

Efficient LLM Decoding

Large Language Models: Introduction and Recent Advances

ELL881 · AIL821



Yatin Nandwani
Research Scientist, IBM Research

Training Vs Inference in LLMs

Forward Pass through an LLM

Transformer based LLM (θ)

<s>	The	cat	sat	on	a	mat	</s>
0	1	2	3	4	5	6	7



Forward Pass through an LLM

Probability distribution over all the tokens at each step (simultaneously)

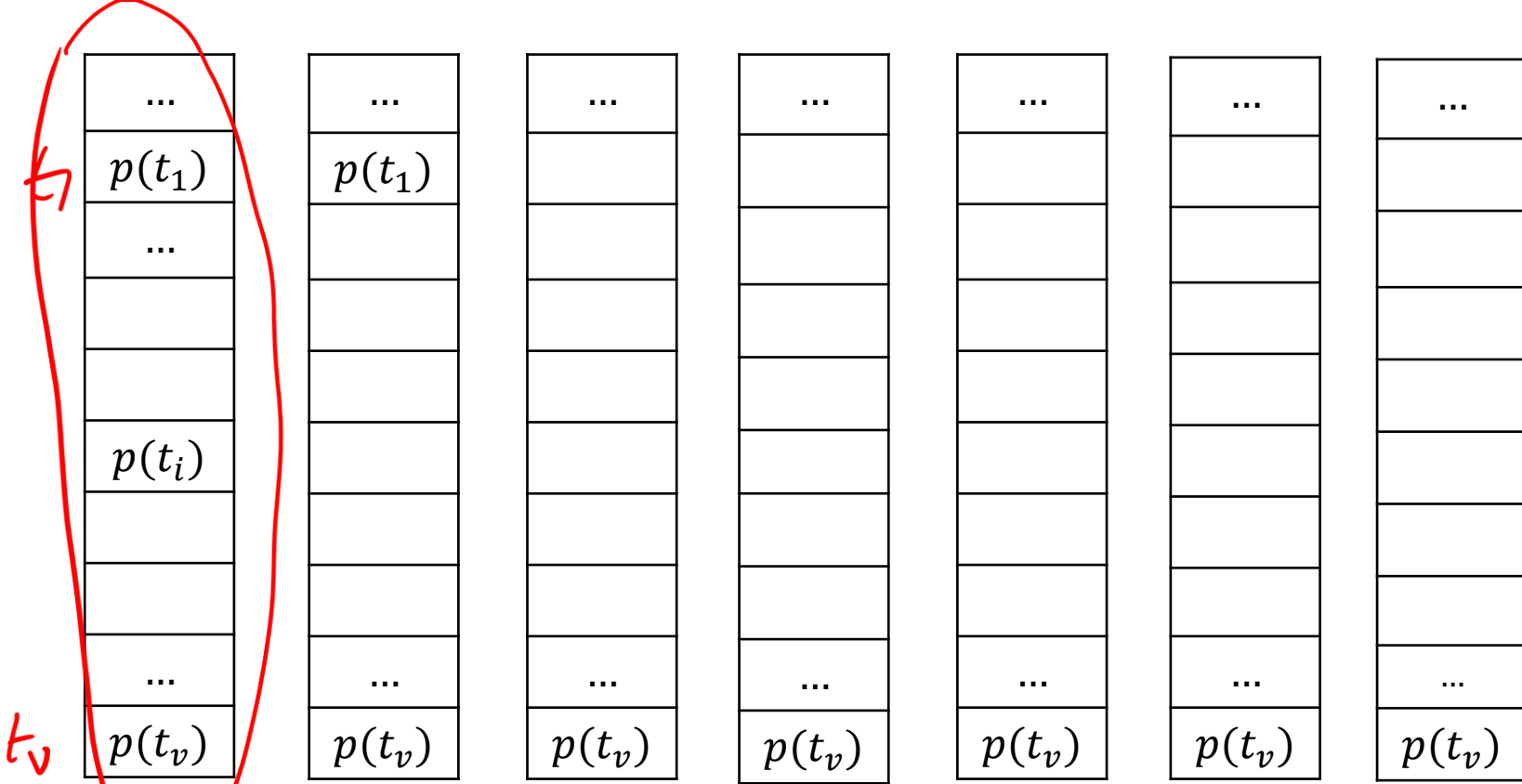
Transformer based LLM (θ)

<s>	The	cat	sat	on	a	mat	</s>
0	1	2	3	4	5	6	7

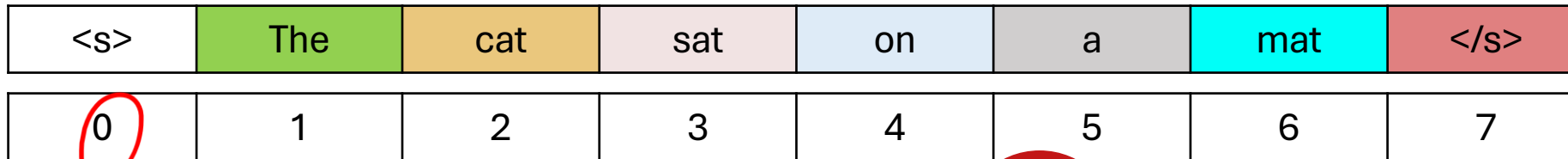


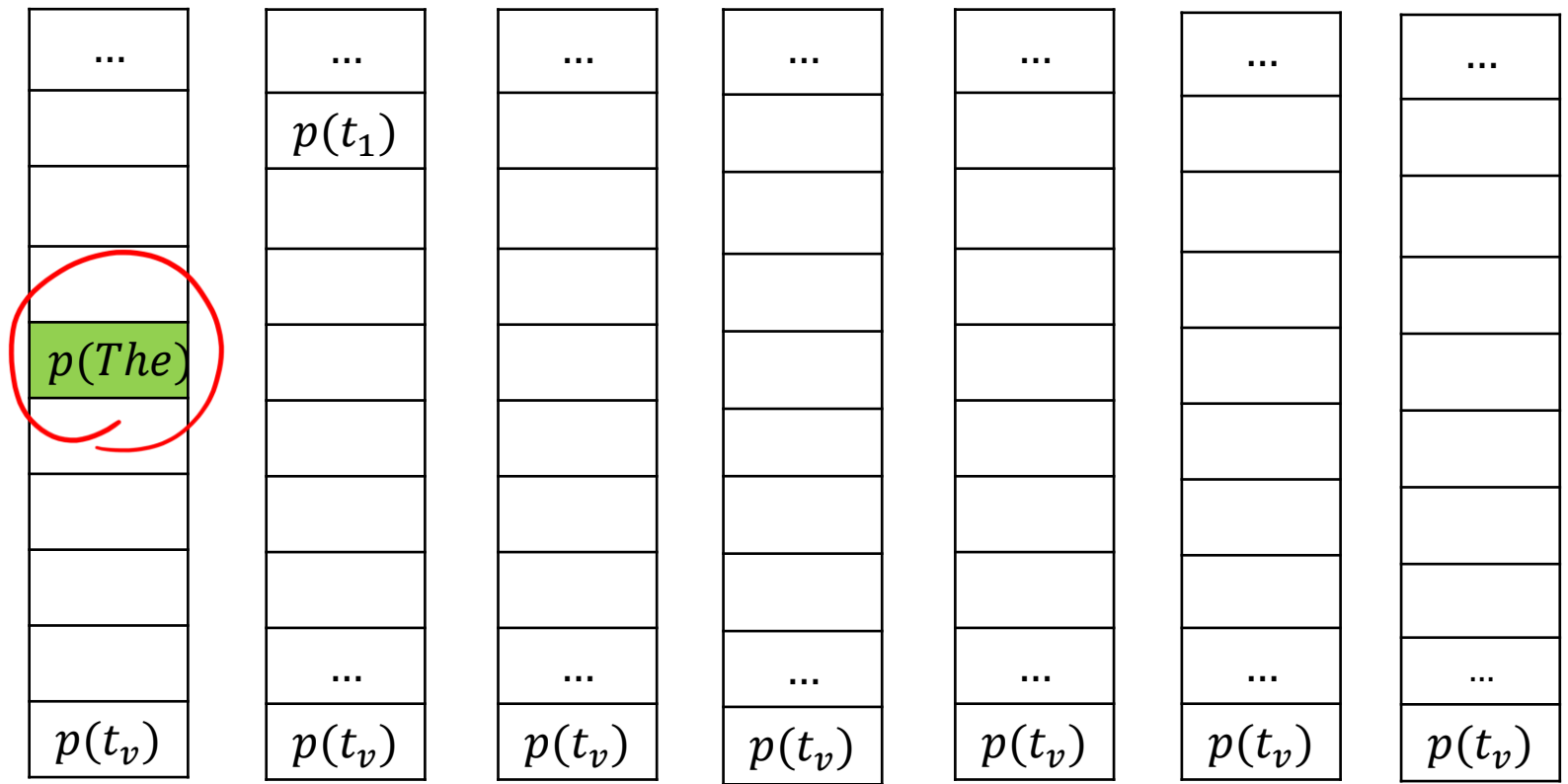
Forward Pass through an LLM

Probability distribution over all the tokens at each step (simultaneously)



Transformer based LLM (θ)

