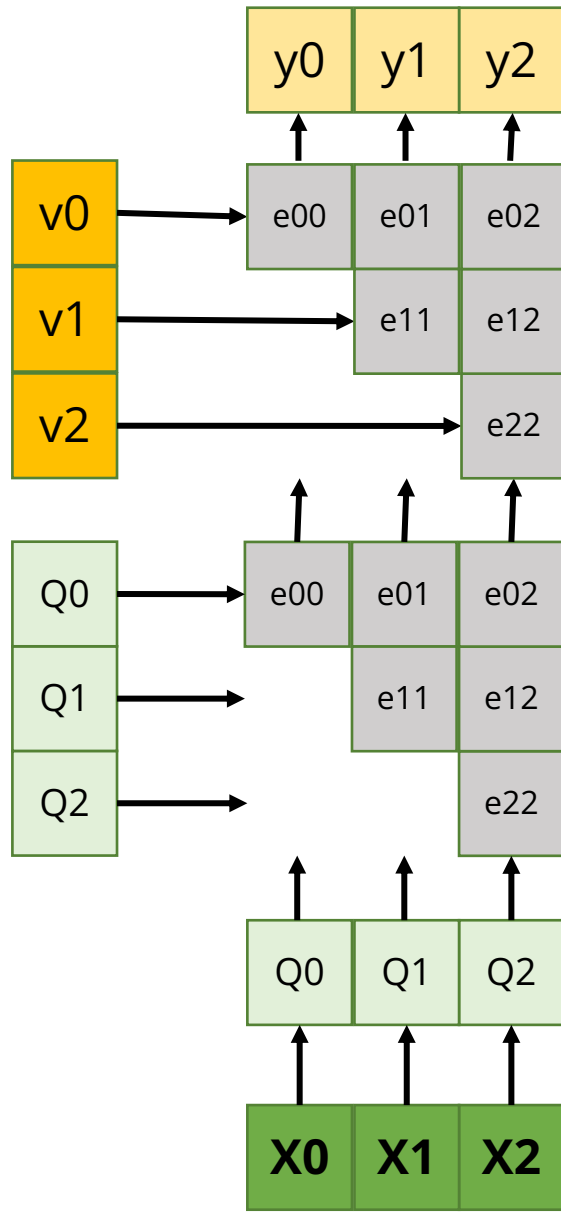
Naïve Approach 1: similarity vector



We can calculate similarities in this way, but

How can we express the meaning of each token?

Naïve Approach 2: similarity vector + **value** vector



multiply respective values

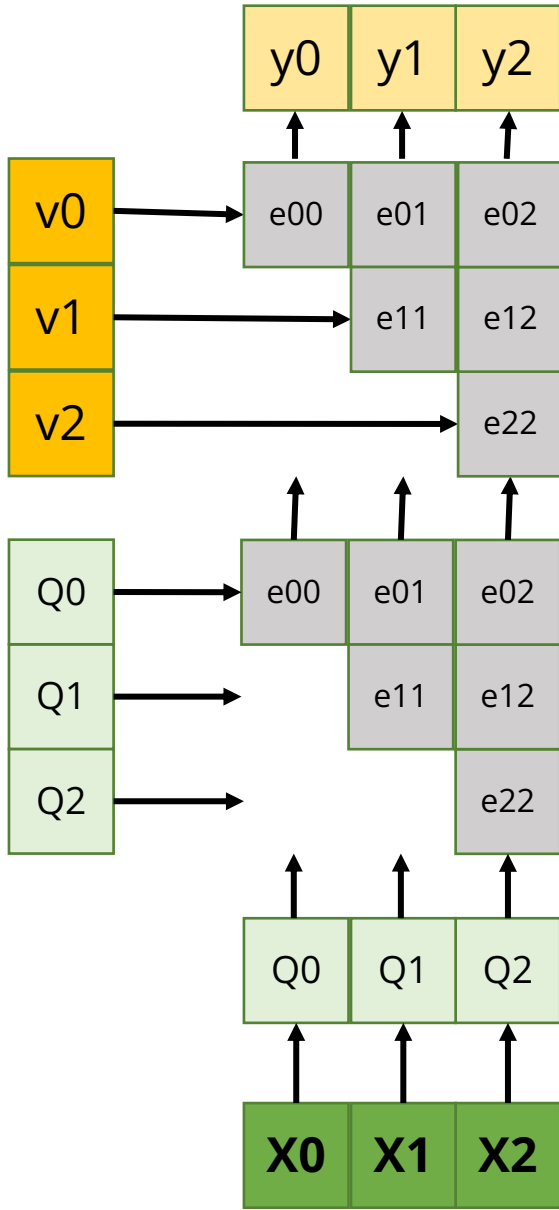
$$y_0 = v_0 * e_{00}$$

$$y_1 = v_0 * e_{01} + v_1 * e_{11}$$

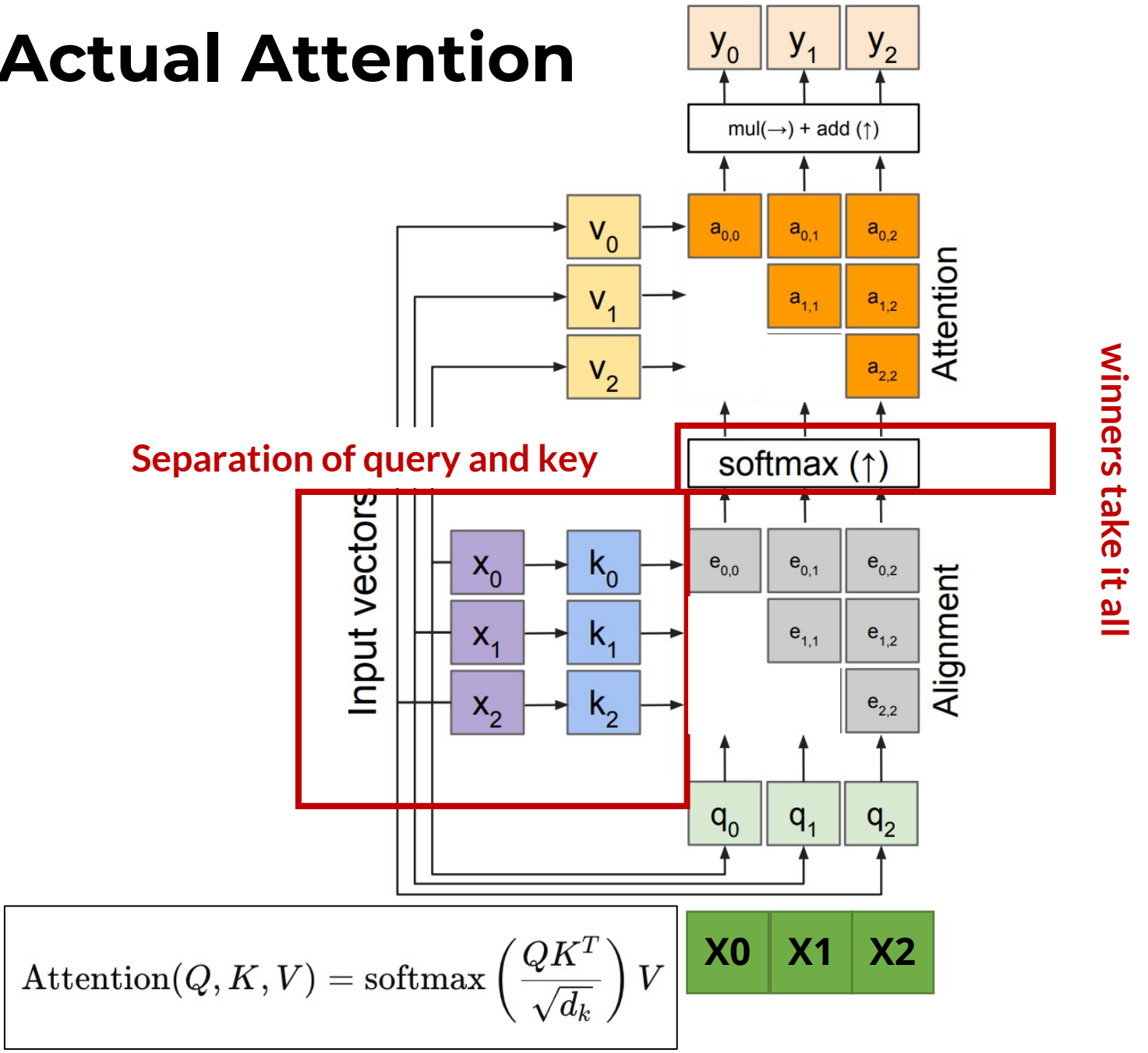
$$y_2 = v_0 * e_{02} + v_1 * e_{12} + v_2 * e_{22}$$

