

Spring25 CS598YP

## 21.2: vLLM: PagedAttention

Yongjoo Park

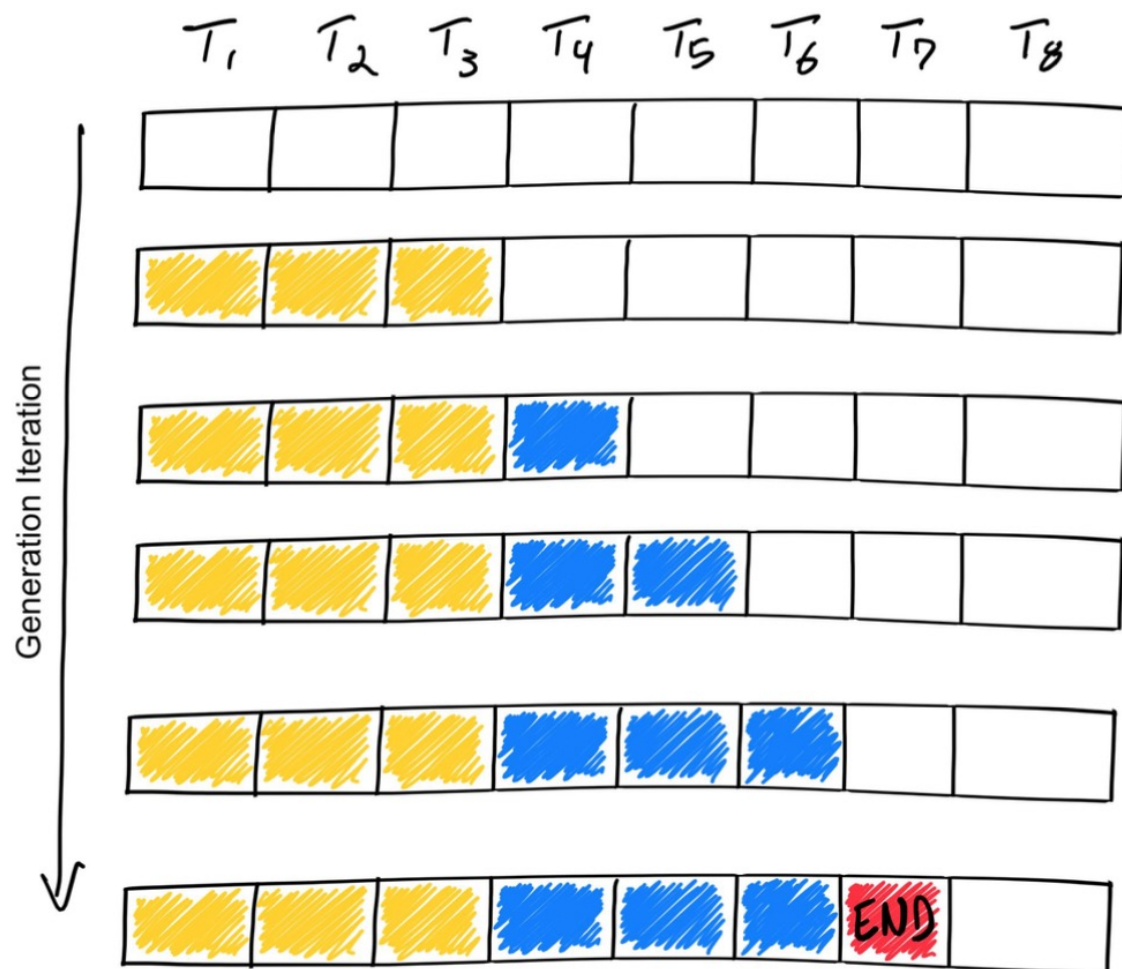
University of Illinois Urbana-Champaign

# Outline

- *Recap: Orca's continuous batching*
- *What is **KV cache**?*
- *vllm's PagedAttention*
- *Scheduling inside vllm*

Recap: Continuous batching

# 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

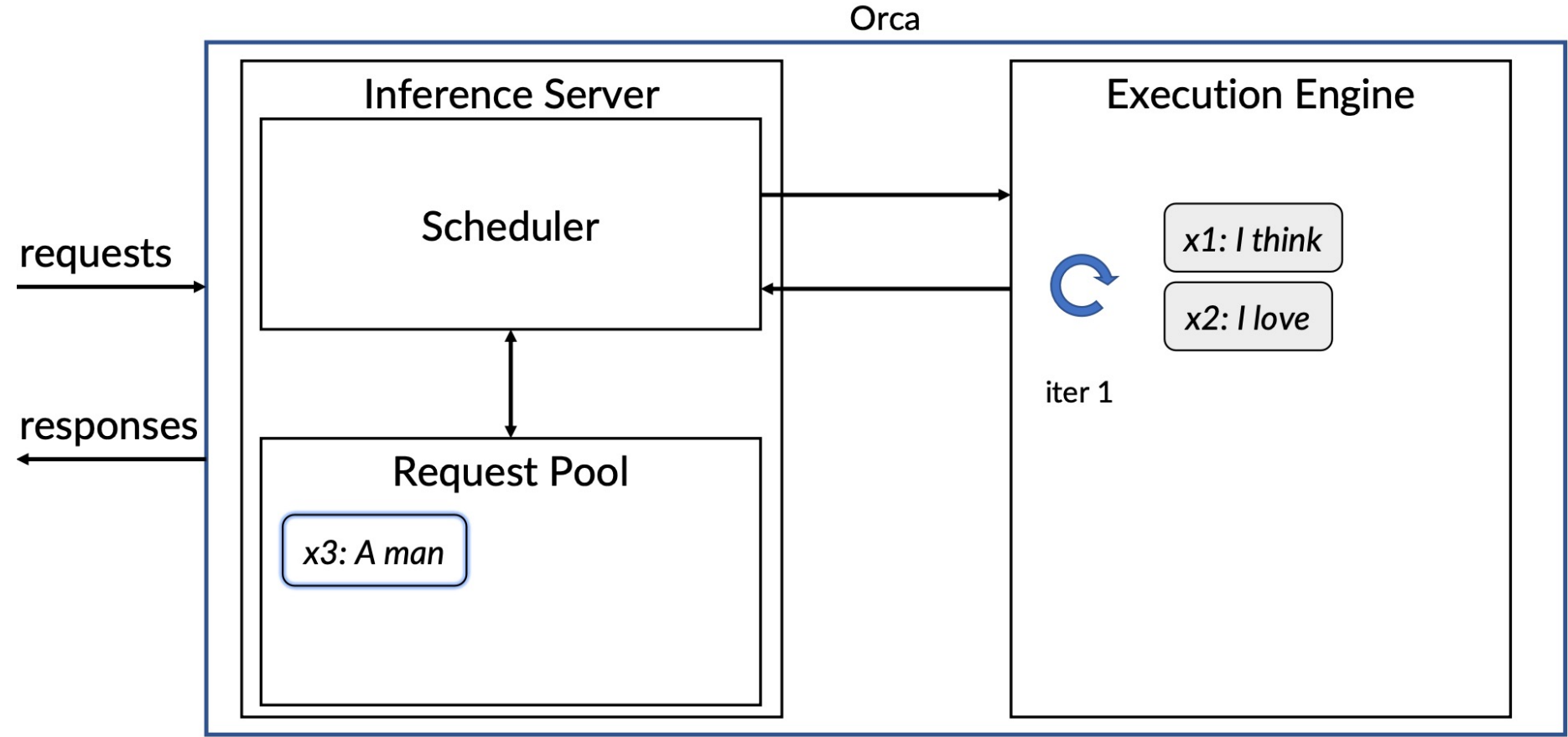
$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$				
$S_2$	$S_2$	$S_2$					
$S_3$	$S_3$	$S_3$					
$S_4$	$S_4$	$S_4$					