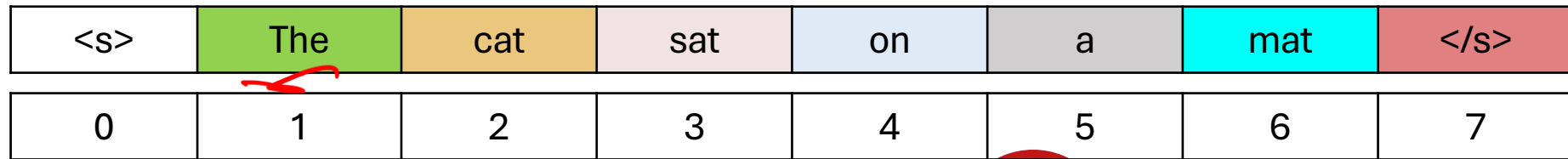


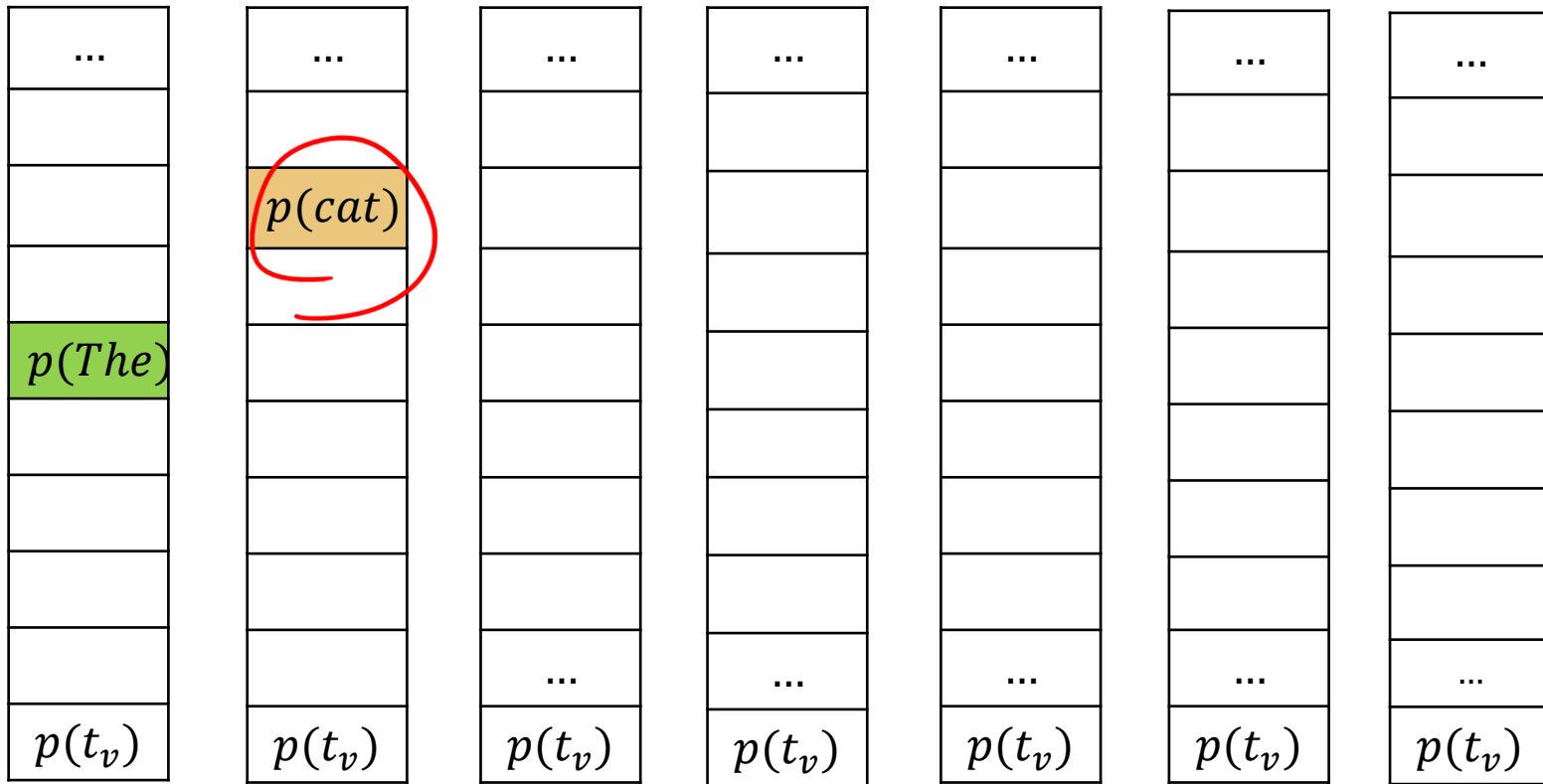


Forward Pass through an LLM

Train to maximize prob. of **The** at step 0

Transformer based LLM (θ)





Forward Pass through an LLM

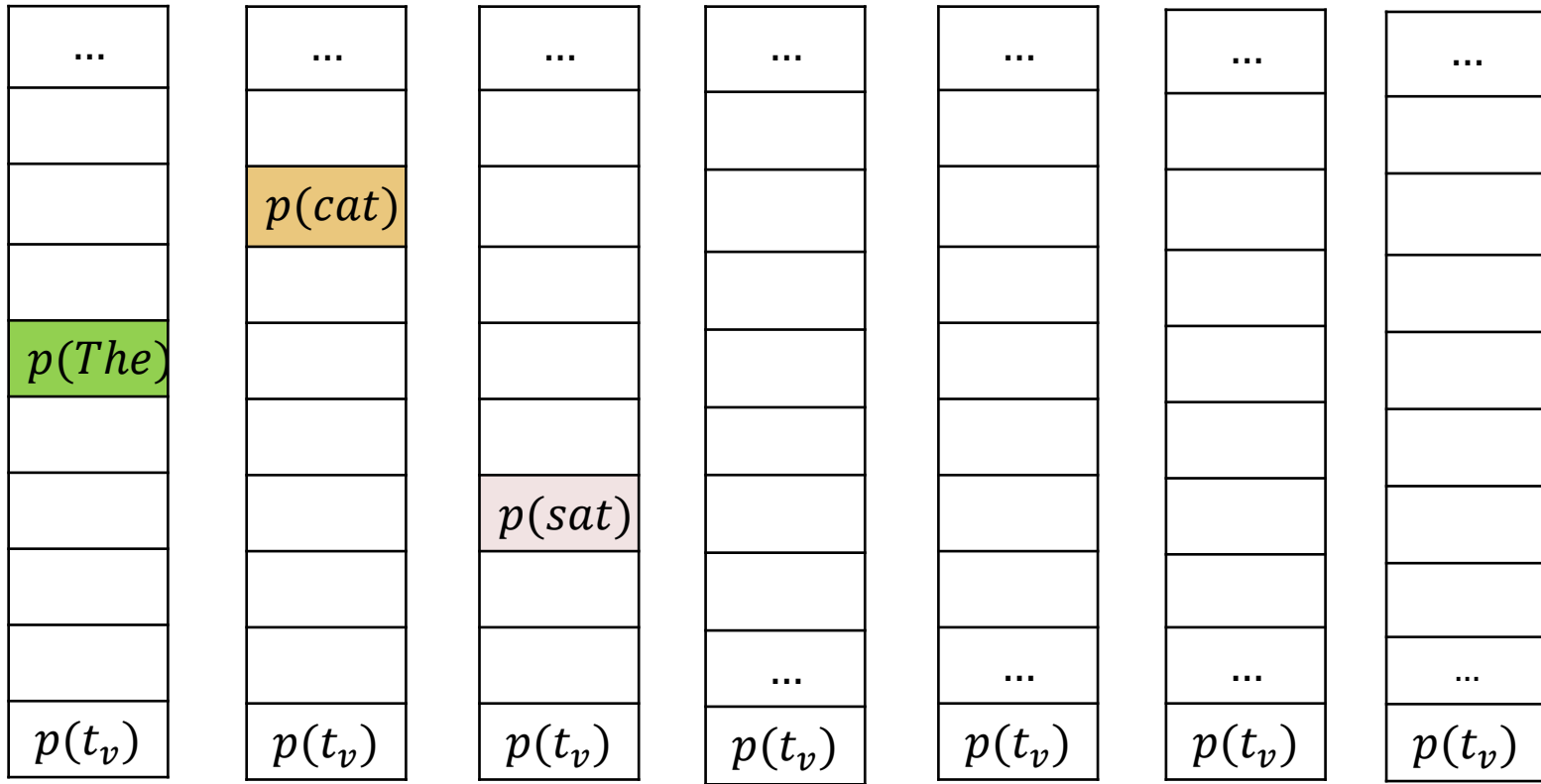
Train to maximize prob. of **cat** at step 1

Transformer based LLM (θ)