pairwise similarity calculations

Our naïve approach



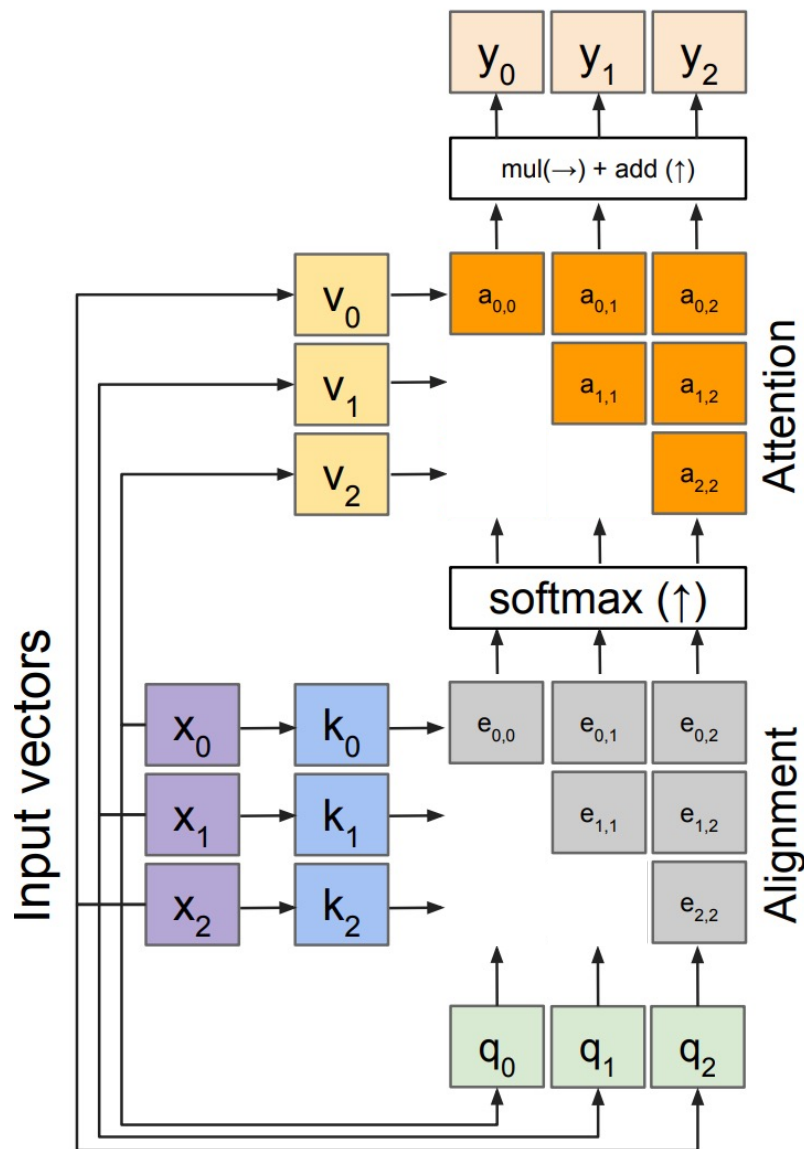
Actual Attention



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

What are the dimensions?

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



For y_0

$$Q = [1 \ H] \quad K = [1 \ H] \quad V = [1 \ H]$$

For y_1

$$Q = [2 \ H] \quad K = [2 \ H] \quad V = [2 \ H]$$

For y_N

$$Q = [N \ H] \quad K = [N \ H] \quad V = [N \ H]$$

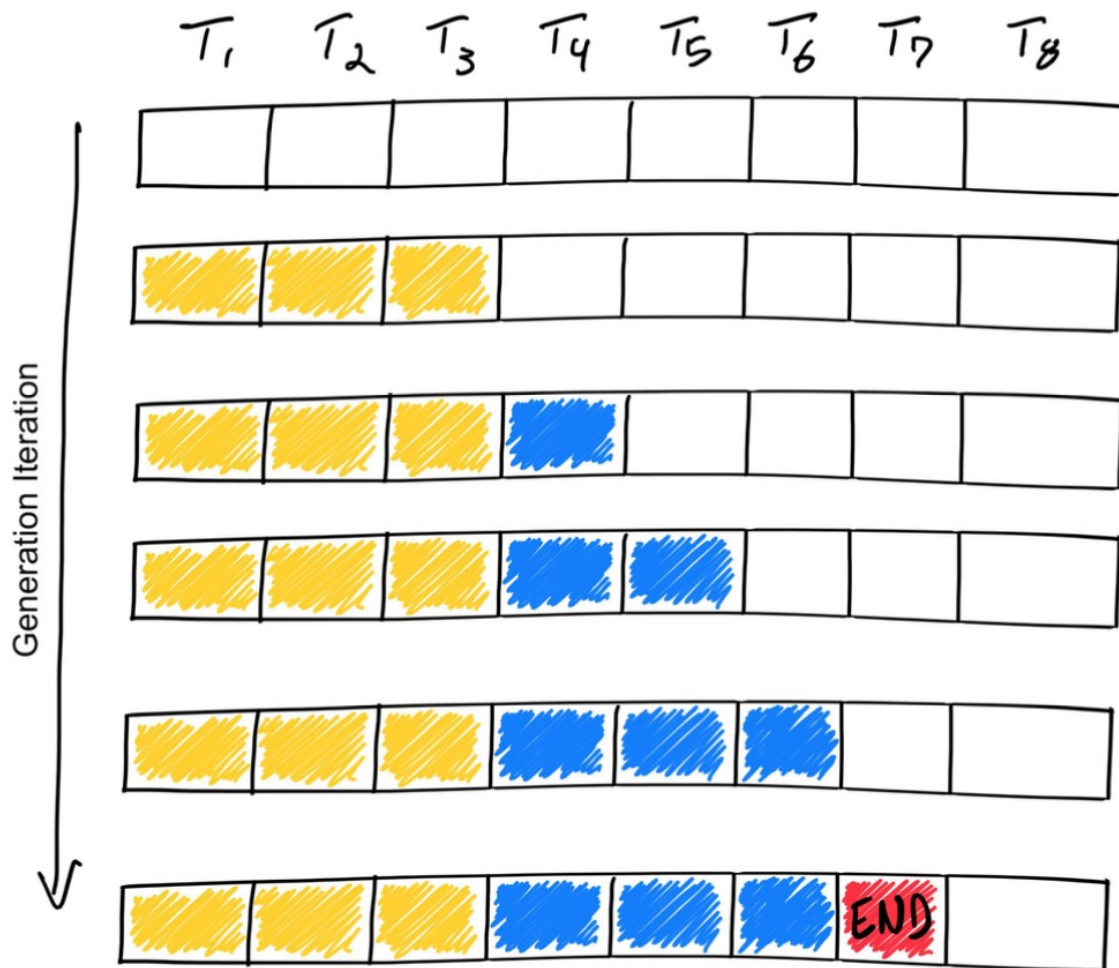
Orca

slide credits:

<https://cseweb.ucsd.edu/~yiying/cse291-winter24/reading/llm-serving.pdf>

https://cs231n.stanford.edu/slides/2022/lecture_11_ruohan.pdf

LLM inference basics



How does text generation work?

Iterative: each forward pass generates a single token

Autoregressive: generation consumes prompt tokens + previously generated tokens

Completion potentially decided by model: A generated token can be the end-of-sequence token

Legend:

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

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

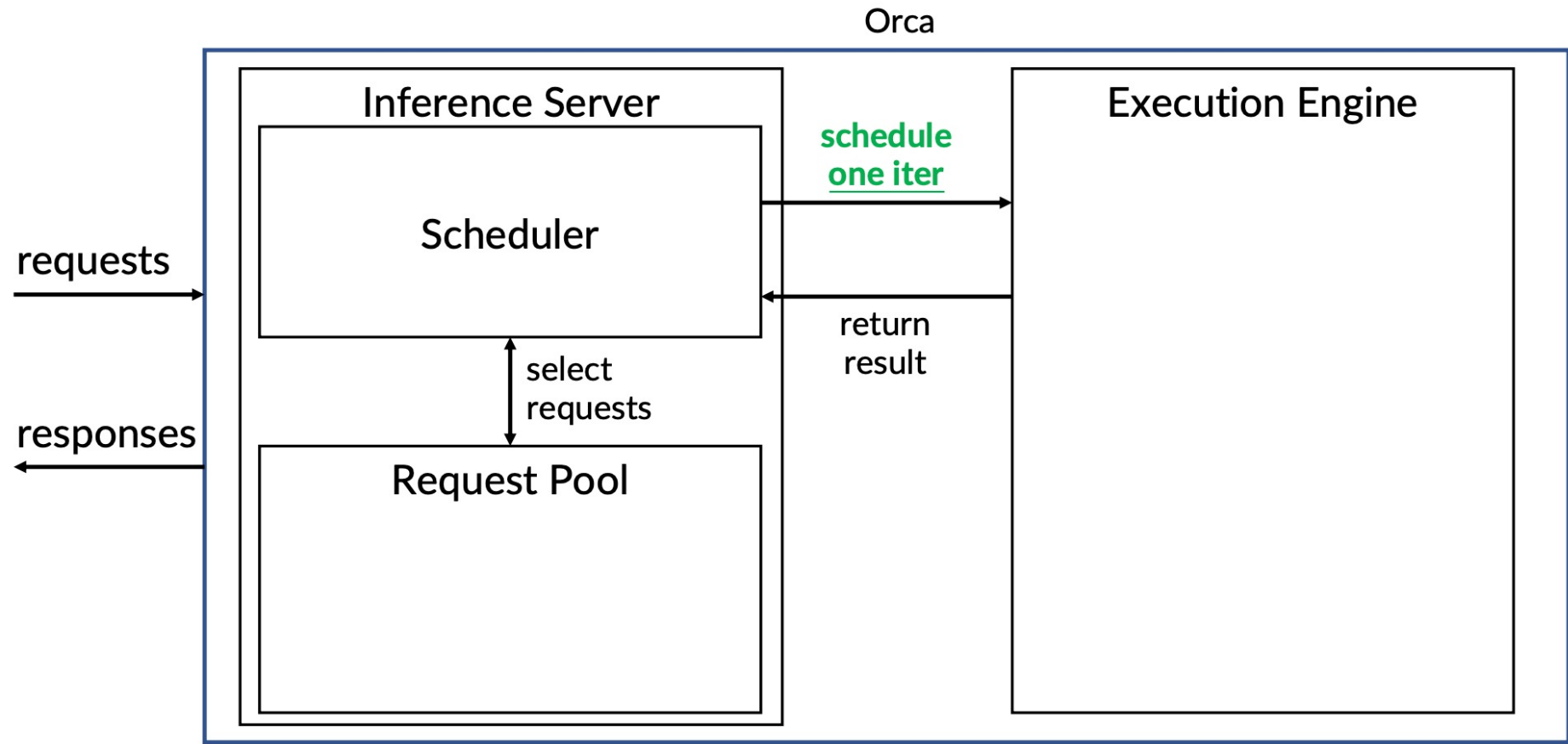
T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1				
S_2	S_2	S_2					
S_3	S_3	S_3					
S_4	S_4	S_4					