$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	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	

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

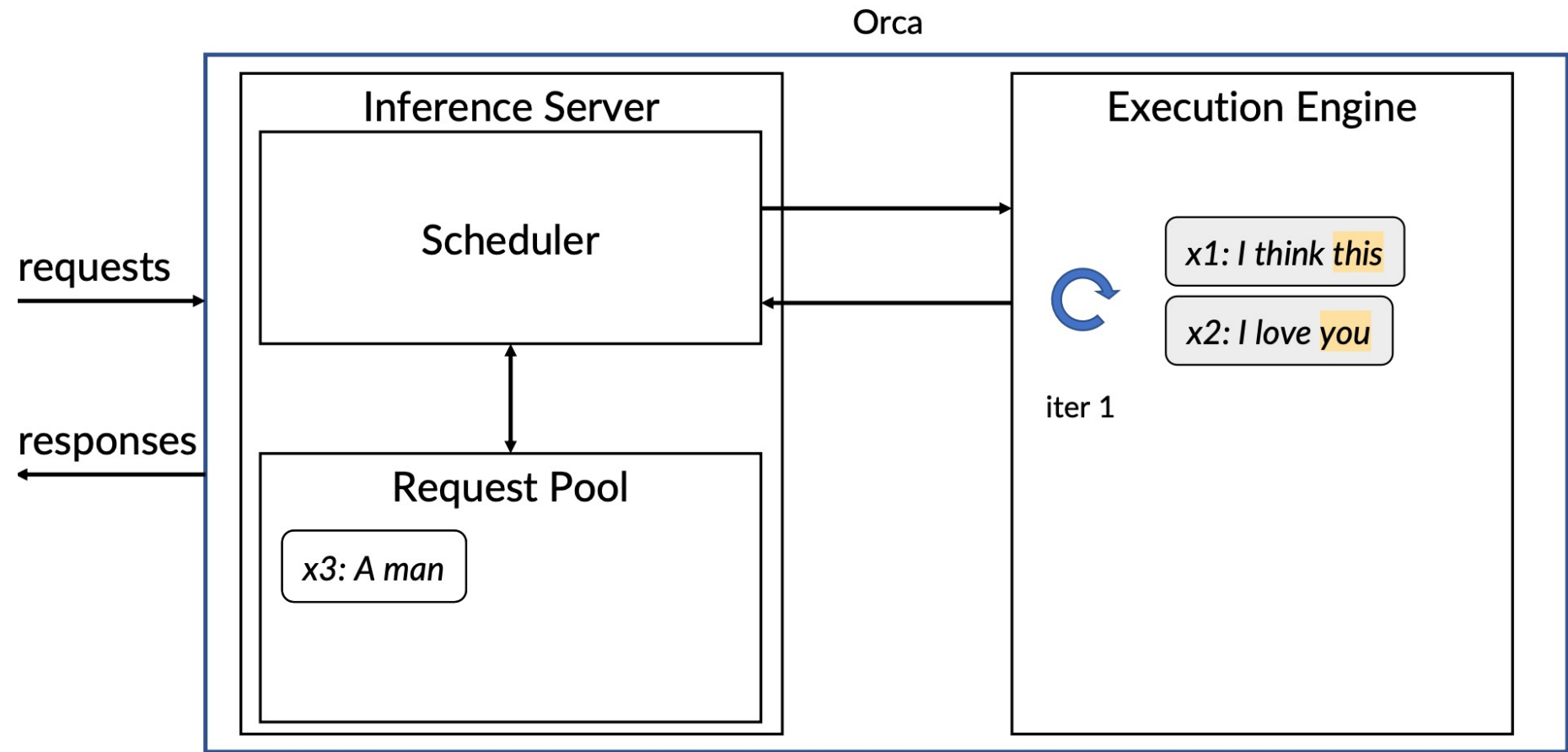
$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END	$S_6$	$S_6$
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END	$S_5$	$S_5$	$S_5$
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	$S_7$