<s> The cat sat on a mat </s>

0 1 2 3 4 5 6 7





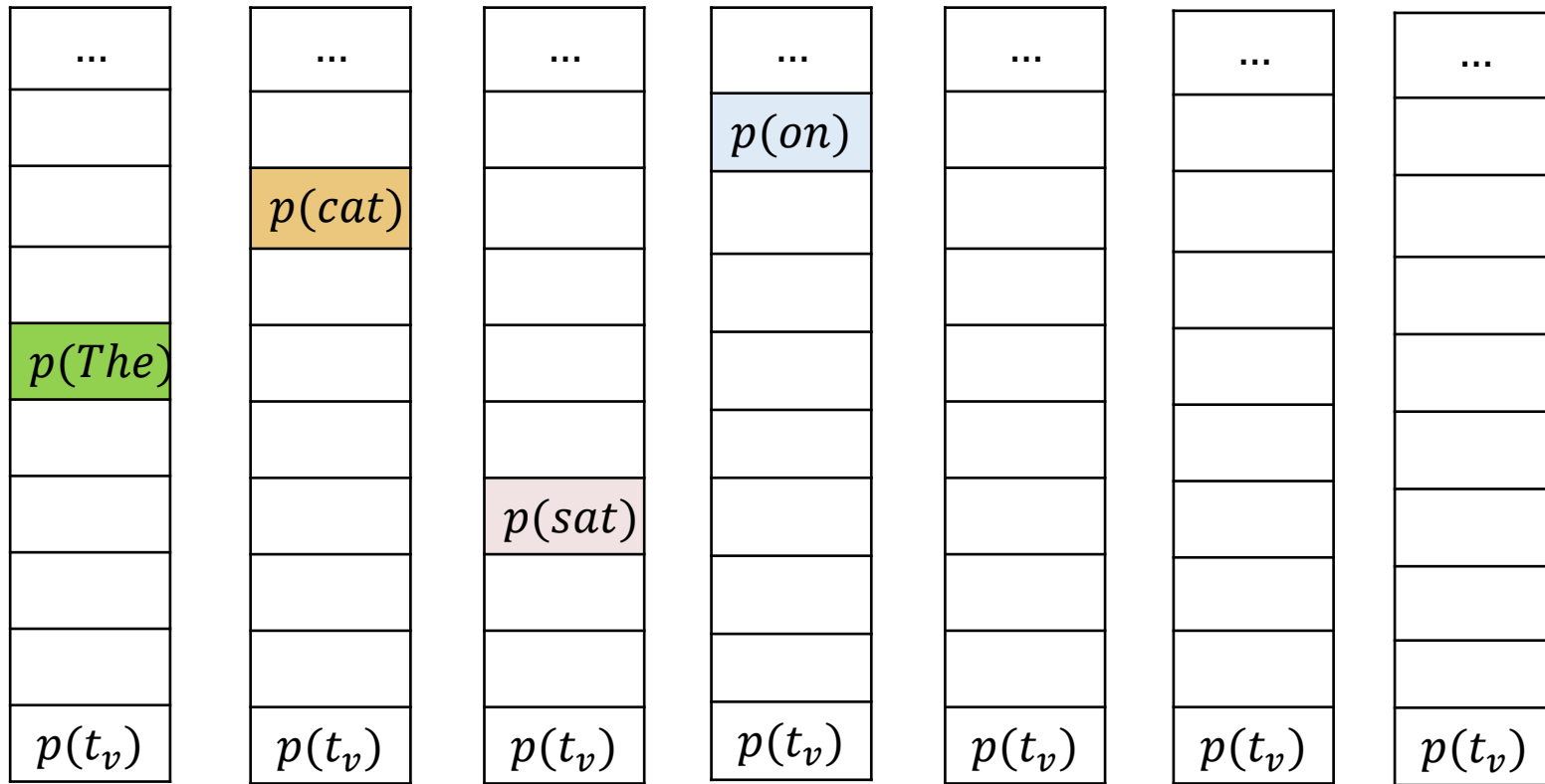
Forward Pass through an LLM

Train to maximize prob. of **sat** at step 2

Transformer based LLM (θ)

<s>	The	cat	sat	on	a	mat	</s>
0	1	2	3	4	5	6	7





Forward Pass through an LLM

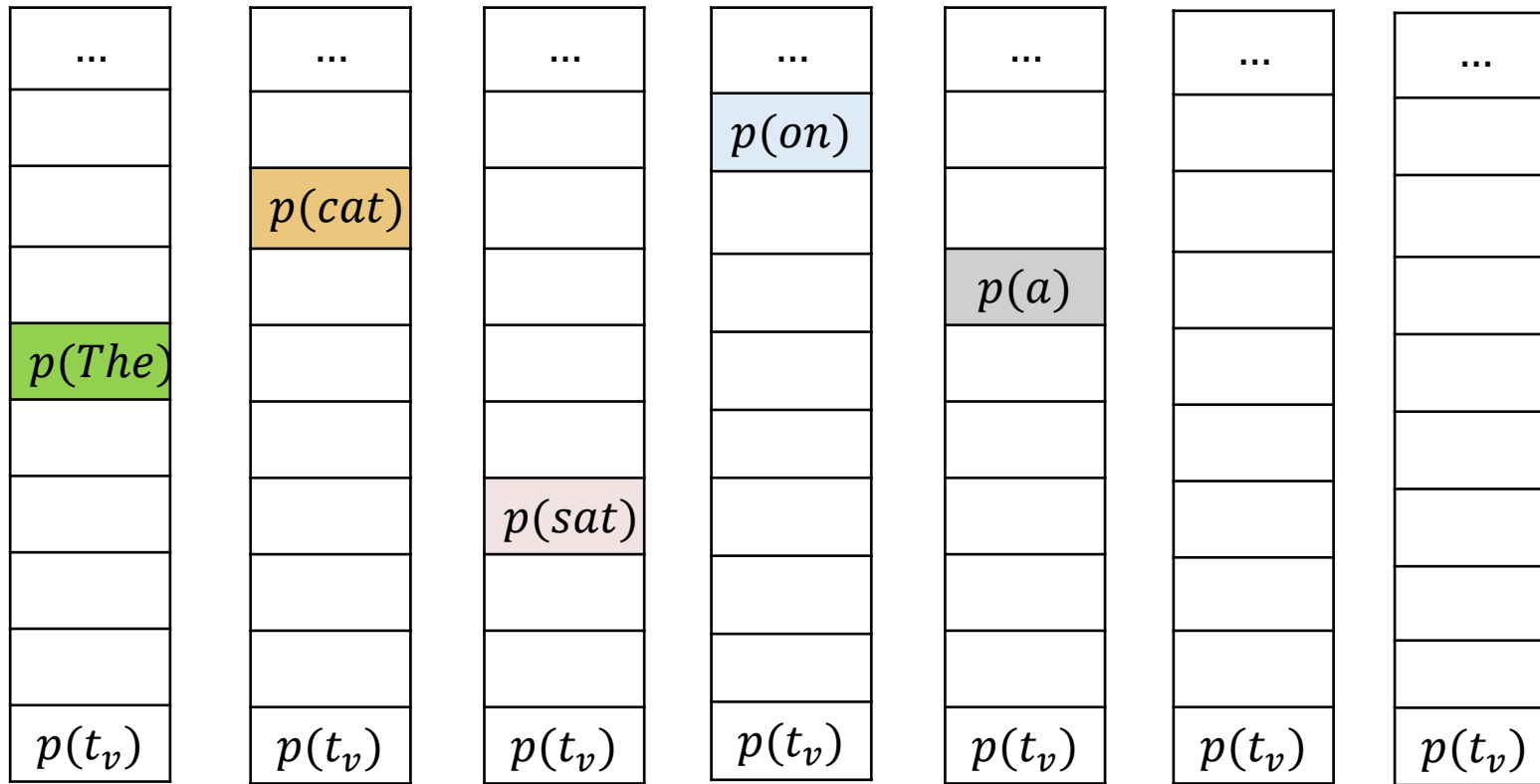
Train to maximize prob. of **on** at step 3

Transformer based LLM (θ)

<s> The cat sat on a mat </s>

0 1 2 3 4 5 6 7





Forward Pass through an LLM

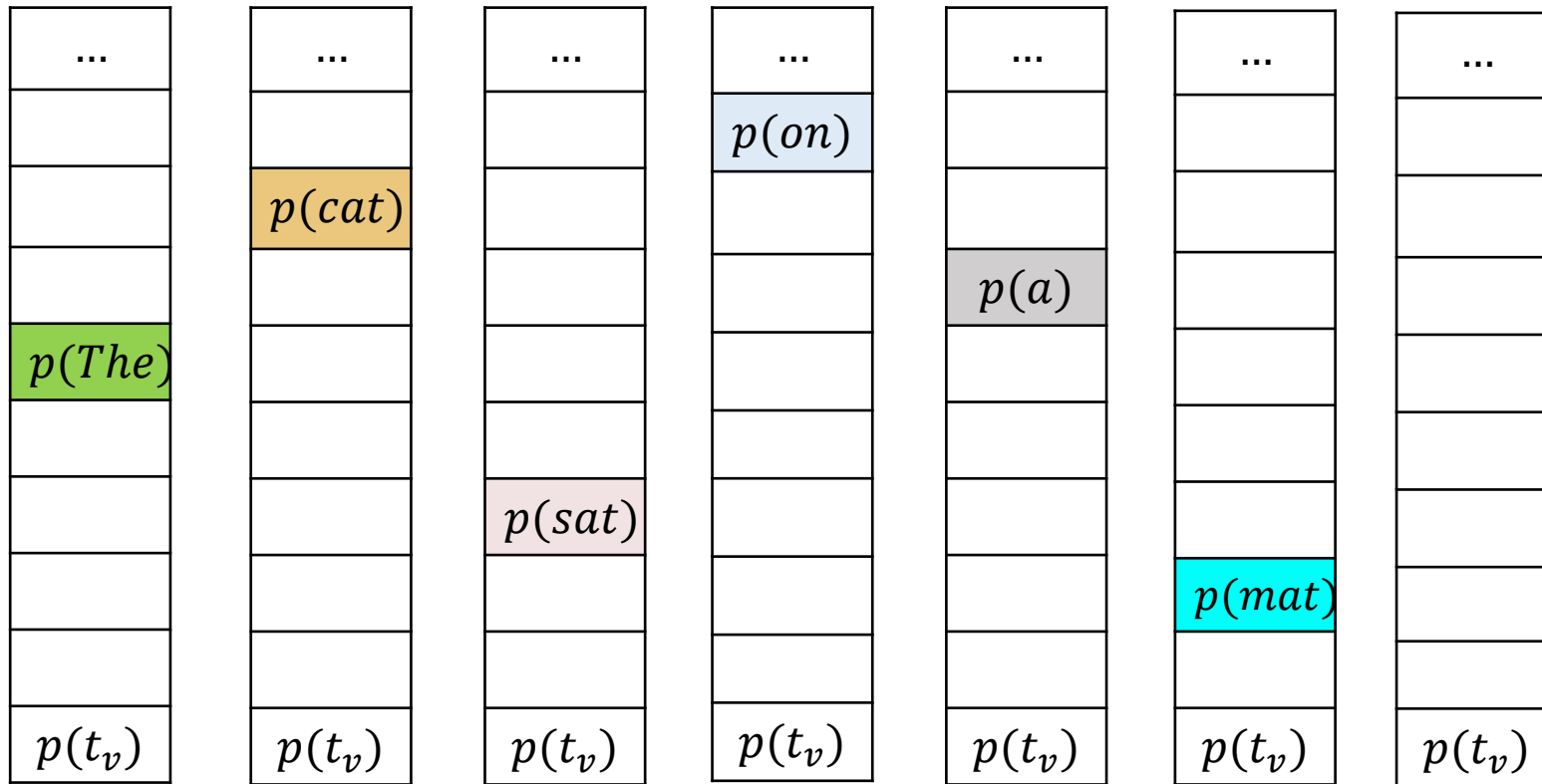
Train to maximize prob. of a at step 4

Transformer based LLM (θ)