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					

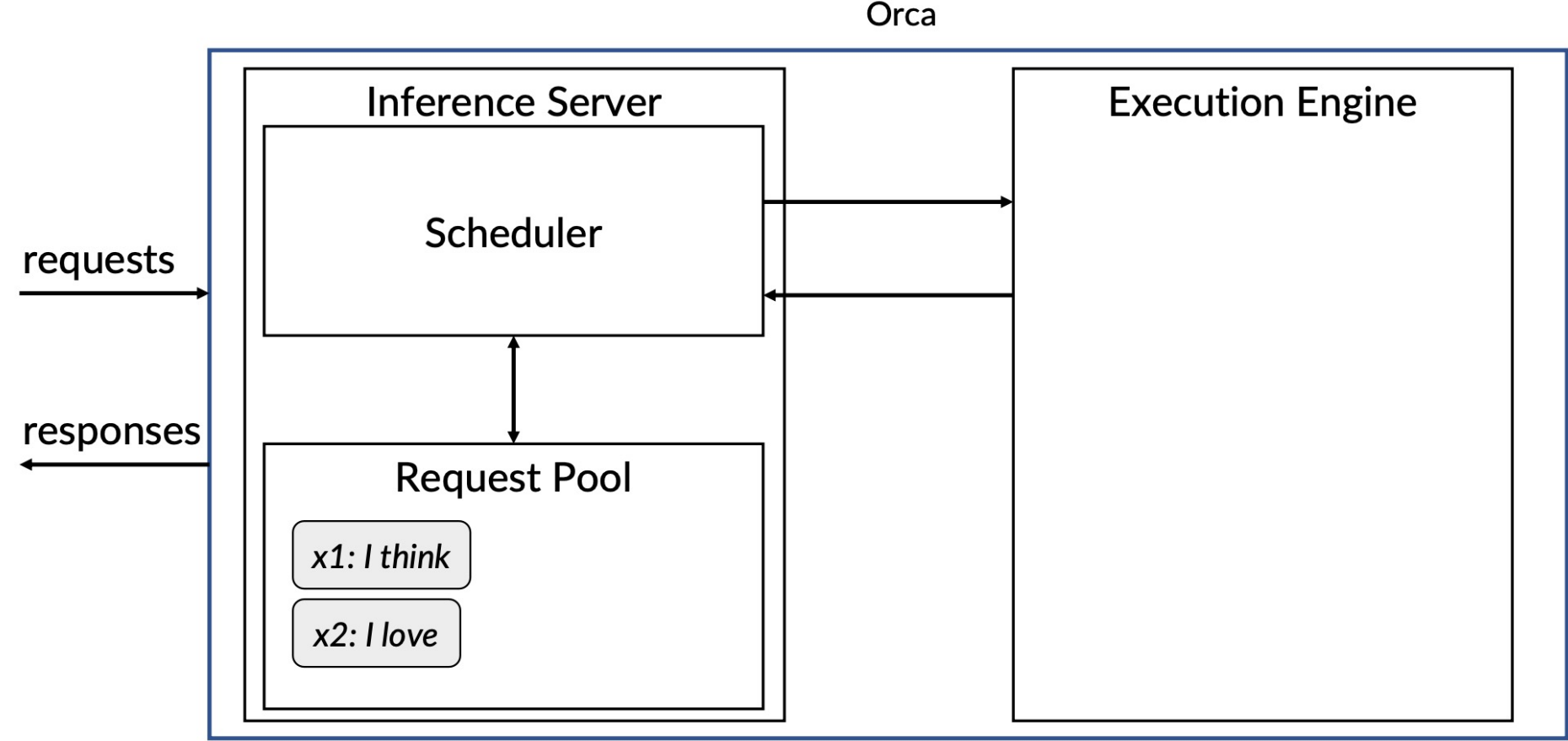
T_1	T_2	T_3	T_4	T_5	T_6	T_7	T_8
S_1	S_1	S_1	S_1	S_1	END	S_6	S_6
S_2	S_2	S_2	S_2	S_2	S_2	S_2	END
S_3	S_3	S_3	S_3	END	S_5	S_5	S_5
S_4	S_4	S_4	S_4	S_4	S_4	END	S_7