# Iteration-level scheduling



\* maximum batch size = 3

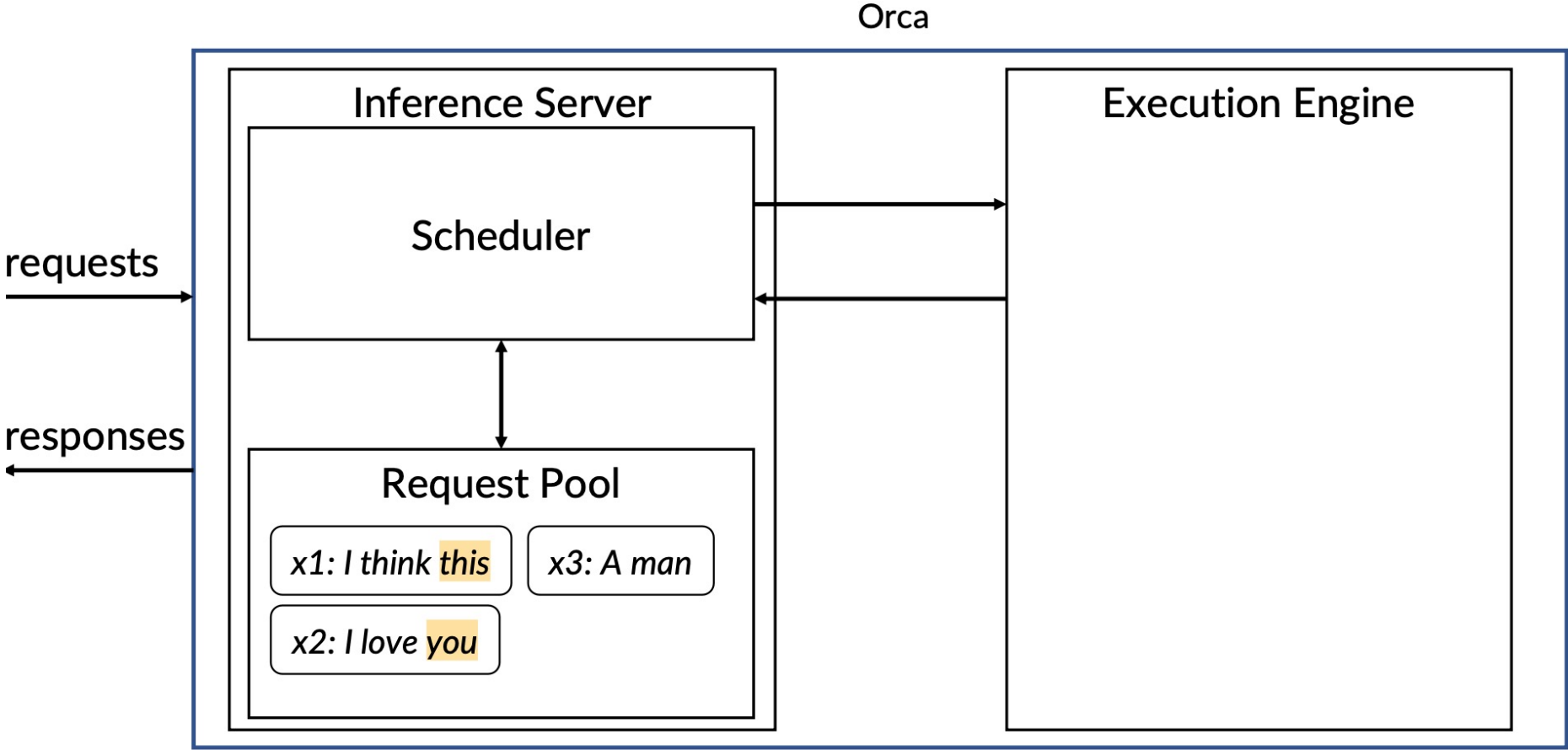
# Iteration-level scheduling



\* maximum batch size = 3

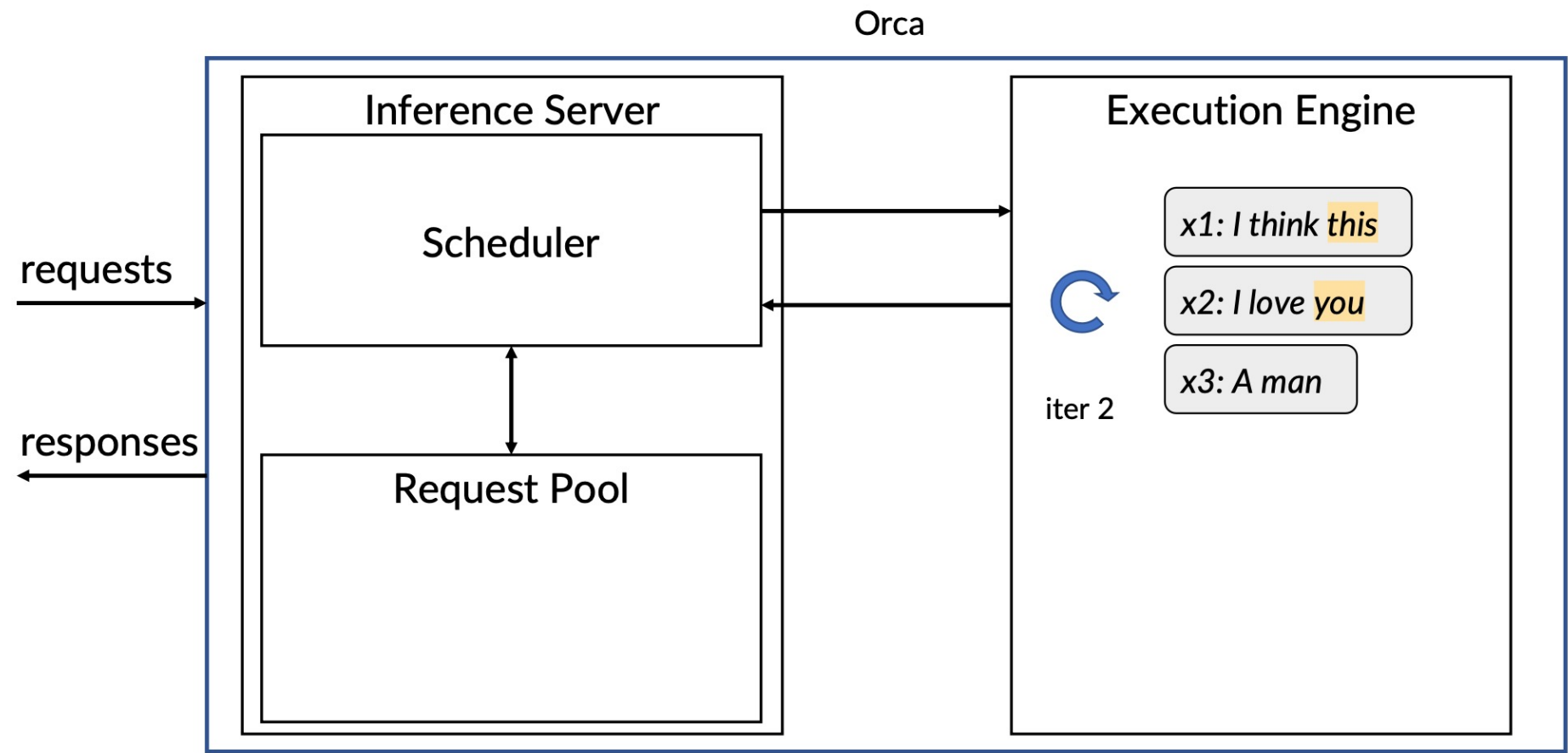


# Iteration-level scheduling



\* maximum batch size = 3

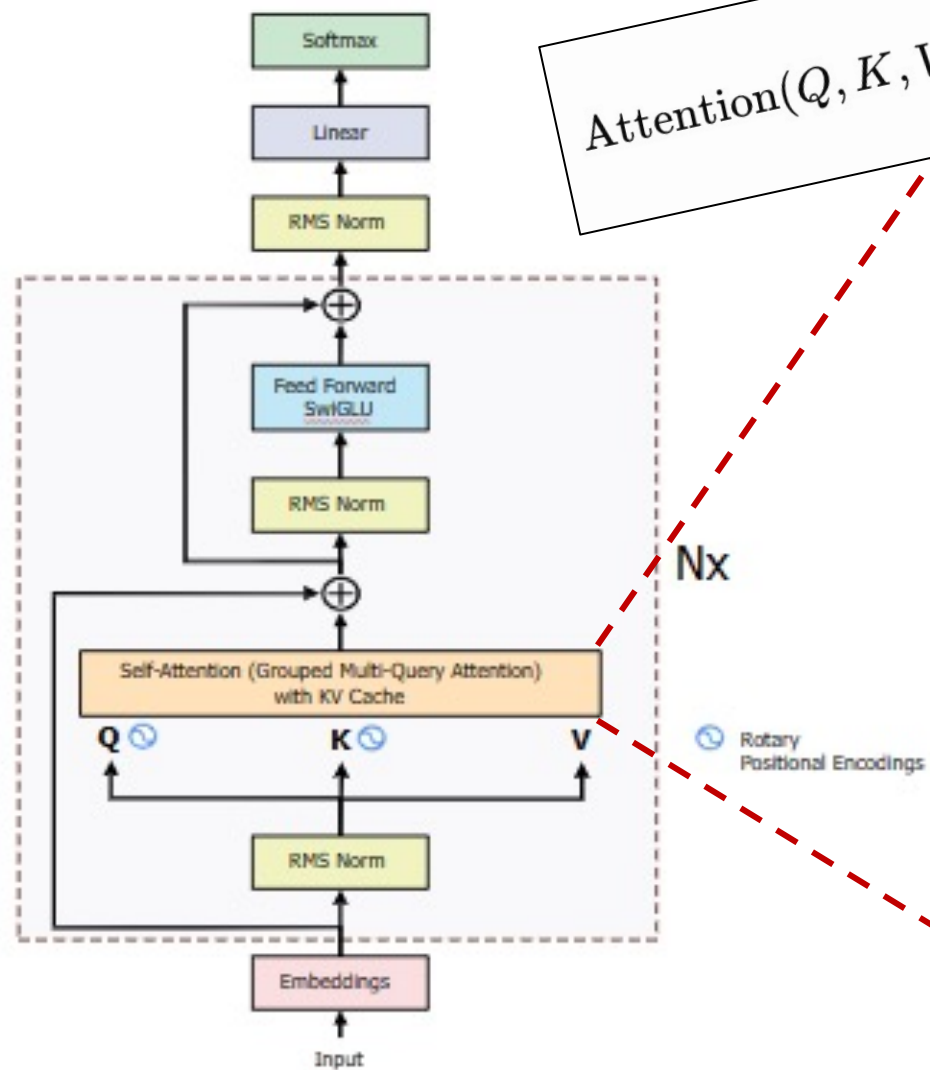