<s> The cat sat on a mat </s>

0 1 2 3 4 5 6 7





Forward Pass through an LLM

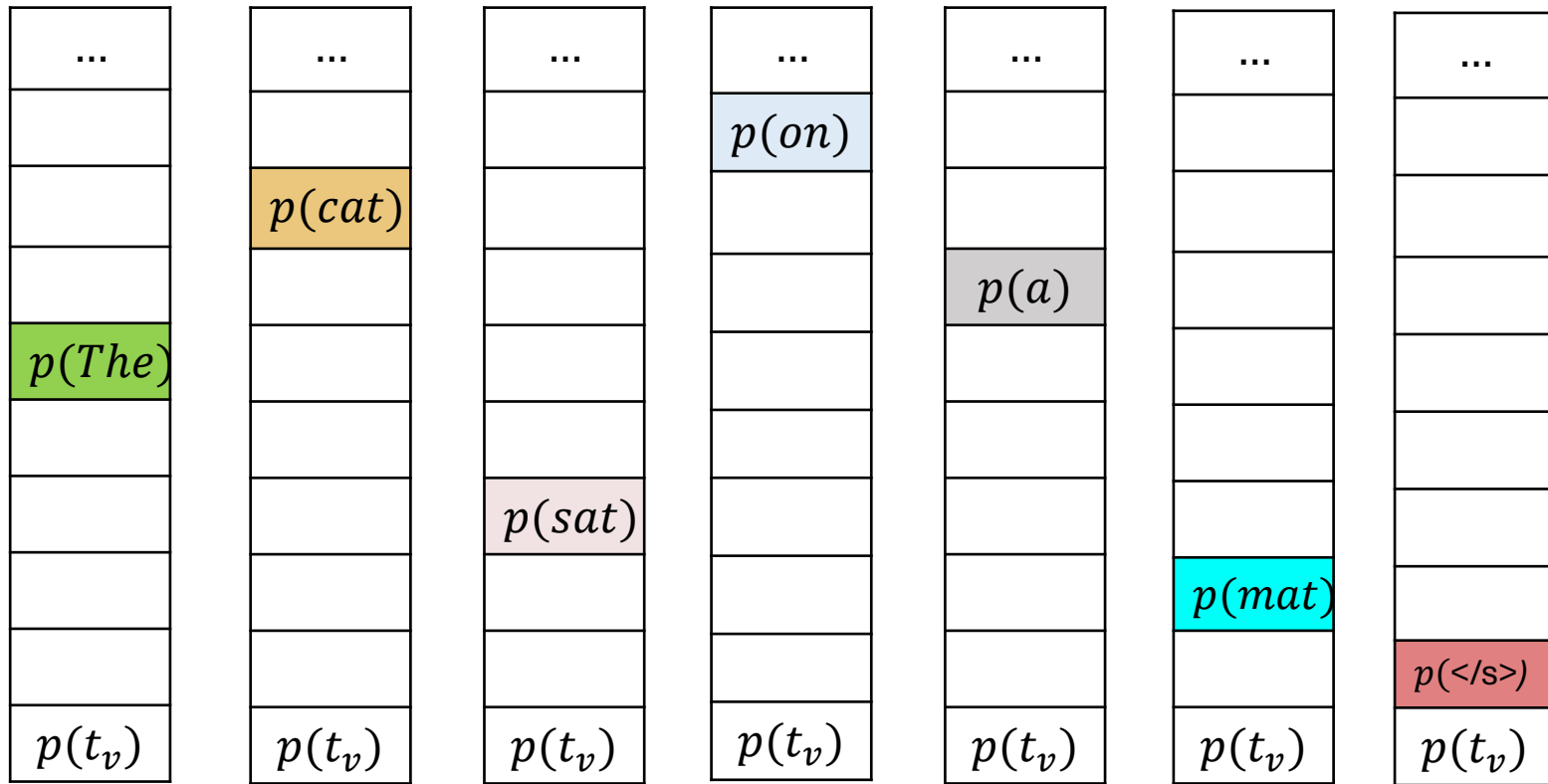
Train to maximize prob. of **mat** at step 5

Transformer based LLM (θ)

<s> The cat sat on a mat </s>

0 1 2 3 4 5 6 7





Forward Pass through an LLM

Train to maximize prob. of **</s>** at step 6

Transformer based LLM (θ)

<s> The cat sat on a mat </s>

0 1 2 3 4 5 6 7



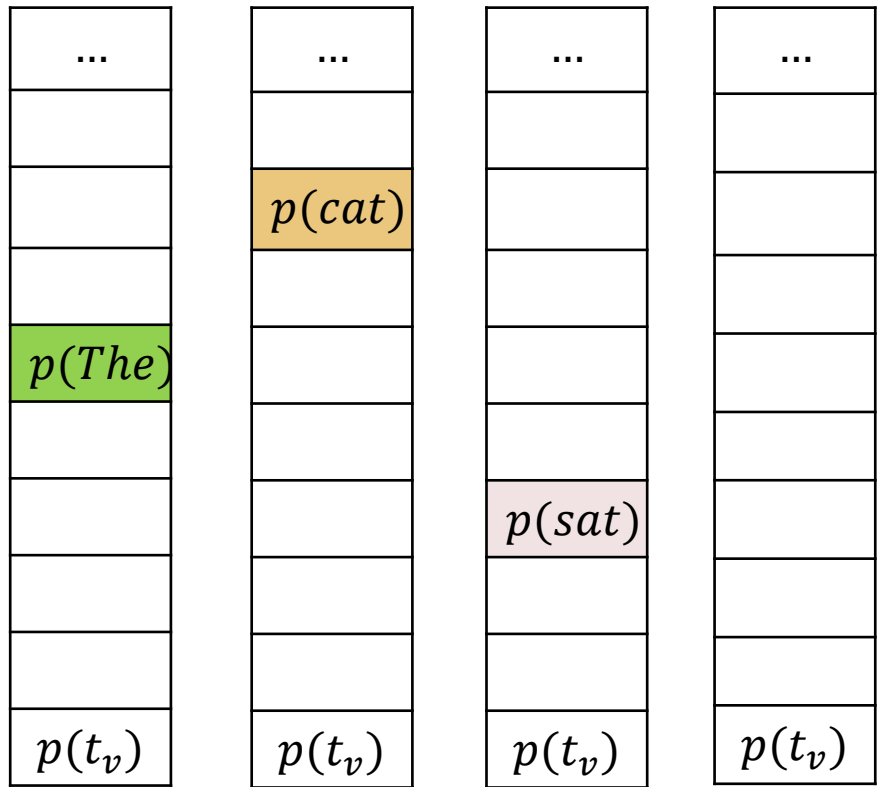
Inference through an LLM

Forward Pass (#1)

Transformer based LLM (θ)