Iteration-level scheduling



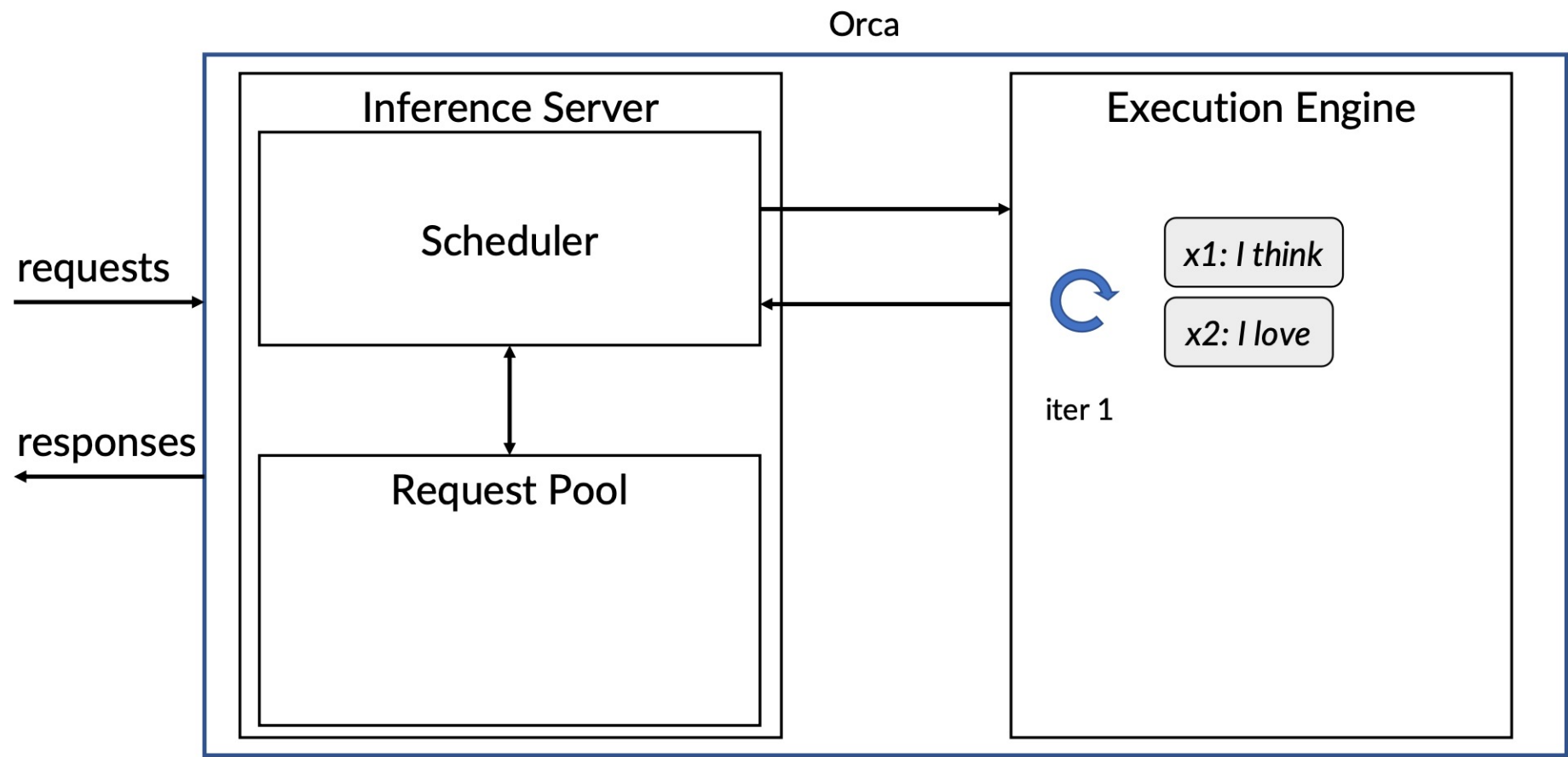
* maximum batch size = 3

Iteration-level scheduling



* maximum batch size = 3

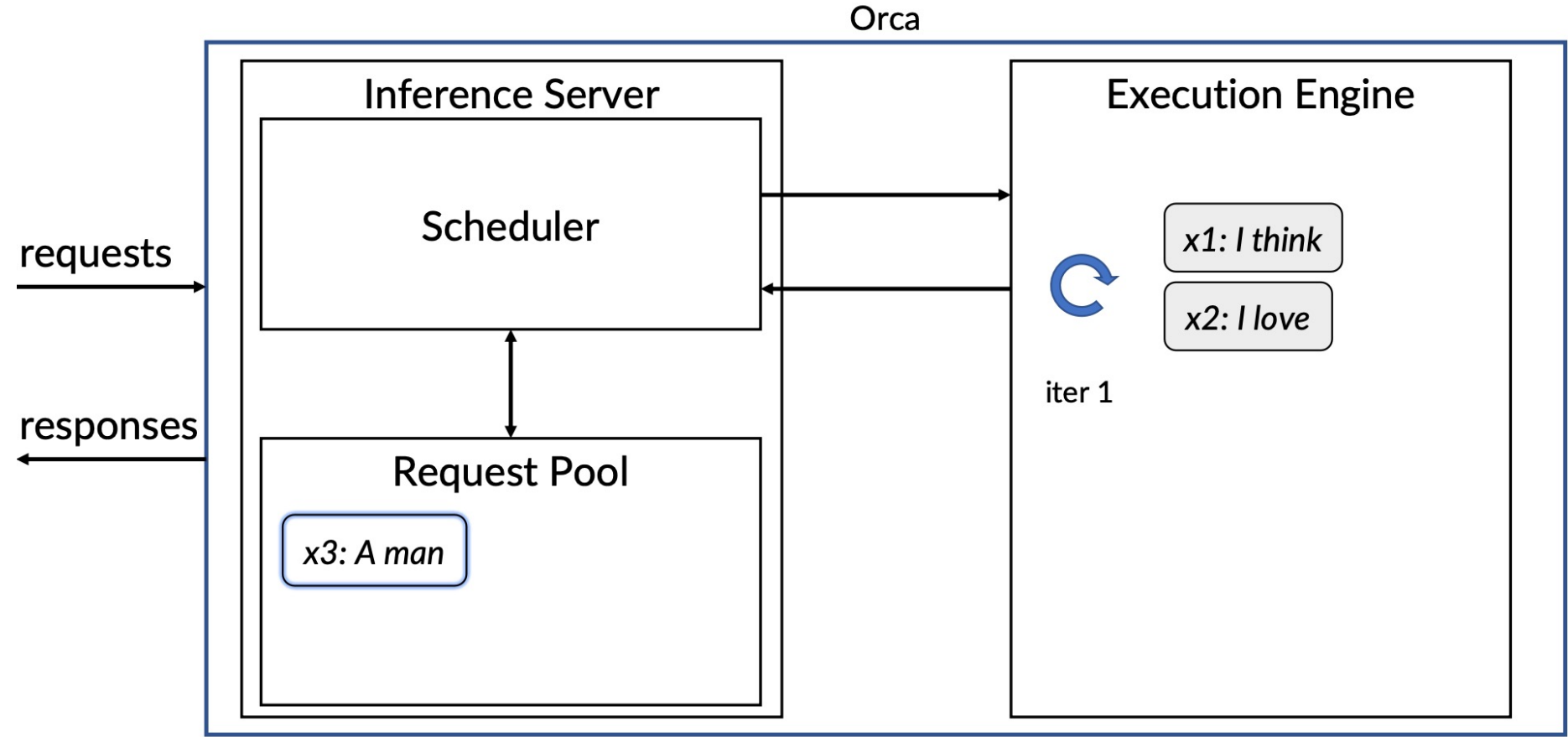
Iteration-level scheduling



* maximum batch size = 3

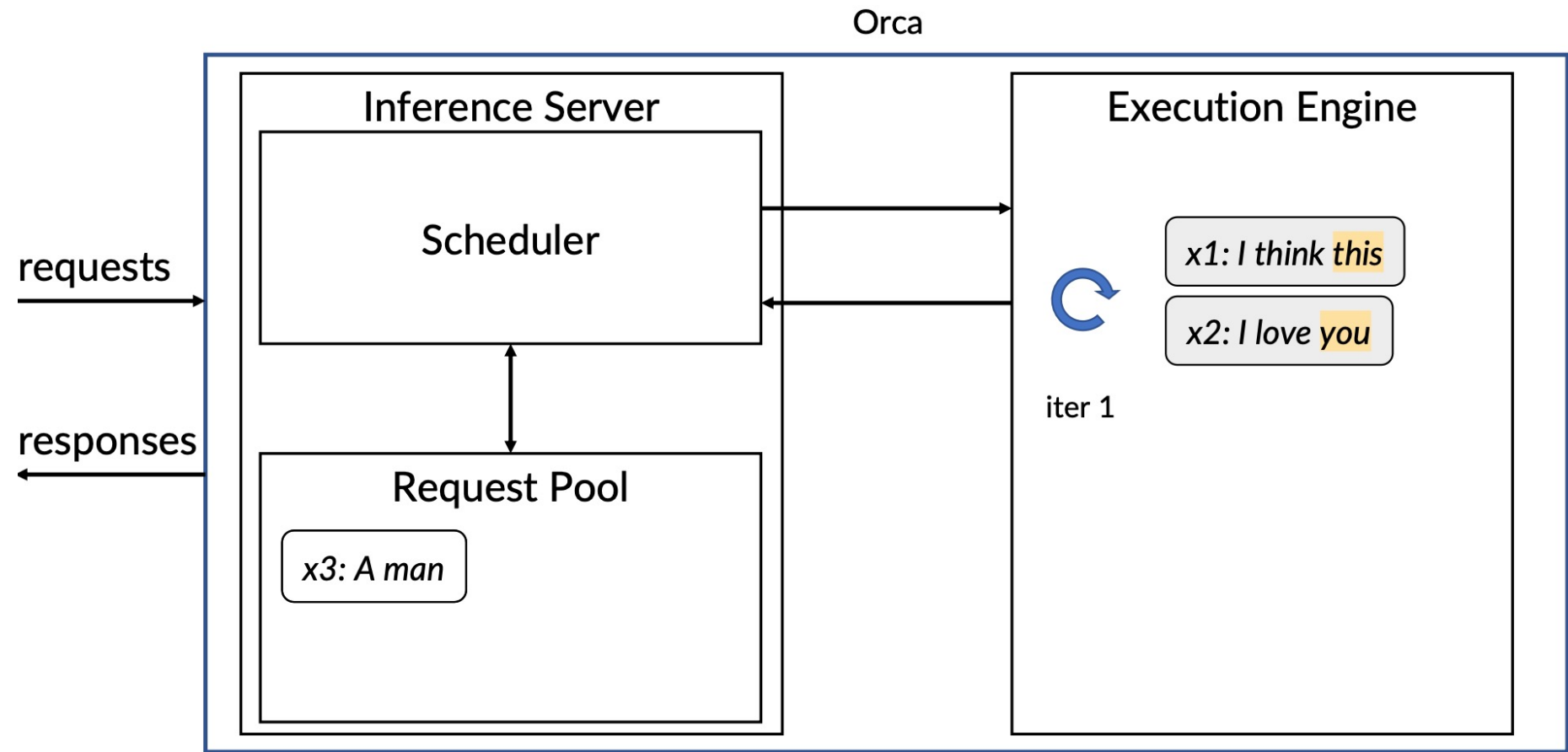
Source: Orca OSDI 2020 Talk

Iteration-level scheduling



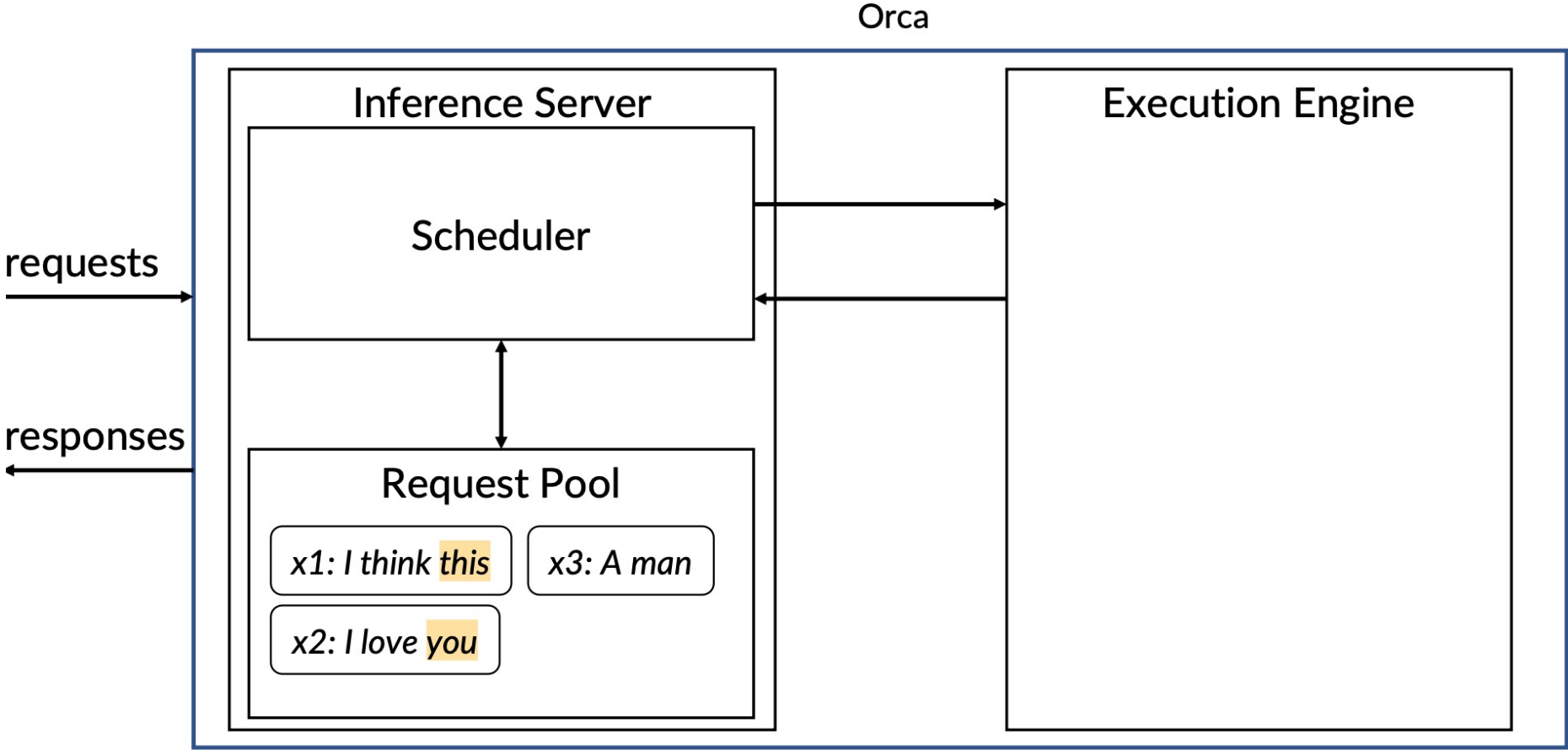
* maximum batch size = 3

Iteration-level scheduling



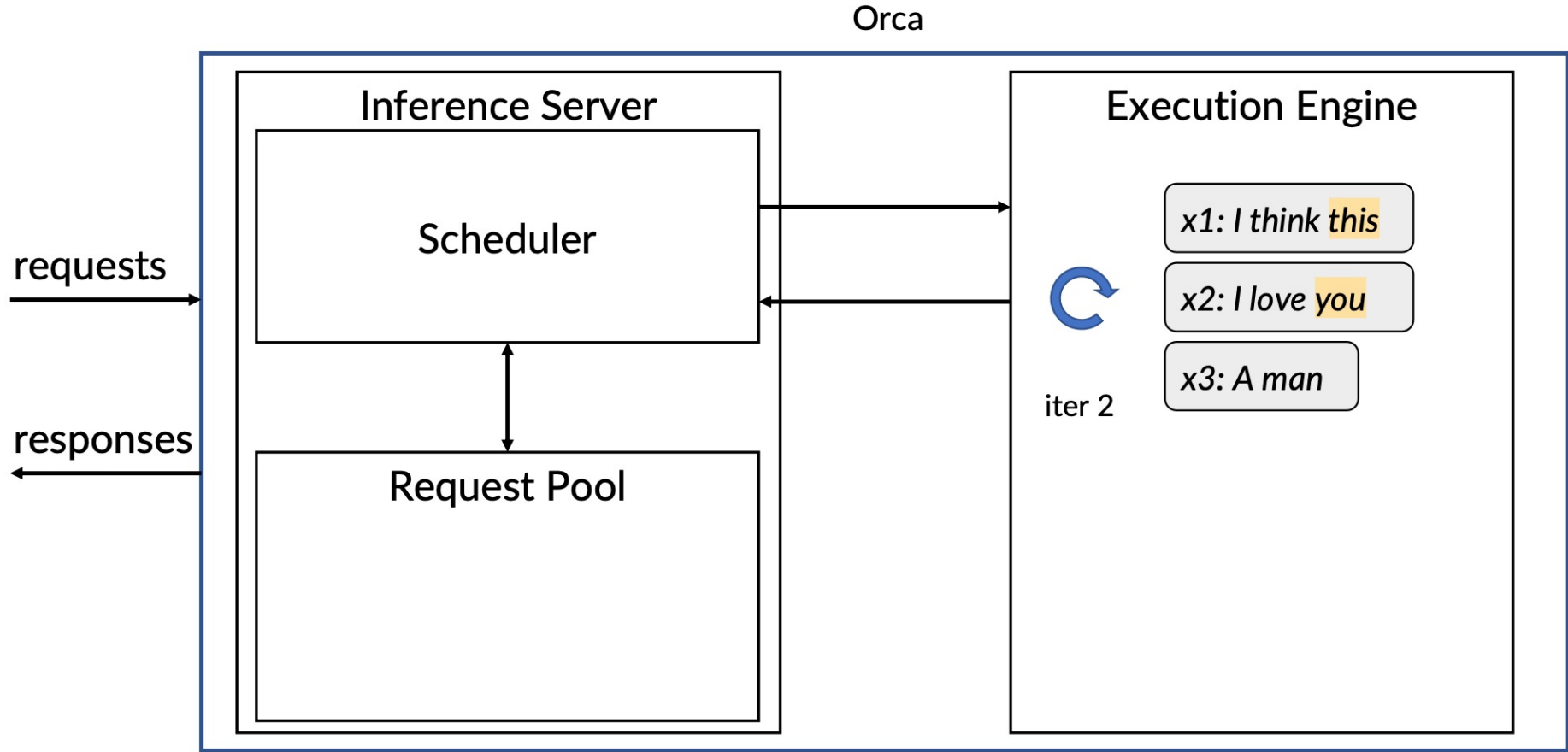