# Iteration-level scheduling



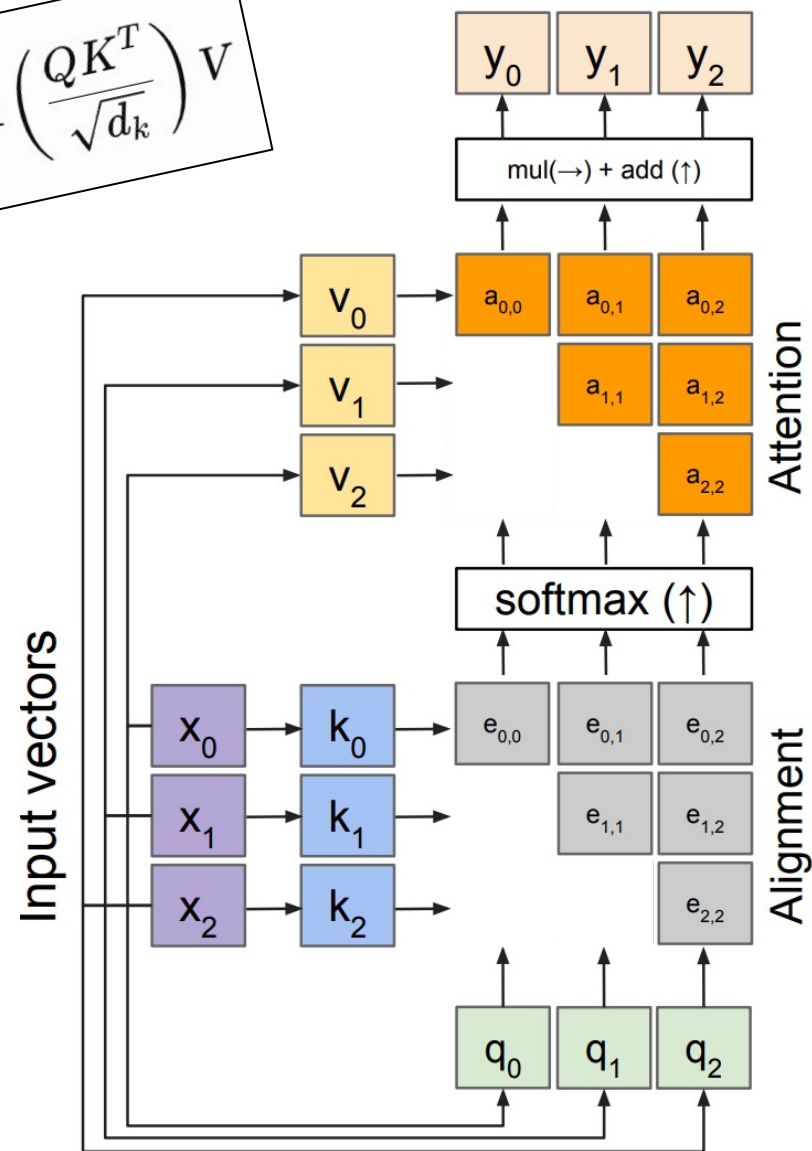
\* maximum batch size = 3

# Attention and KV Cache

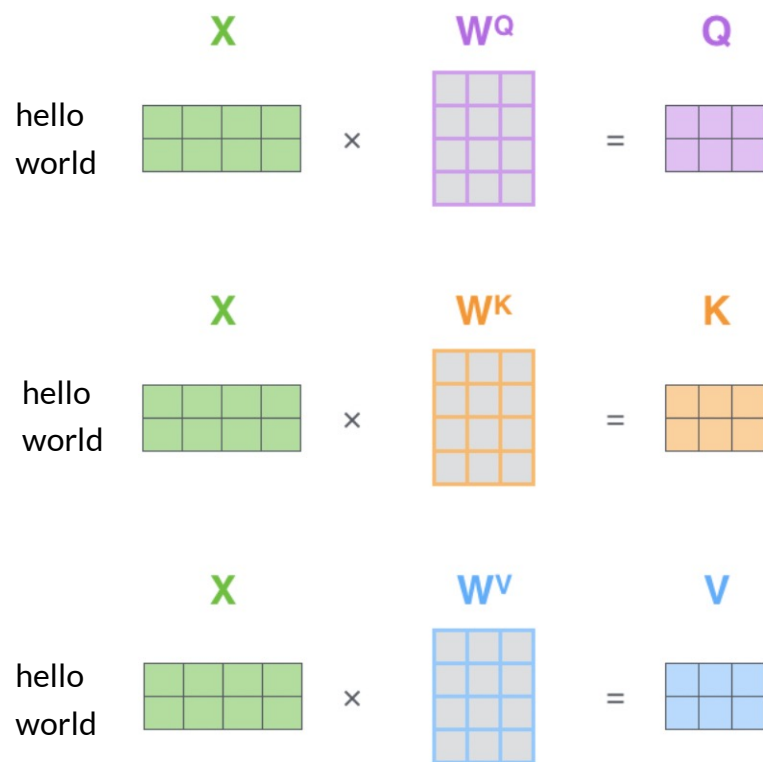
# Attention zoomed in



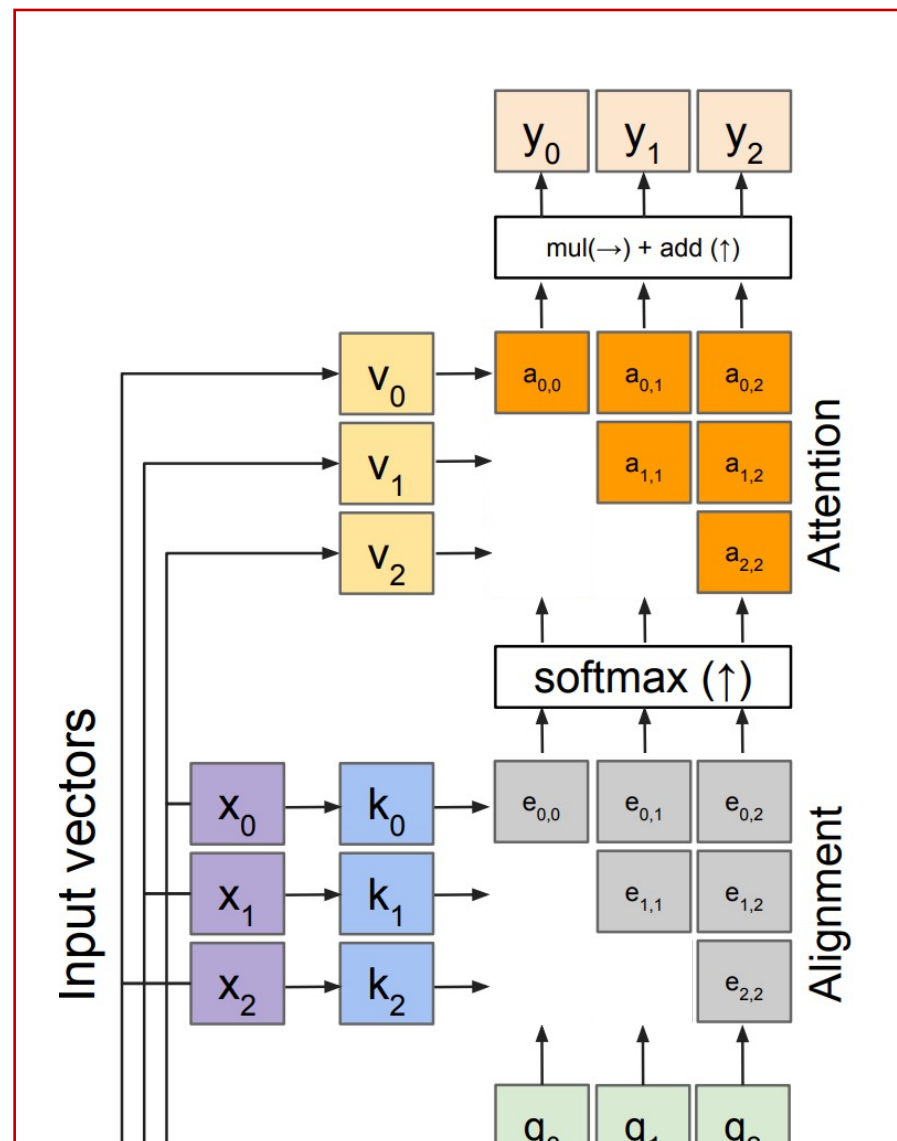
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