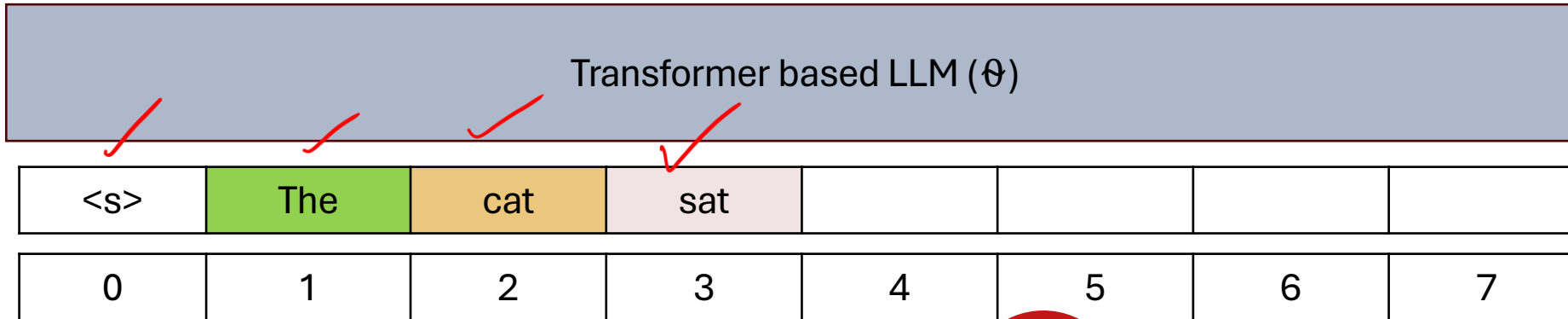
<s>	The	cat	sat				
0	1	2	3	4	5	6	7

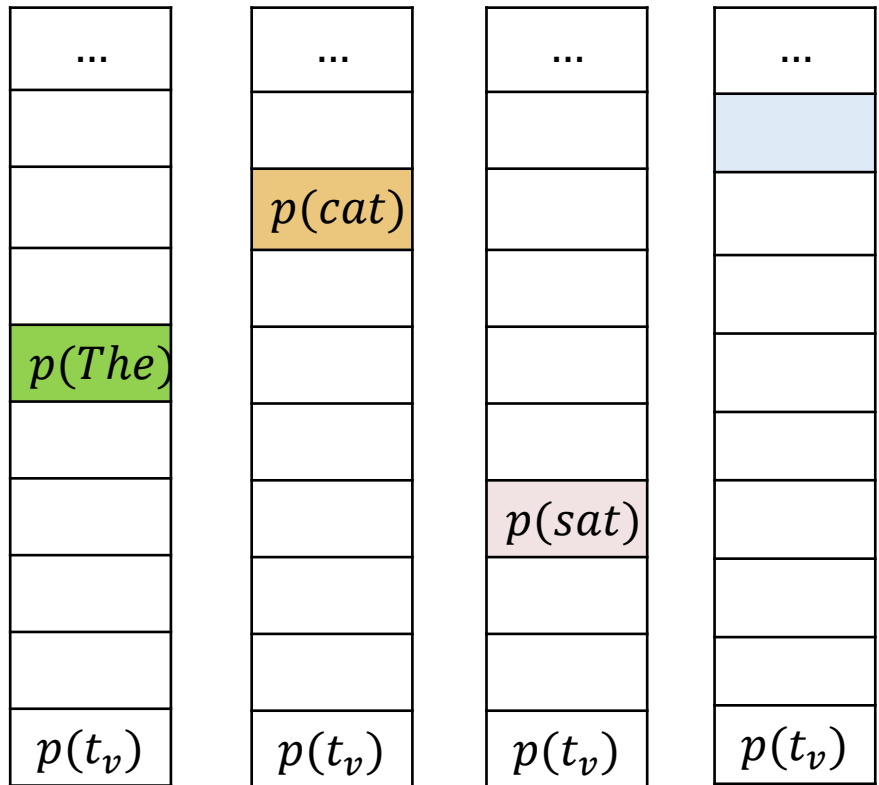




Inference through an LLM

Prob. Dist. at all steps

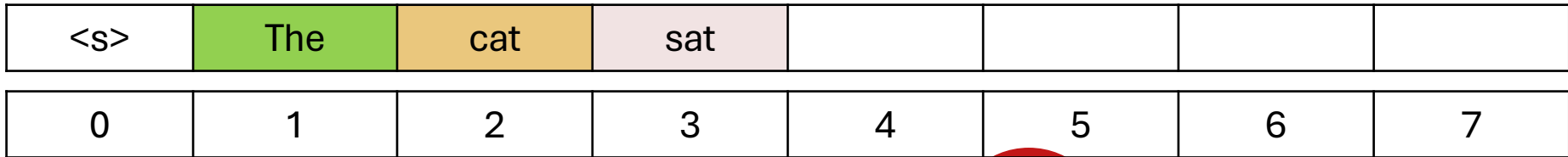


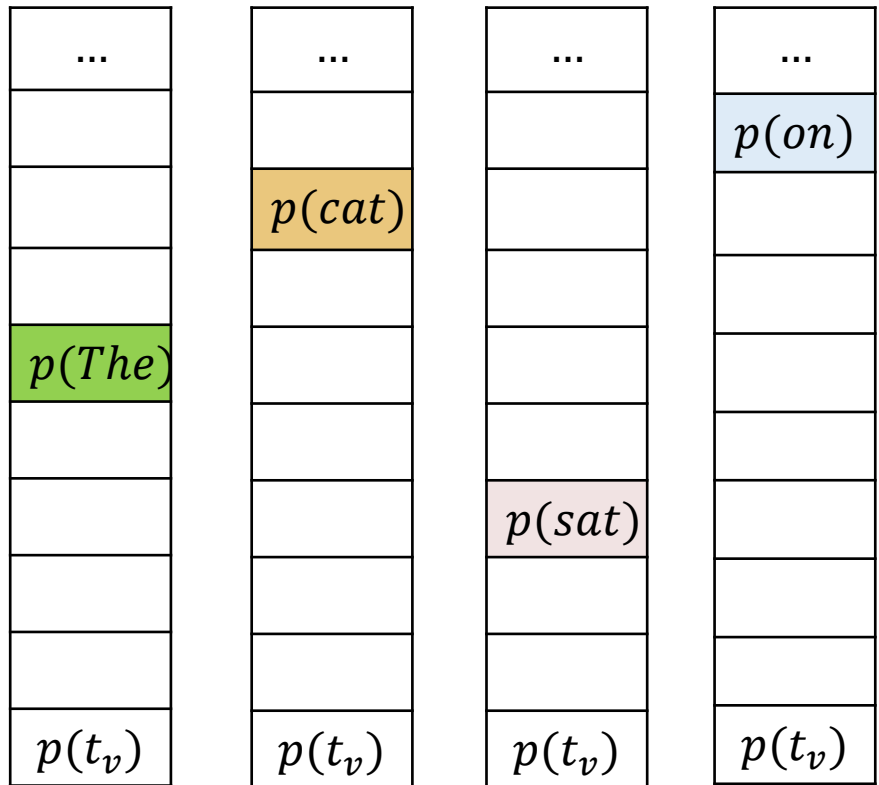


Inference through an LLM

Pick the token having max. probability at step 3

Transformer based LLM (θ)

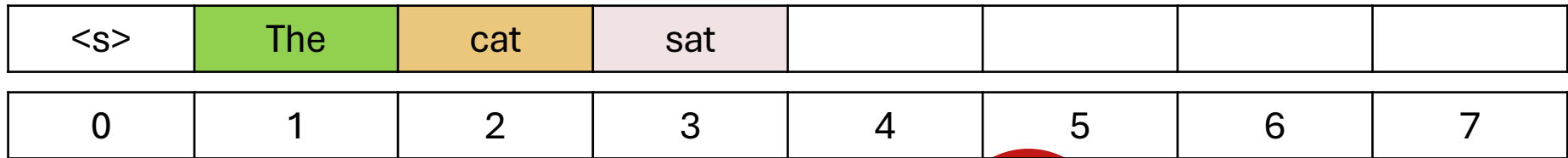


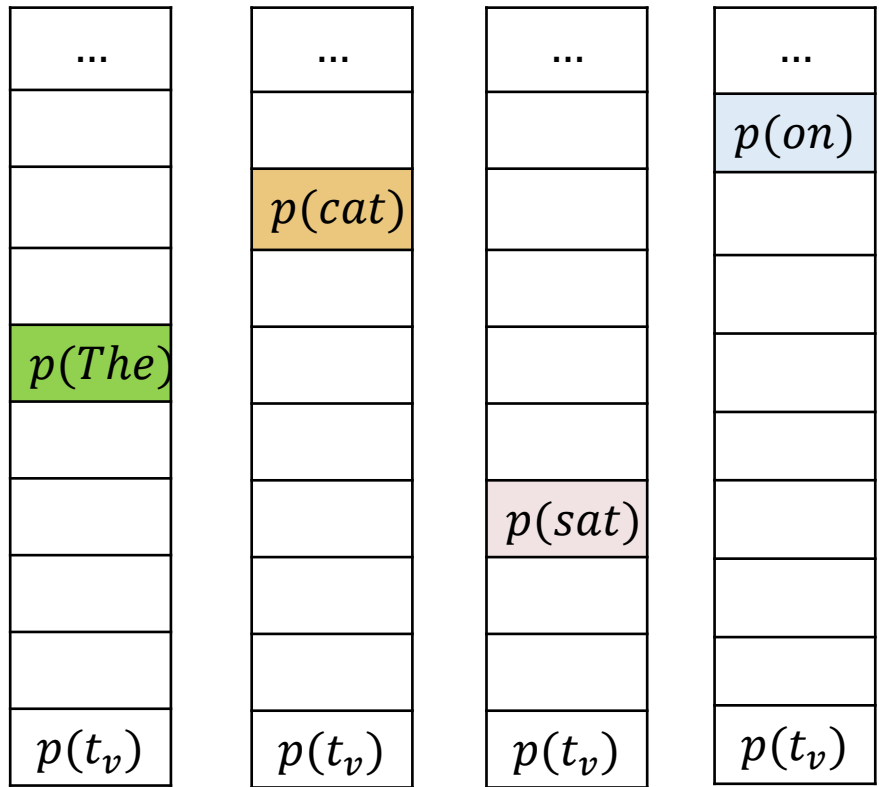


Inference through an LLM

Pick the token having max. probability at step 3

Transformer based LLM (θ)

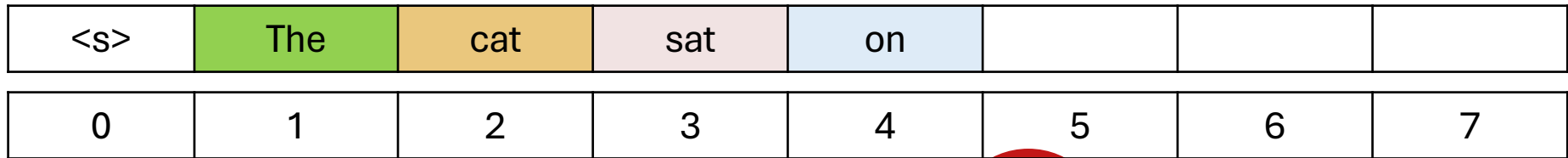




Inference through an LLM

Fill at step 4

Transformer based LLM (θ)



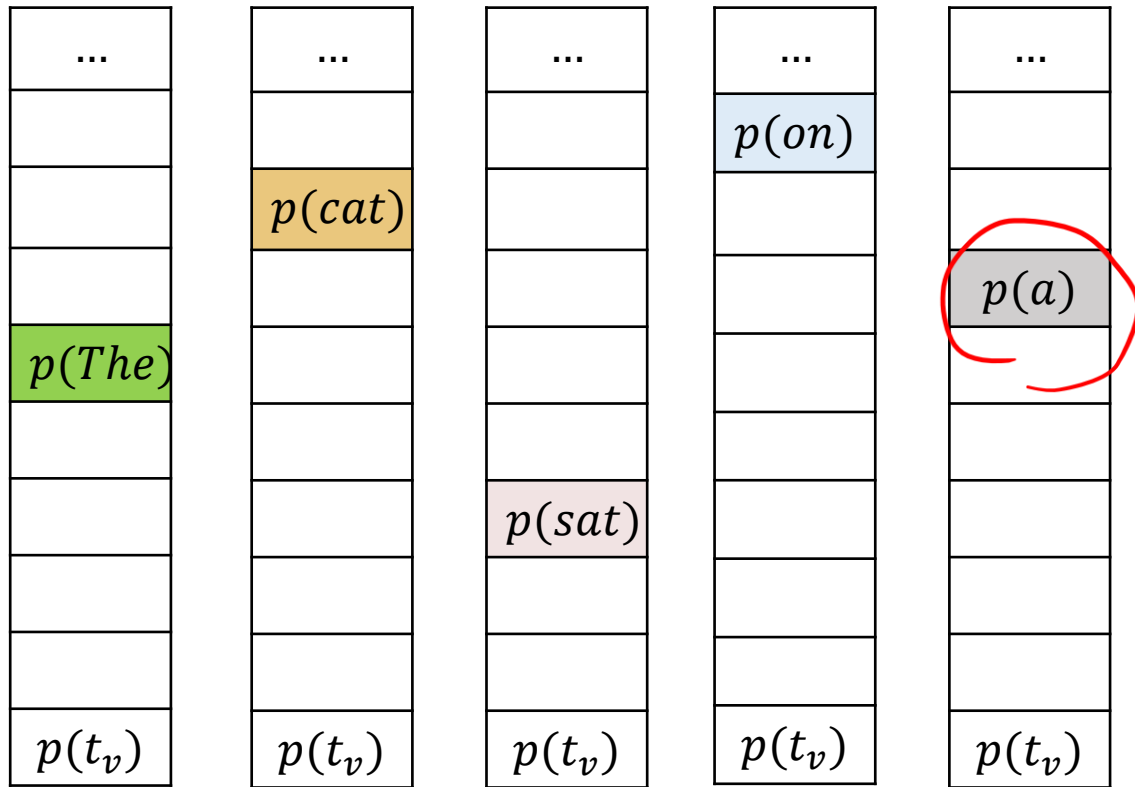
Inference through an LLM

Fwd. Pass (#2)