* maximum batch size = 3

Iteration-level scheduling



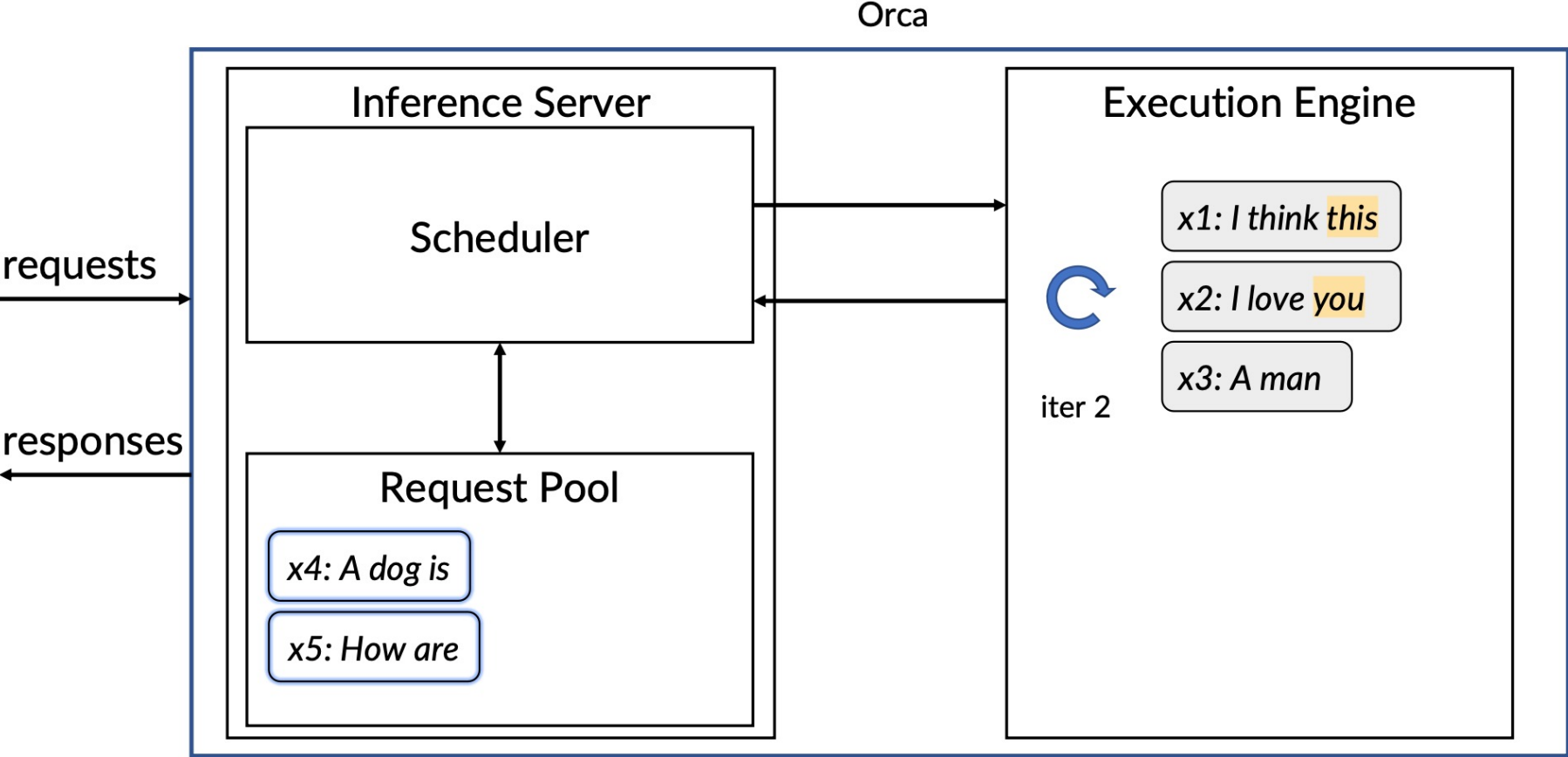
* maximum batch size = 3

Iteration-level scheduling



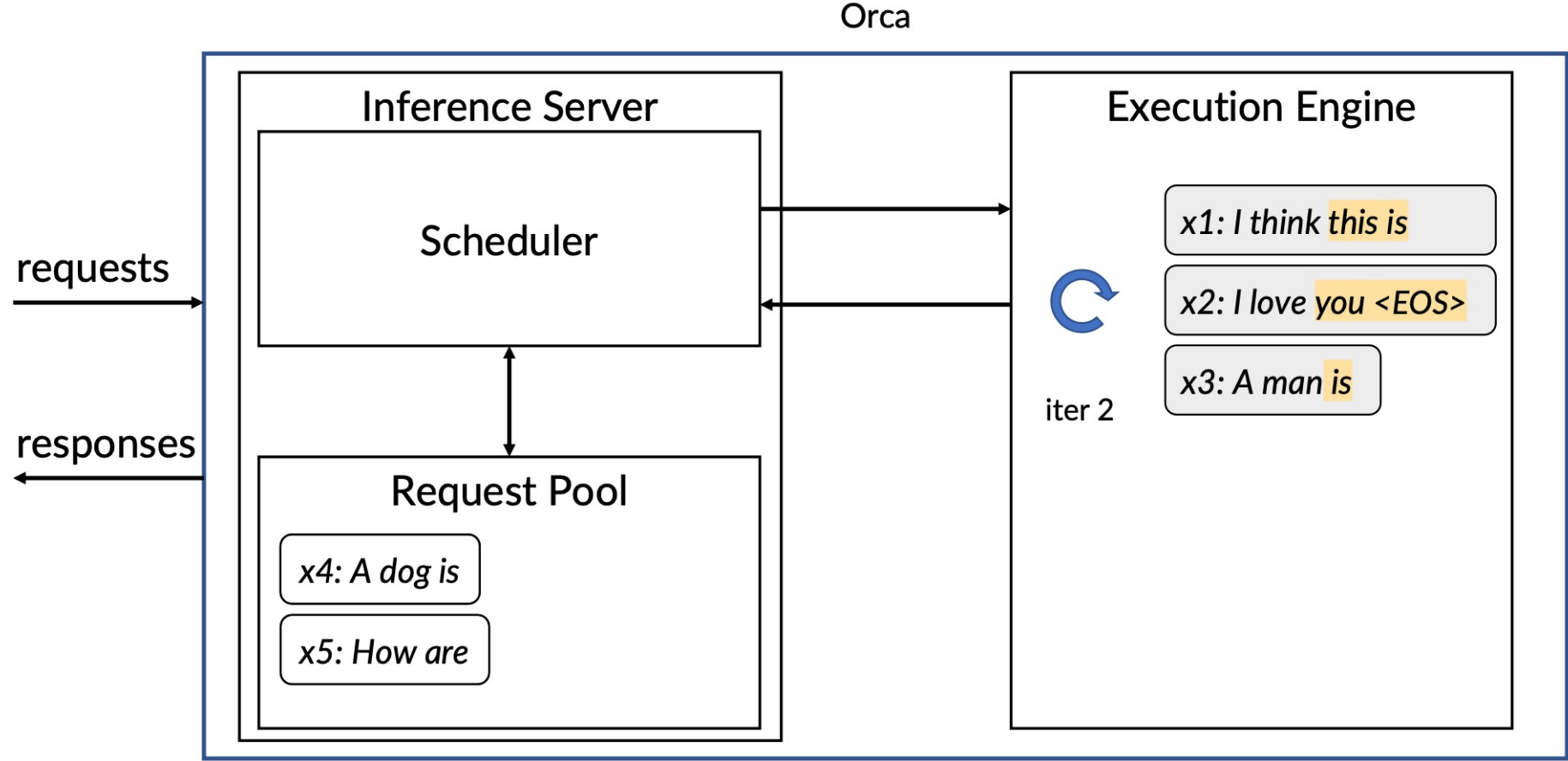
* maximum batch size = 3

Iteration-level scheduling



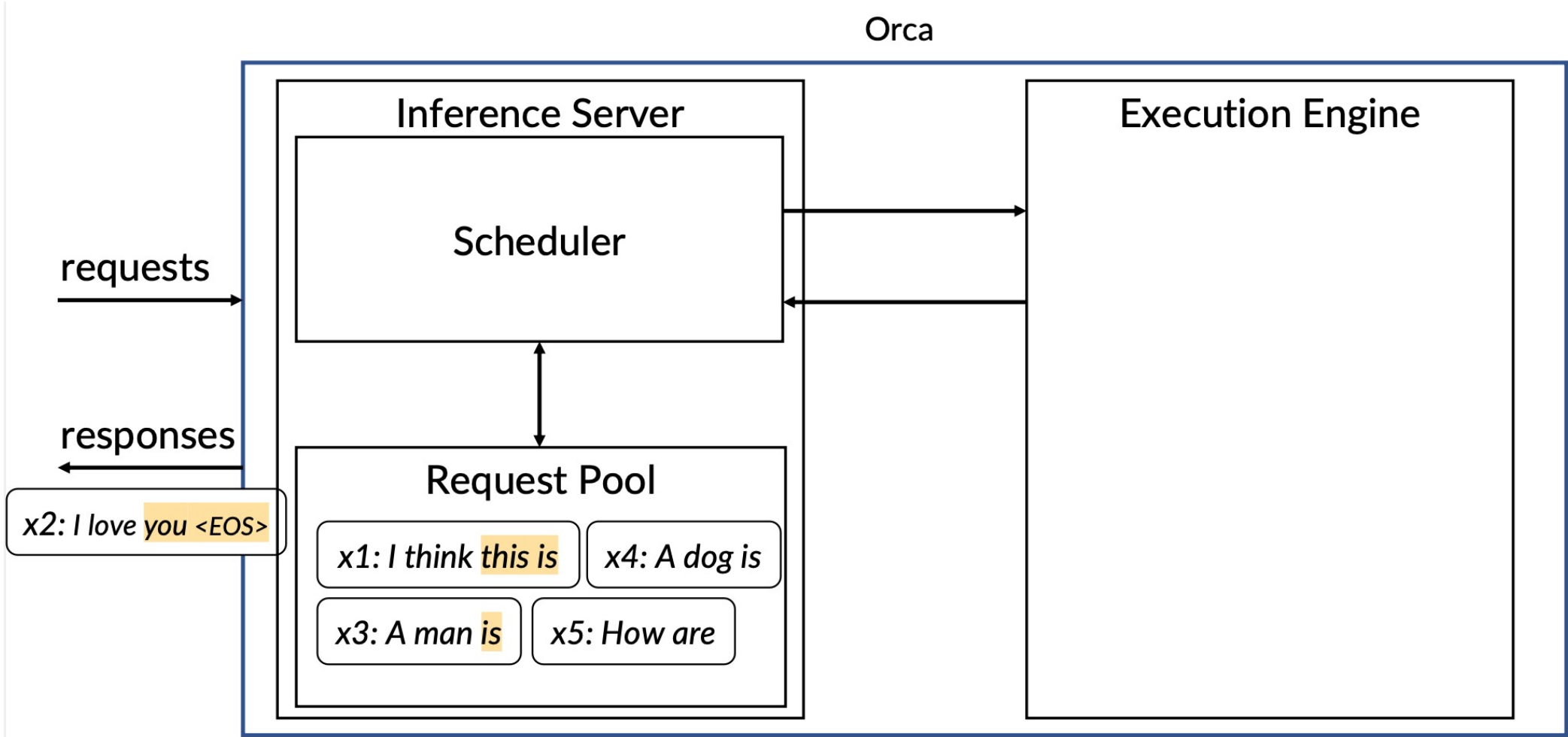
* maximum batch size = 3

Iteration-level scheduling



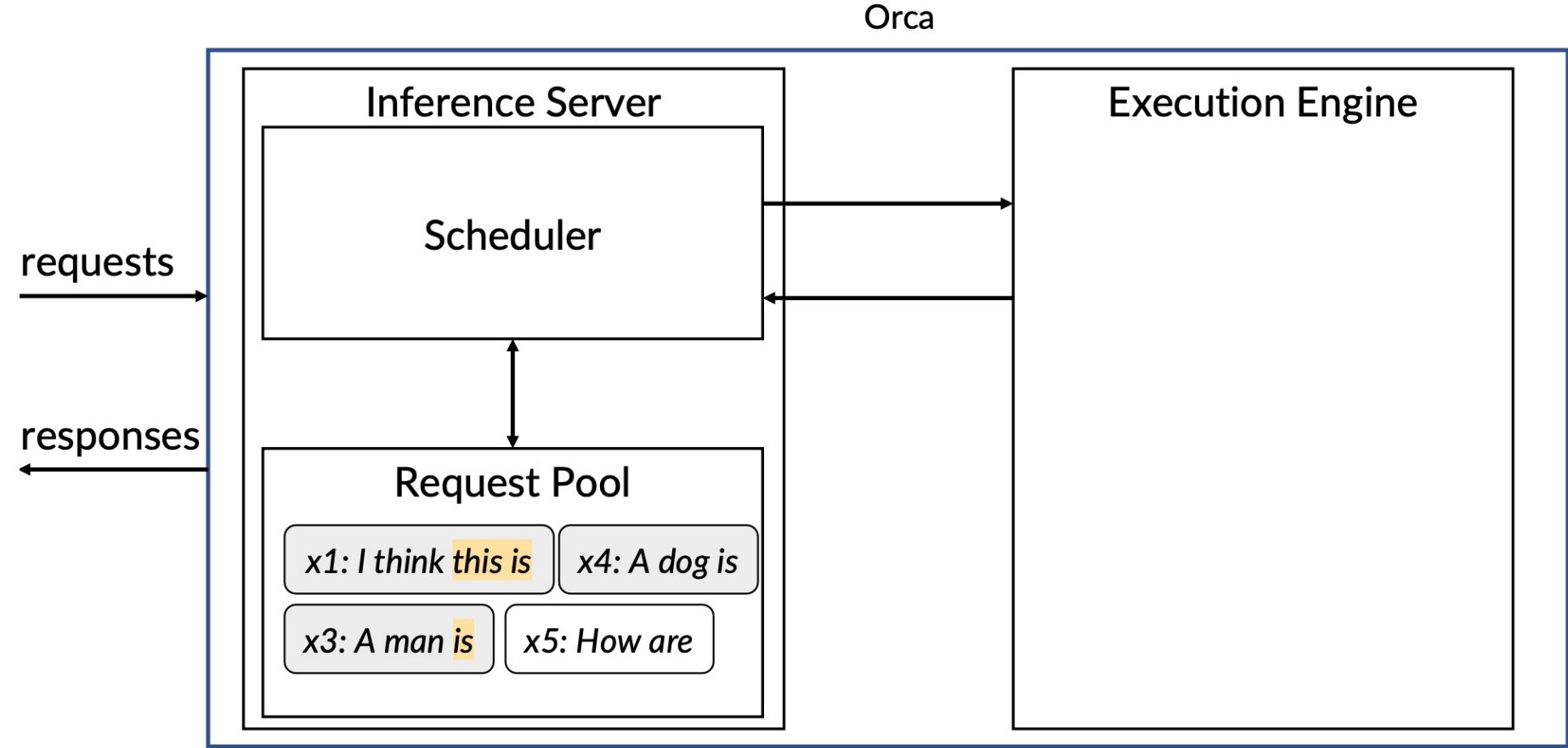
* maximum batch size = 3

Iteration-level scheduling



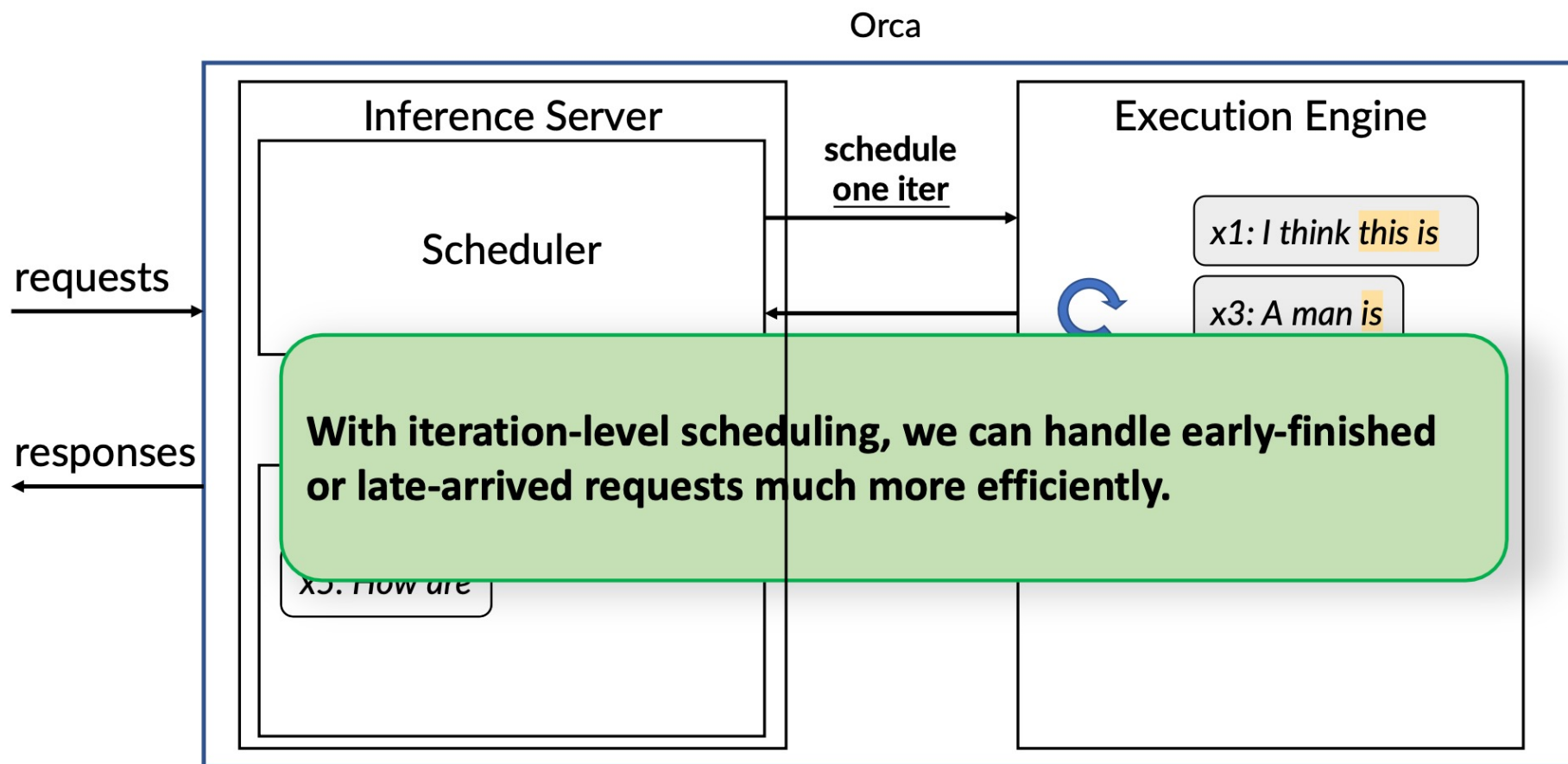