# Getting K, Q, V is expensive



Dimension: 4,096 for Llama3-8B



We can re-use K and V for previous tokens -> **KV Cache**

# 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	

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

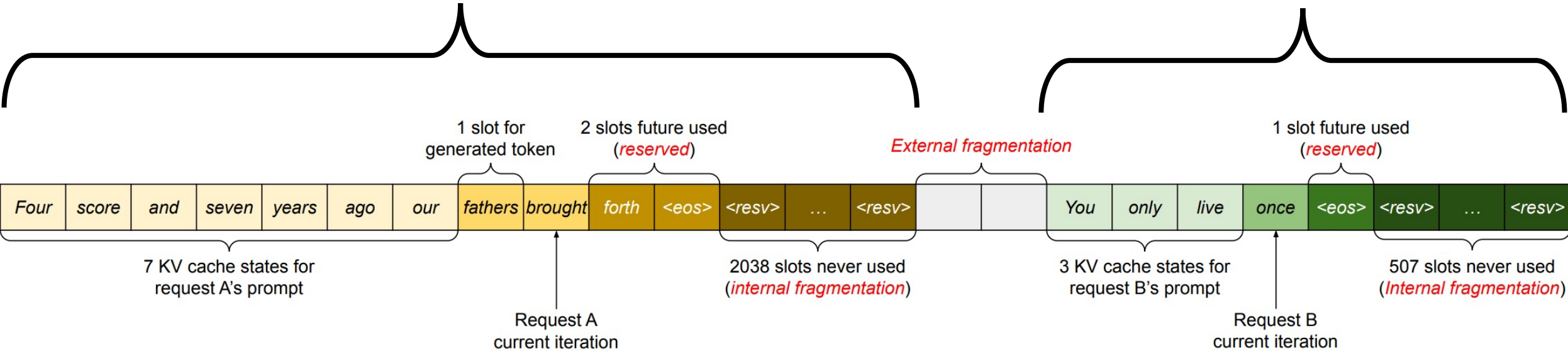
$T_1$	$T_2$	$T_3$	$T_4$	$T_5$	$T_6$	$T_7$	$T_8$
$S_1$	$S_1$	$S_1$	$S_1$	$S_1$	END	$S_6$	$S_6$
$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	$S_2$	END
$S_3$	$S_3$	$S_3$	$S_3$	END	$S_5$	$S_5$	$S_5$
$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	$S_4$	END	$S_7$



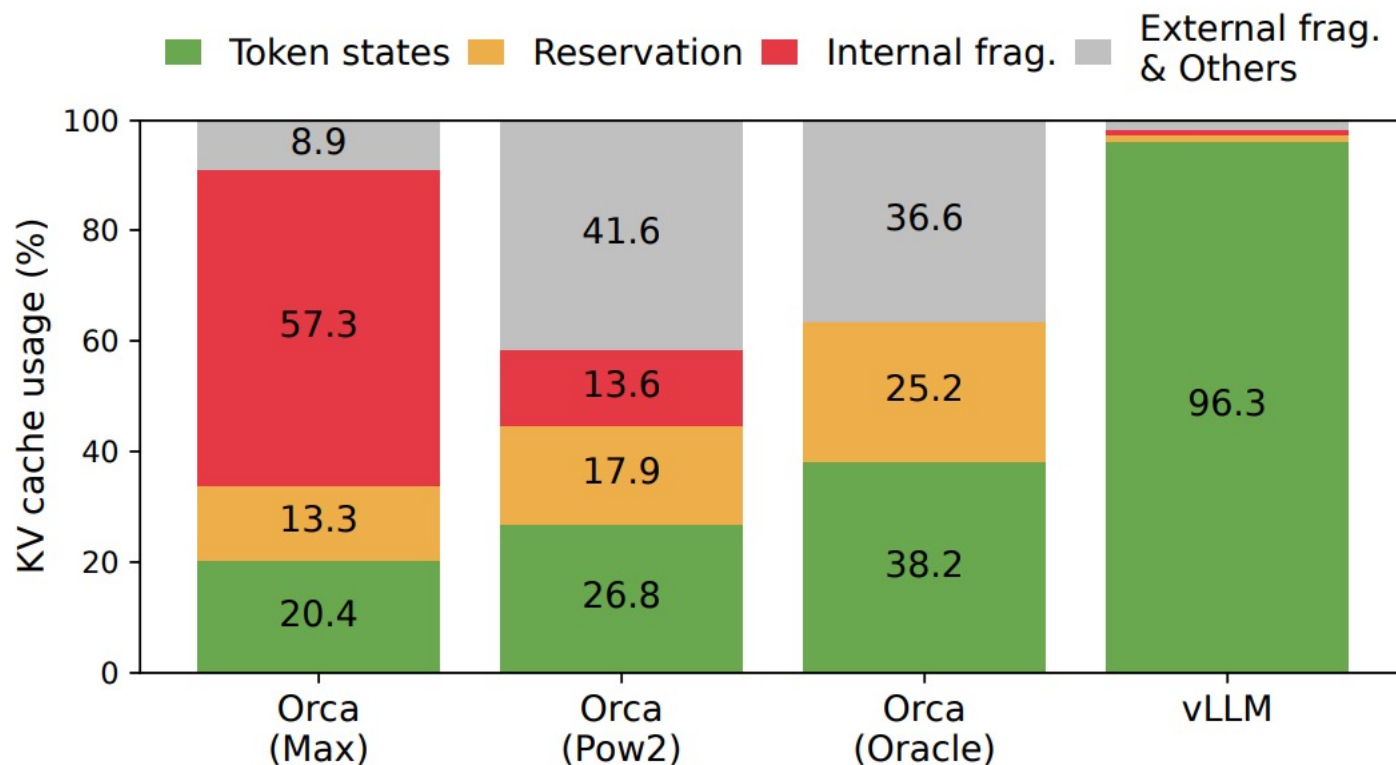
# Fragmentation -> Wasted GPU memory

Max 2048 tokens  
(e.g., Llama 2 context size: 4096 tokens)