Transformer based LLM (θ)

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7

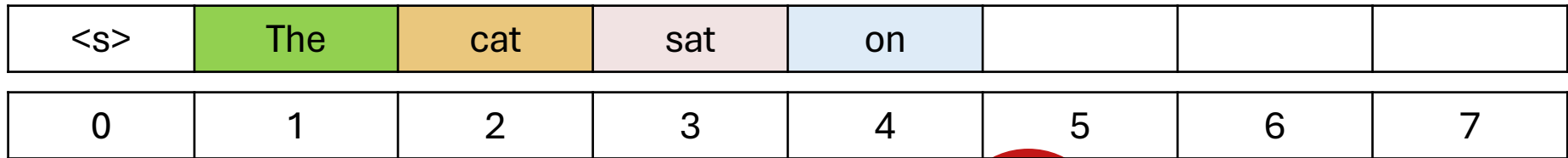


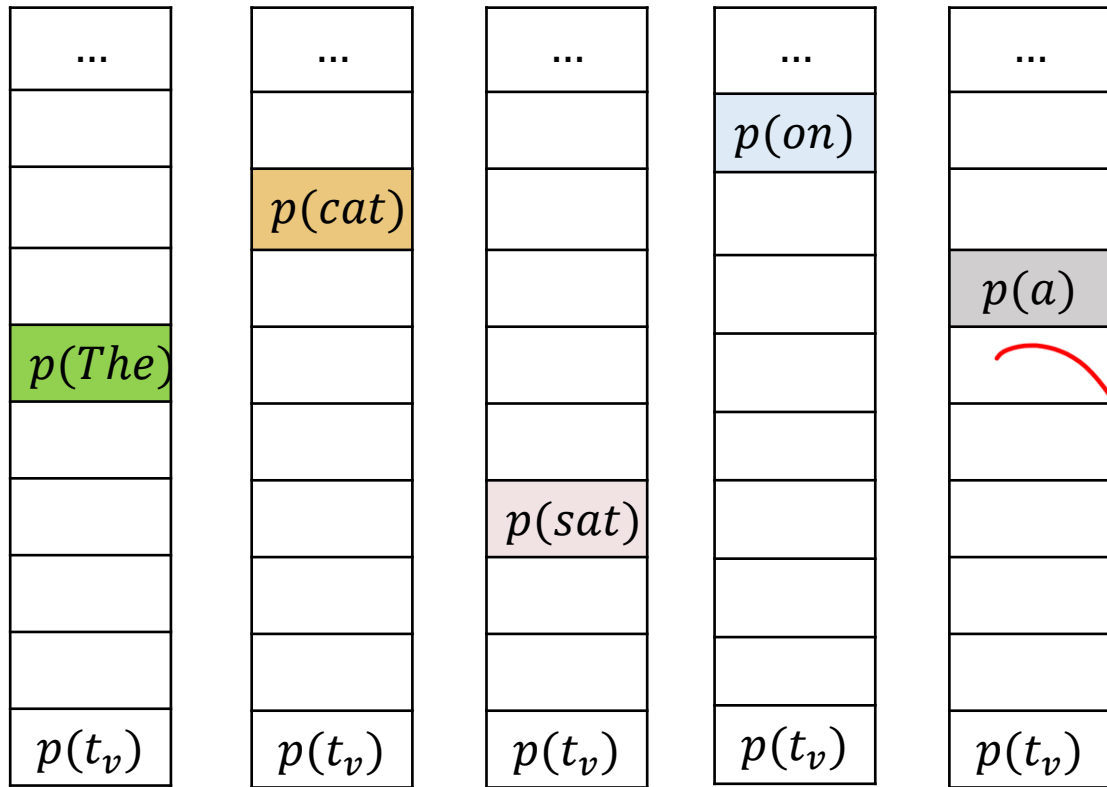


Inference through an LLM

Fwd. pass (#2) to get distribution at step 4

Transformer based LLM (θ)





Inference through an LLM

Fill at step 5

Transformer based LLM (θ)

<s>	The	cat	sat	on	a		
0	1	2	3	4	5	6	7



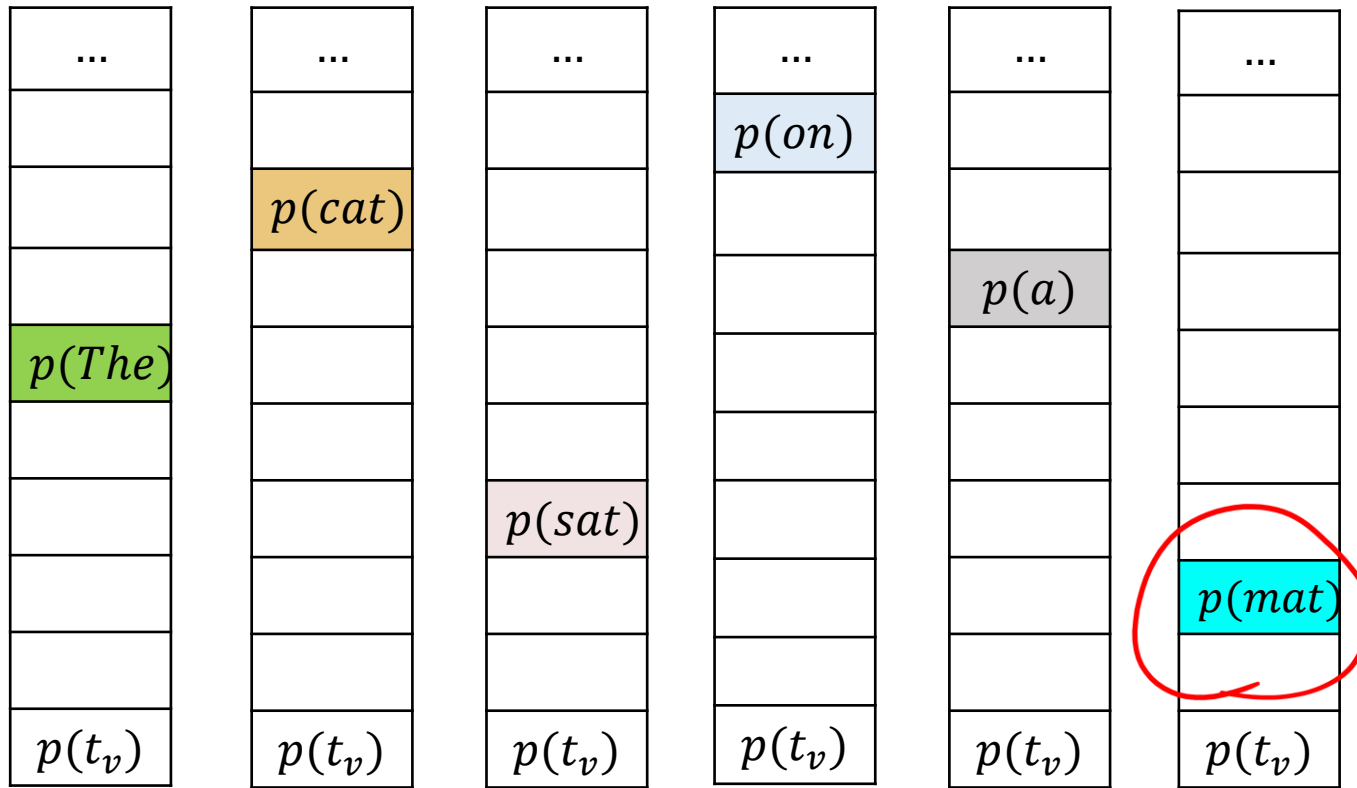
Inference through an LLM

Fwd. pass again (#3)

Transformer based LLM (θ)