* maximum batch size = 3

Iteration-level scheduling



* maximum batch size = 3

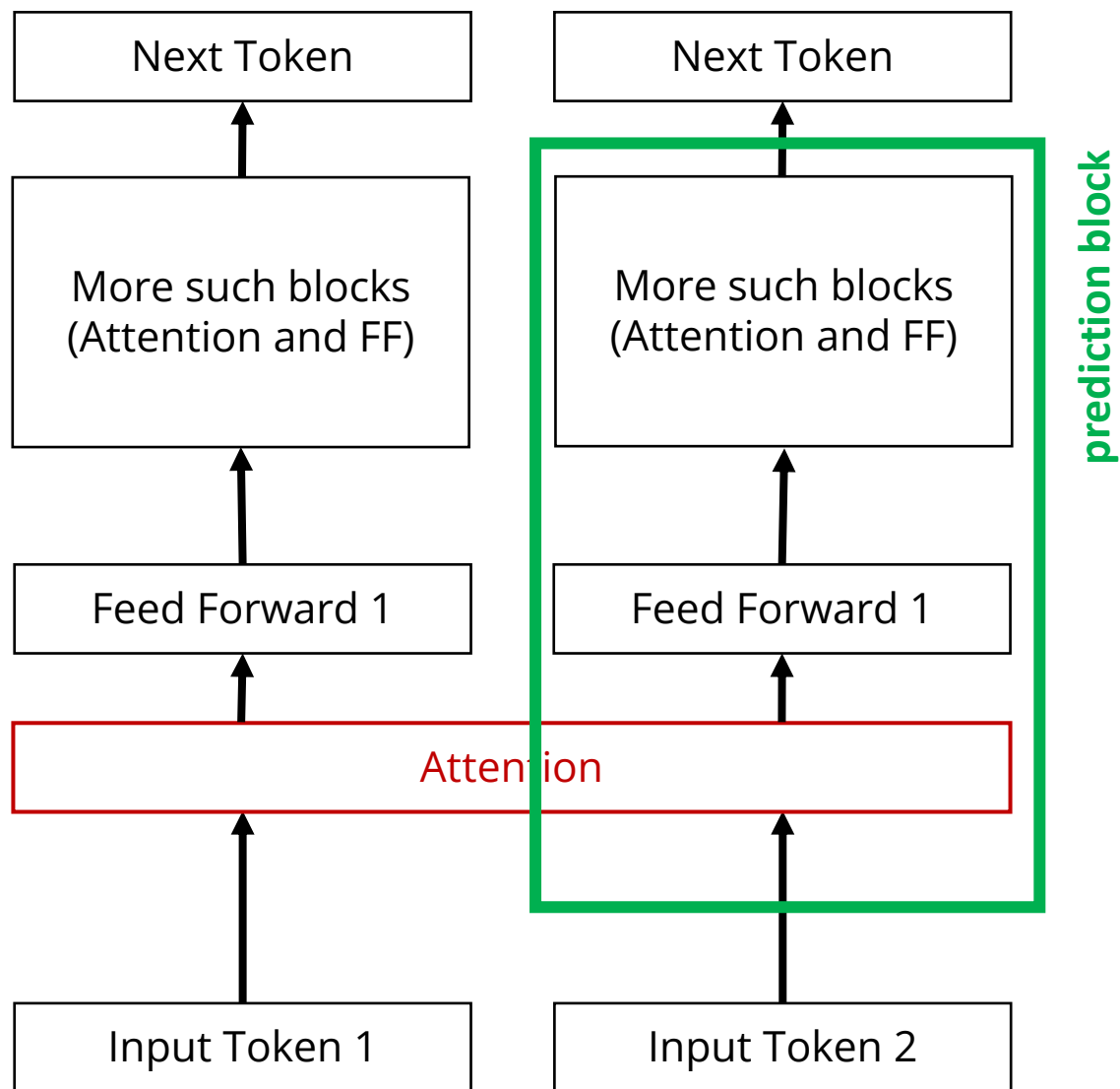
Iteration-level scheduling



* maximum batch size = 3

source: Orca OSDI 2022 Talk

Feed-forward is independent, given attention



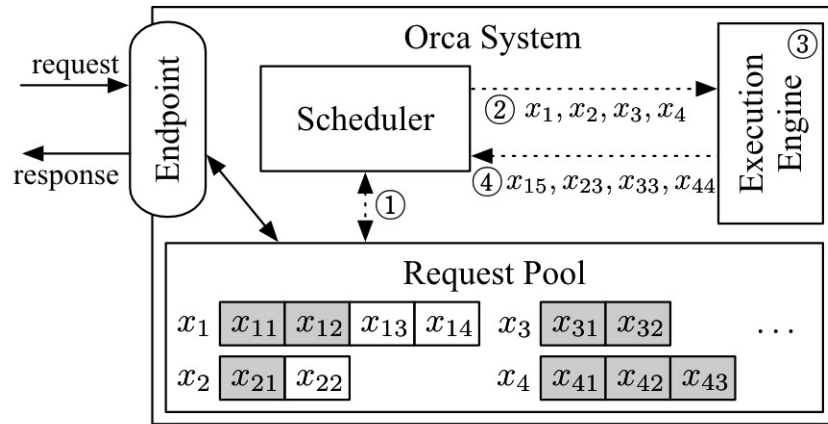
$$y_1 = f(Wx_1 + b), \quad y_2 = f(Wx_2 + b)$$

$$\hat{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^{2n}$$

$$\hat{W} = \begin{bmatrix} W & W \end{bmatrix} \in \mathbb{R}^{m \times 2n}$$