Max 512 tokens

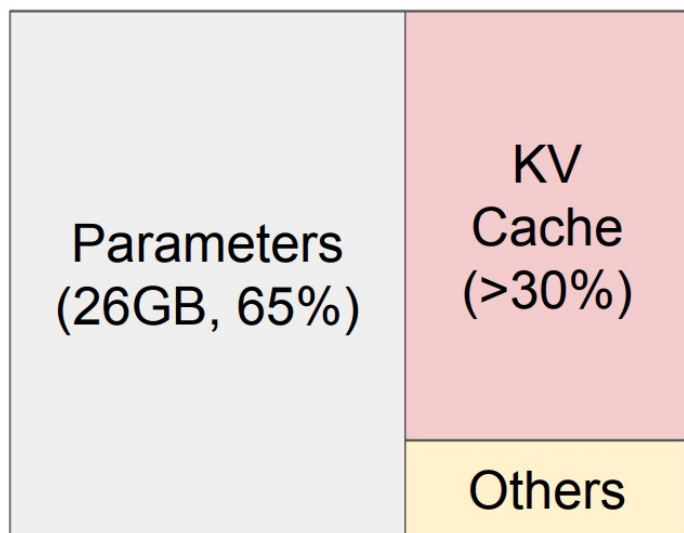


# Much GPU memory is wasted

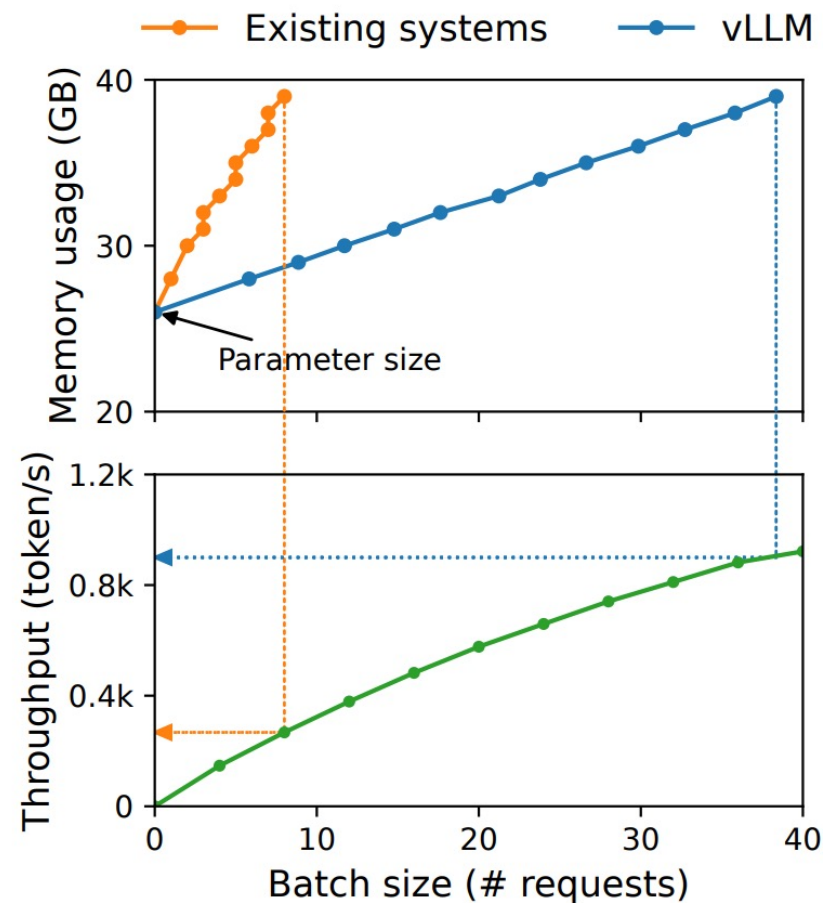


- **Orca (Oracle).** We assume the system has the knowledge of the lengths of the outputs that will be actually generated for the requests. This shows the upper-bound performance of Orca, which is infeasible to achieve in practice.
- **Orca (Pow2).** We assume the system over-reserves the space for outputs by at most 2×. For example, if the true output length is 25, it reserves 32 positions for outputs.
- **Orca (Max).** We assume the system always reserves the space up to the maximum sequence length of the model, i.e., 2048 tokens.

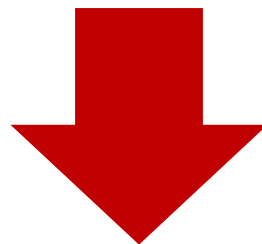
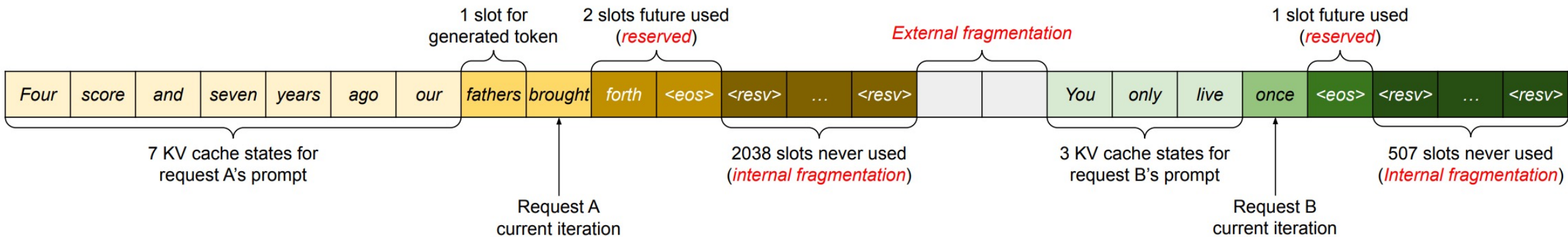
# GPU memory is insufficient