<s>	The	cat	sat	on	a		
0	1	2	3	4	5	6	7





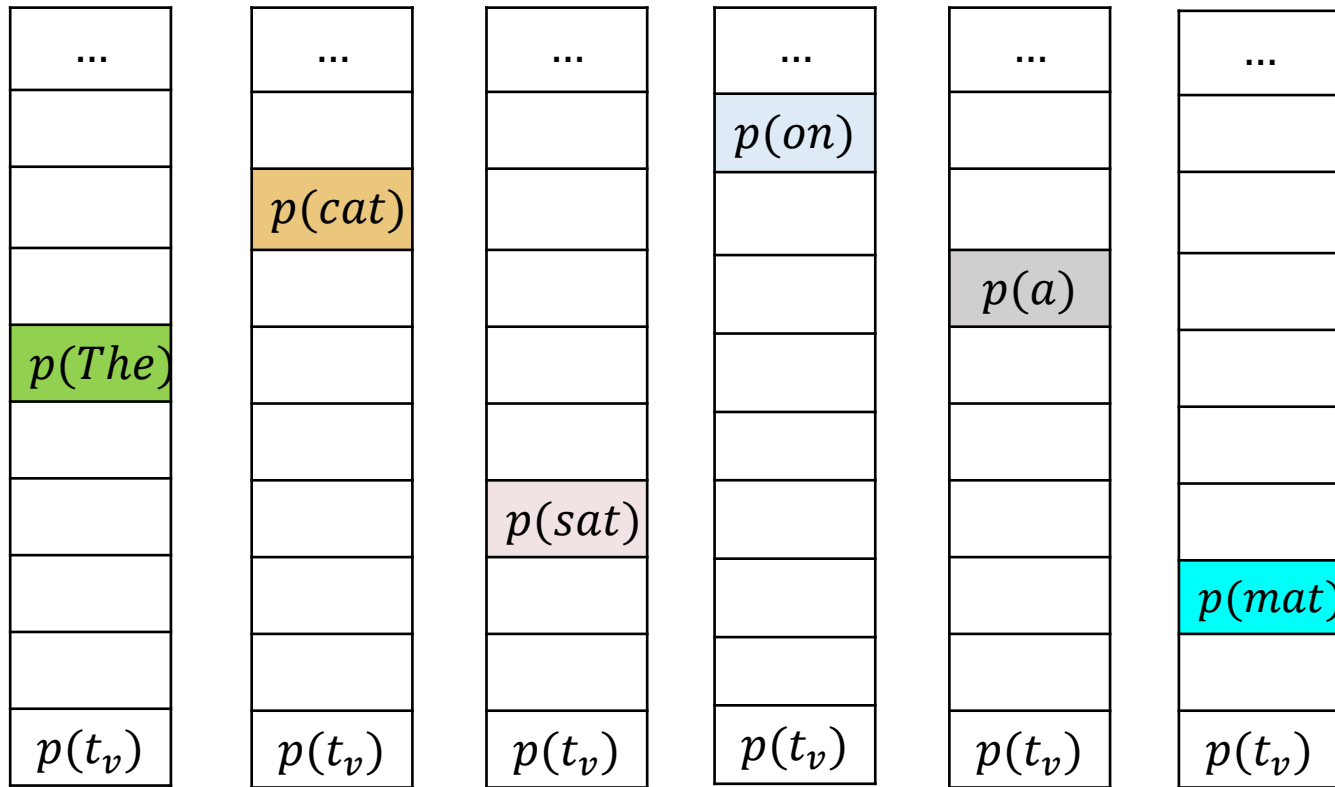
Inference through an LLM

Fwd. pass again (#3)

Transformer based LLM (θ)

<s>	The	cat	sat	on	a		
0	1	2	3	4	5	6	7





Inference through an LLM

Fill at step 6

Transformer based LLM (θ)

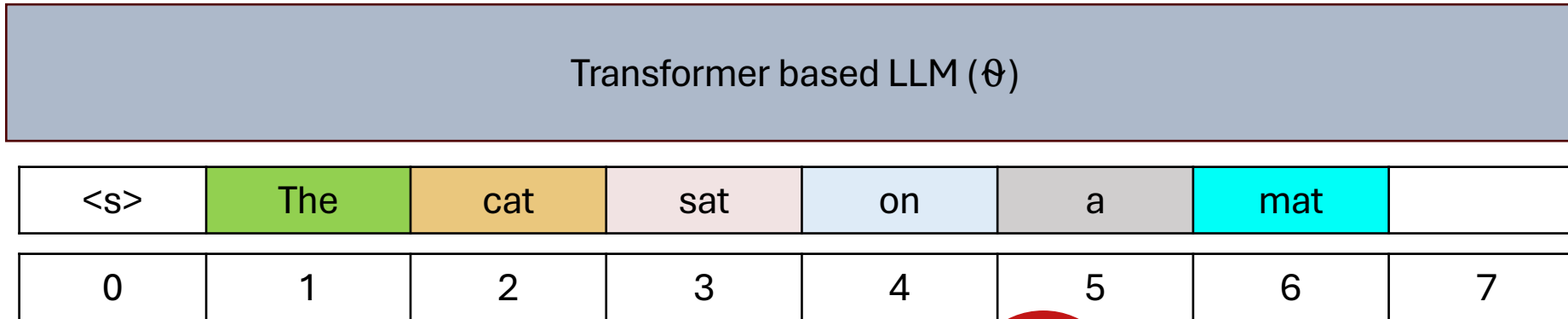
<s> The cat sat on a mat

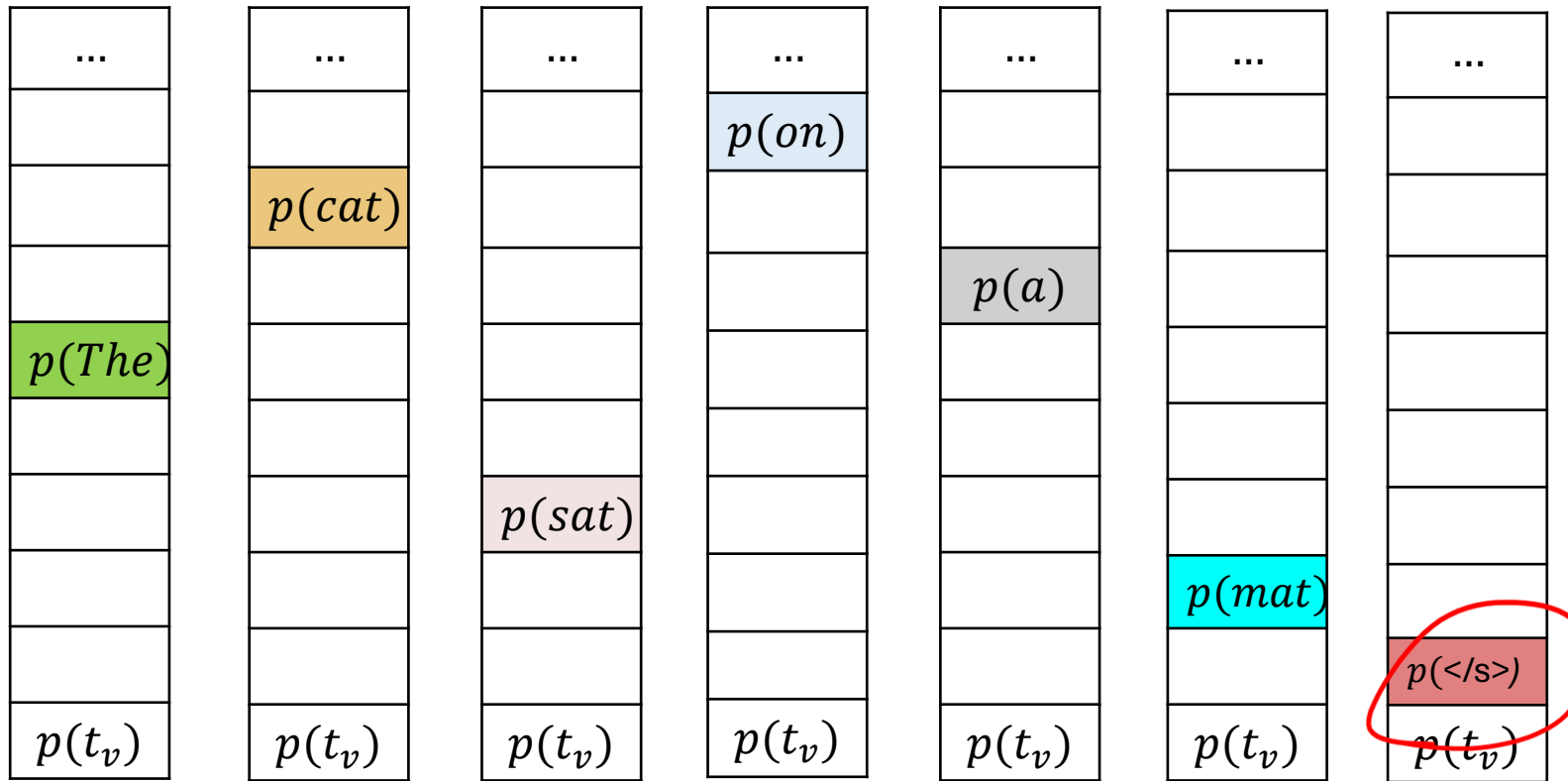
0 1 2 3 4 5 6 7



Inference through an LLM

Fwd. pass again (#4)





Inference through an LLM

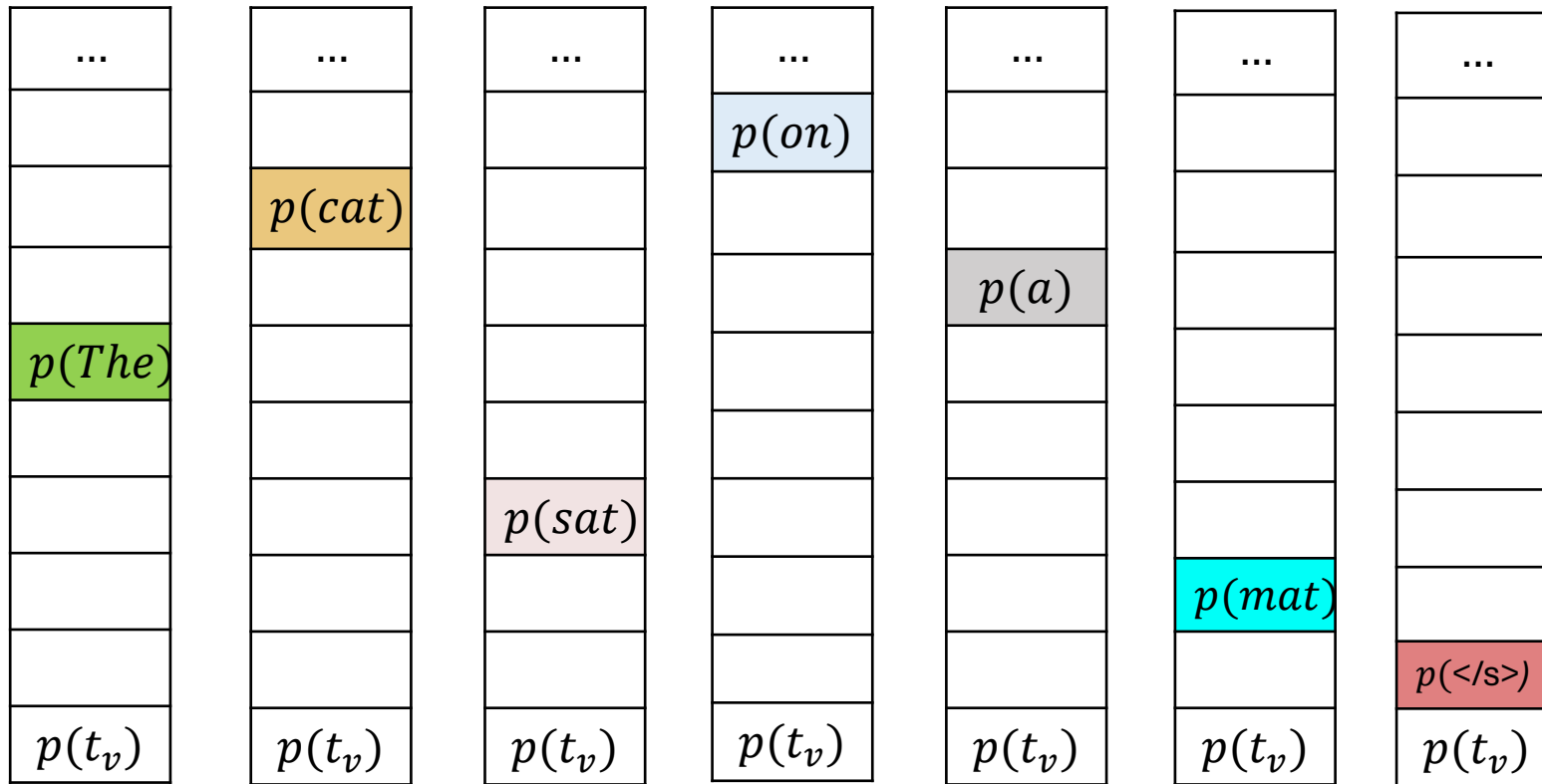
Fwd. pass again (#4)