$$\hat{y} = \hat{W}\hat{x} + \hat{b} = Wx_1 + Wx_2 + 2b = y_1 + y_2$$

Incompatible scheduling: Attention sizes are different



- Two prefills of different lengths (x_3 and x_4)
- Two decoding at different indexes (x_1 and x_2)
- Prefill and decoding (x_1 and x_3)

Mismatch in attention sizes

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

Non-trivial to merge in matrices

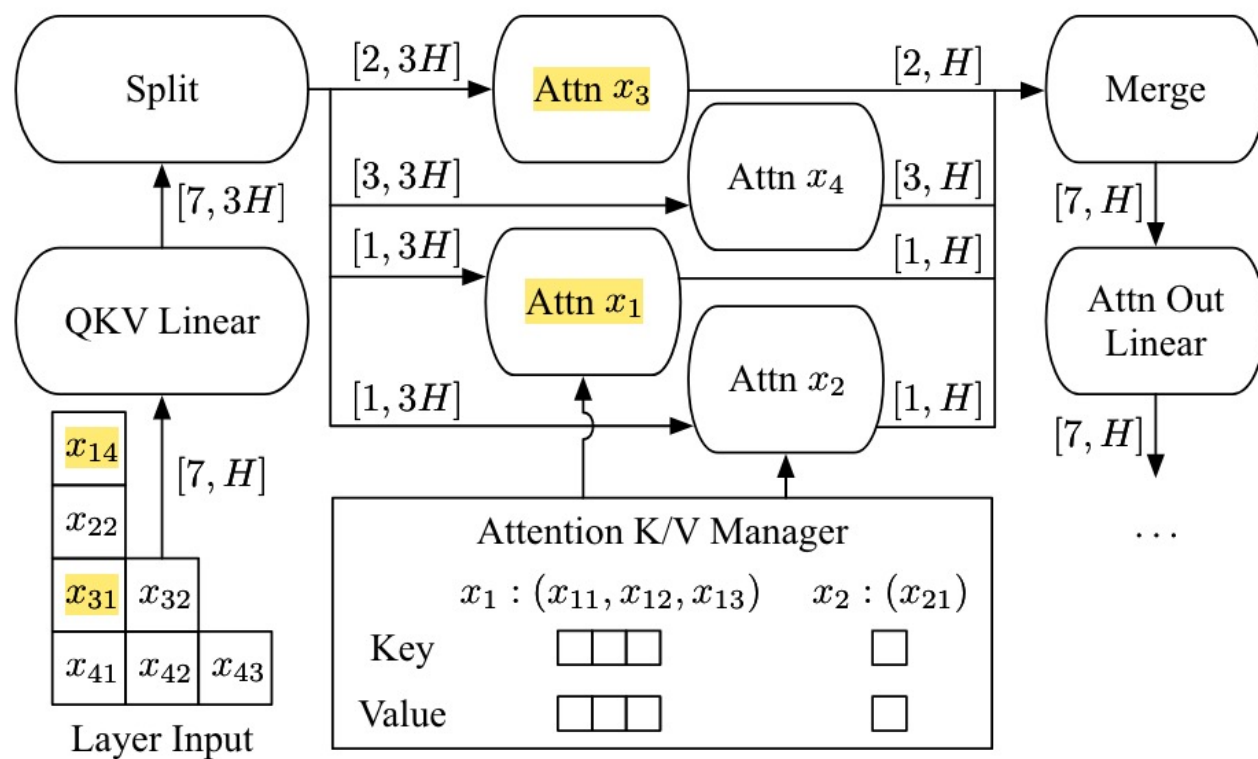
prompt 1:

$$Q = \begin{bmatrix} n1 * d \end{bmatrix} \quad K^T = \begin{bmatrix} d * n1 \end{bmatrix}$$

prompt 2:

$$Q = \begin{bmatrix} n2 * d \end{bmatrix} \quad K^T = \begin{bmatrix} d * n2 \end{bmatrix}$$

Selective batching: *separate* Attention computing



Summary

- **Attention** inside Transformer: is a fundamental mechanism for capturing semantic relatedness
- **Continuous batching:** *allows efficient GPU resource utilization*
- **Selective batching:** *merged feed-forward + separate attention computations*

Questions?