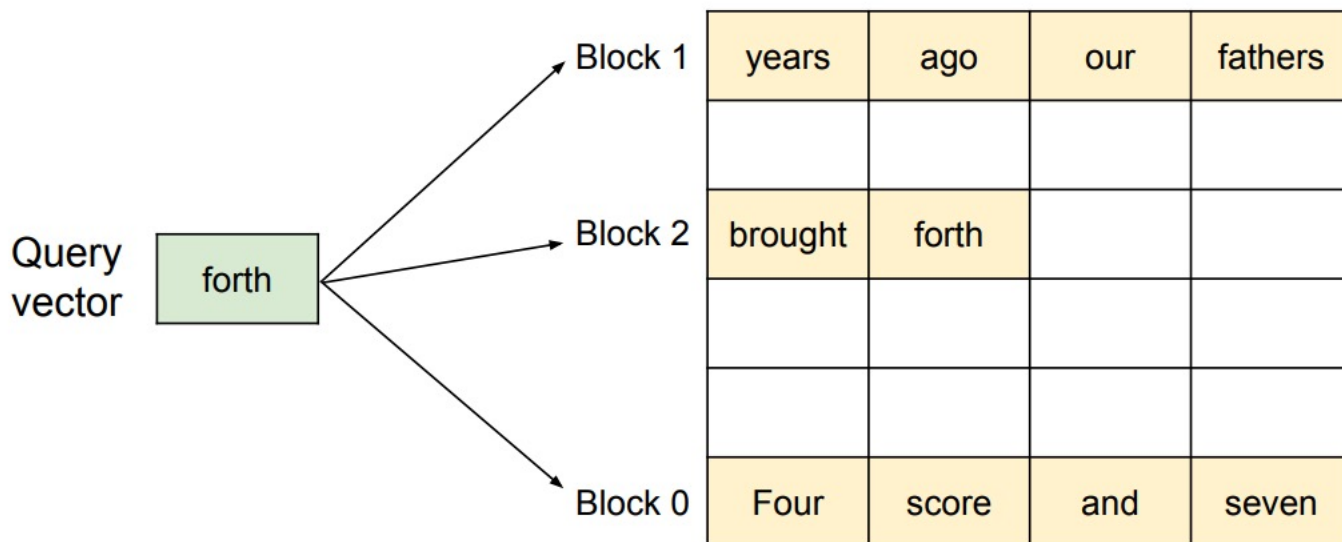
NVIDIA A100 40GB



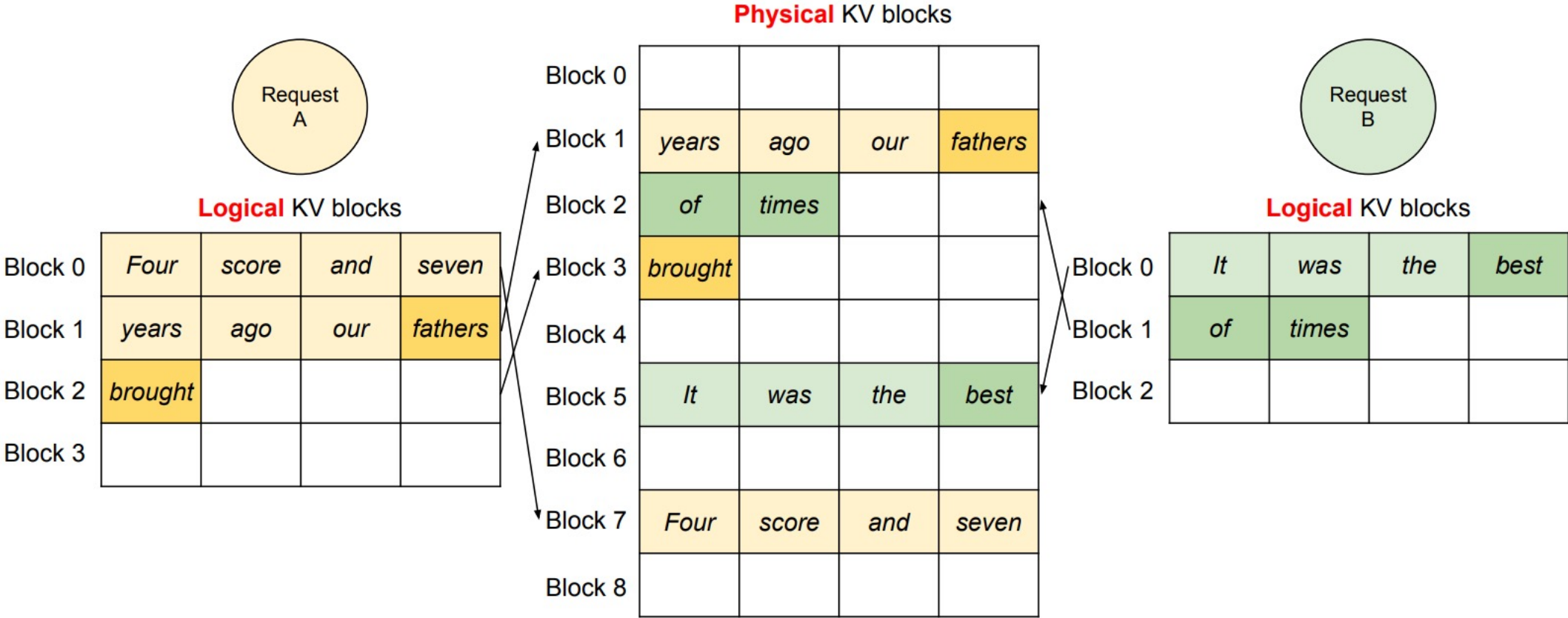
# PagedAttention and Scheduling



### Key and value vectors



# Pages can be shared by multiple prompts





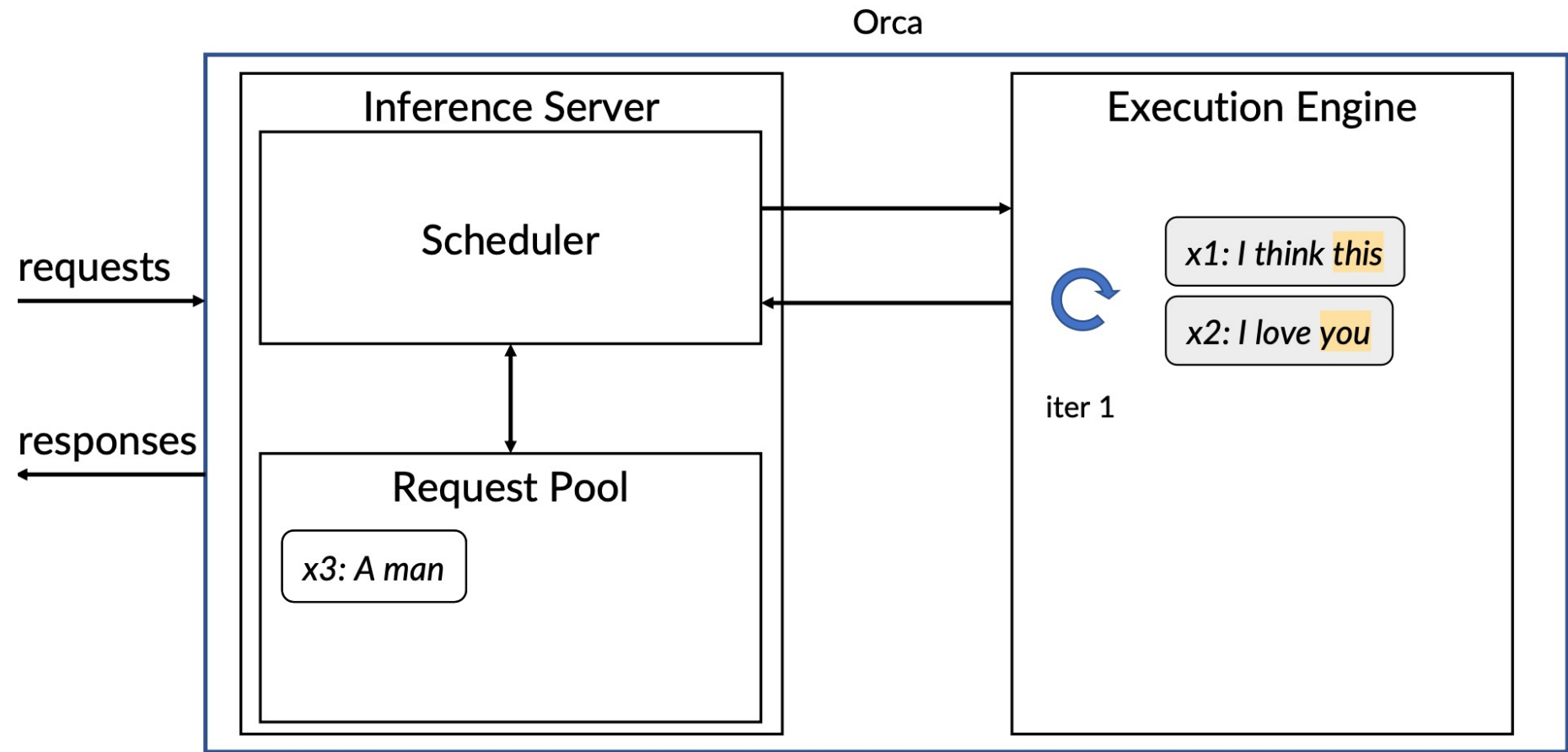
# Attention computation needs additional lookups

```
46 namespace vllm {  
89     int PARTITION_SIZE = 0; // zero means no partitioning.  
90     __device__ void paged_attention_kernel(  
91         float* __restrict__ exp_sums, // [num_seqs, num_heads, max_num_partitions]  
92         float* __restrict__ max_logits, // [num_seqs, num_heads,  
93         // max_num_partitions]  
94         scalar_t* __restrict__ out, // [num_seqs, num_heads, max_num_partitions,  
95         // head_size]  
96         const scalar_t* __restrict__ q, // [num_seqs, num_heads, head_size]  
97         const cache_t* __restrict__ k_cache, // [num_blocks, num_kv_heads,  
98         // head_size/x, block_size, x]  
99         const cache_t* __restrict__ v_cache, // [num_blocks, num_kv_heads,  
100        // head_size, block_size]  
101        const int num_kv_heads, // [num_heads]  
102        const float scale,  
103        const int* __restrict__ block_tables, // [num_seqs, max_num_blocks_per_seq]  
104        const int* __restrict__ seq_lens, // [num_seqs]  
105        const int max_num_blocks_per_seq,  
106        const float* __restrict__ alibi_slopes, // [num_heads]  
107        const int q_stride, const int kv_block_stride, const int kv_head_stride,  
108        const float* k_scale, const float* v_scale, const int tp_rank,  
109        const int blockspare_local_blocks, const int blockspare_vert_stride,  
110        const int blockspare_block_size, const int blockspare_head_sliding_step) {
```

# No more blocks -> Preempt (i.e., kick out) requests

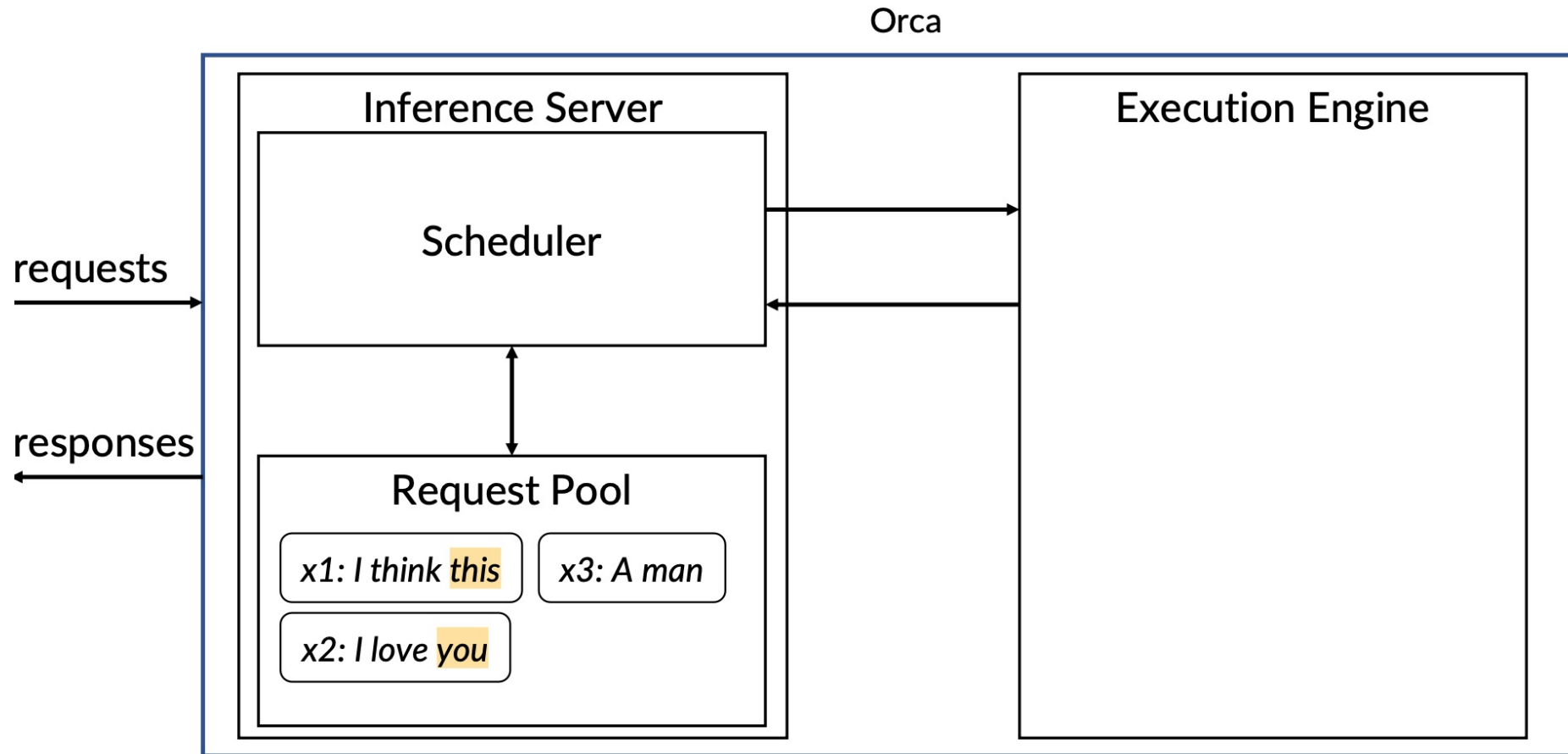
```
vllm > v1 > core > sched > scheduler.py > Scheduler > schedule
32 class Scheduler(SchedulerInterface):
115     def schedule(self) -> SchedulerOutput:
188
189         while True:
190             new_blocks = self.kv_cache_manager.allocate_slots(
191                 request, num_new_tokens)
192             if new_blocks is None:
193                 # The request cannot be scheduled.
194                 # Preempt the lowest-priority request.
195                 preempted_req = self.running.pop()
196                 self.kv_cache_manager.free(preempted_req)
197                 preempted_req.status = RequestStatus.PREEMPTED
198                 preempted_req.num_computed_tokens = 0
199                 if self.log_stats:
200                     preempted_req.record_event(
201                         EngineCoreEventType.PREEMPTED, scheduled_timestamp)
202
203                 self.waiting.appendleft(preempted_req)
204                 preempted_reqs.append(preempted_req)
205                 if preempted_req == request:
206                     # No more request to preempt.
207                     can_schedule = False
208                     break
209             else:
210                 # The request can be scheduled.
211                 can_schedule = True
212                 break
213         if not can_schedule:
214             break
215         assert new_blocks is not None
```

# Iteration-level scheduling



\* maximum batch size = 3

# Iteration-level scheduling



***X2 may be pre-empted for the next iteration, processing only X1***

# Summary

- **KV cache** allows avoiding expensive operations for previous tokens
- vllm proposes ***PagedAttention***, which incrementally allocates memory
- For each iteration, determines what requests to ***preempt*** / run

Questions?