Transformer based LLM (θ)

<s> The cat sat on a mat

0 1 2 3 4 5 6 7





Inference through an LLM

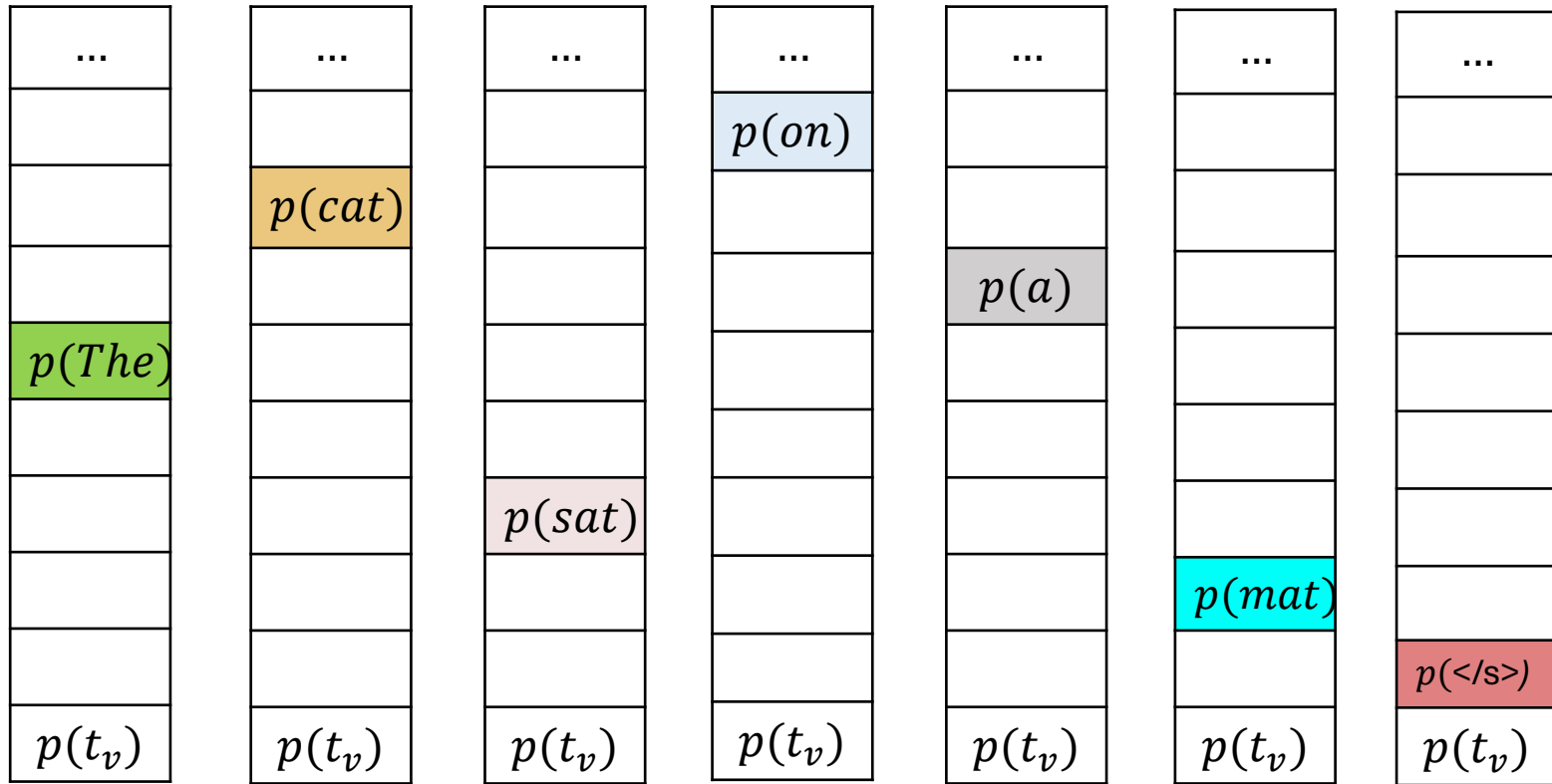
Stop at end of seq. token:
</s>

Transformer based LLM (θ)

<s> The cat sat on a mat </s>

0 1 2 3 4 5 6 7





Inference through an LLM

Fwd Passes: 4
#Tokens: 4

Transformer based LLM (θ)

<s>	The	cat	sat	on	a	mat	</s>
0	1	2	3	4	5	6	7



Inference through an LLM

- ❑ 4 forward passes for 4 tokens
- ❑ Not feasible at production scale
- ❑ Let us revisit forward pass through and see if we can optimize
- ❑ We will focus on attention layer as that is the bottleneck

Fwd Passes: 4
#Tokens: 4

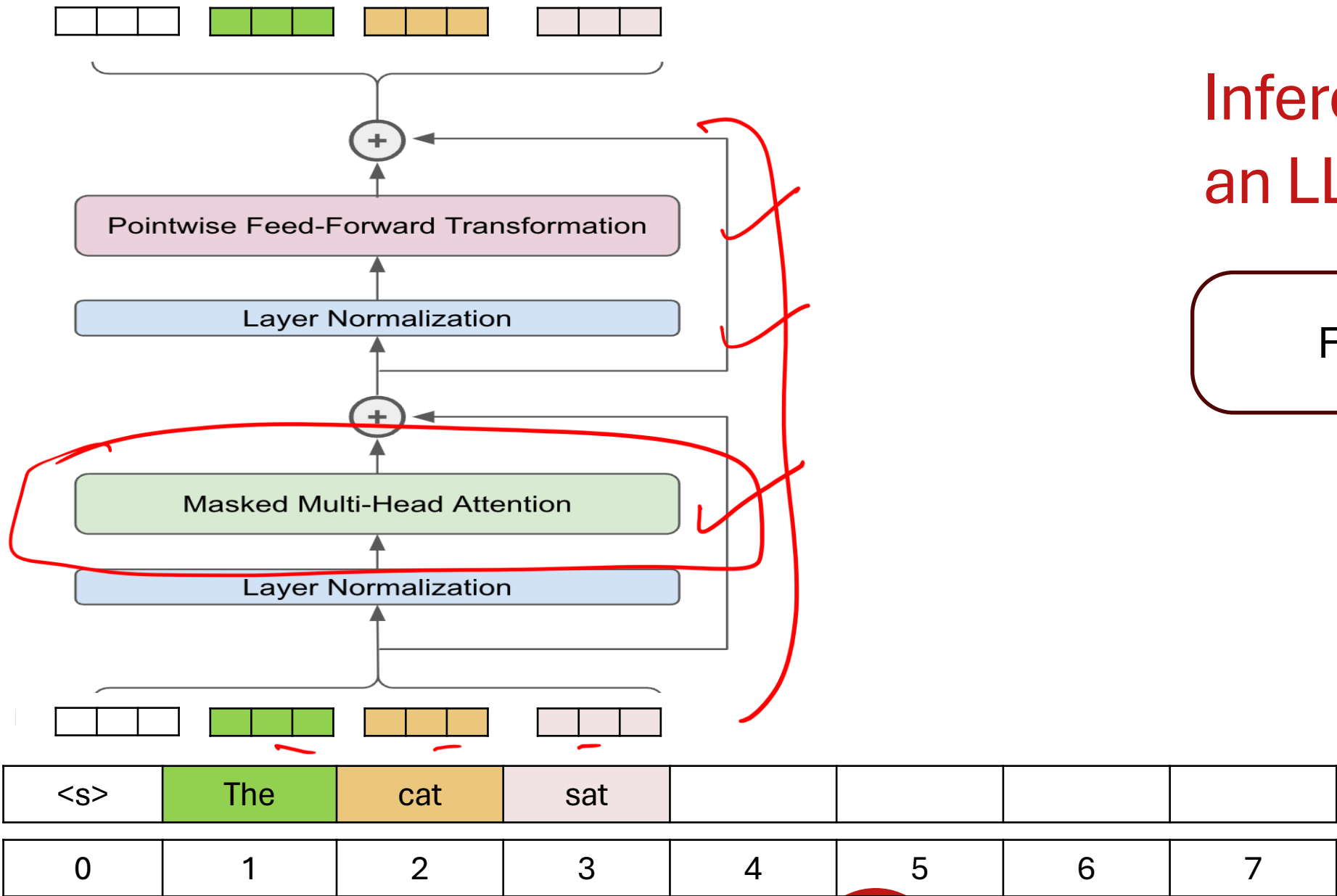
Transformer based LLM (θ)

<s>	The	cat	sat	on	a	mat	</s>
0	1	2	3	4	5	6	7



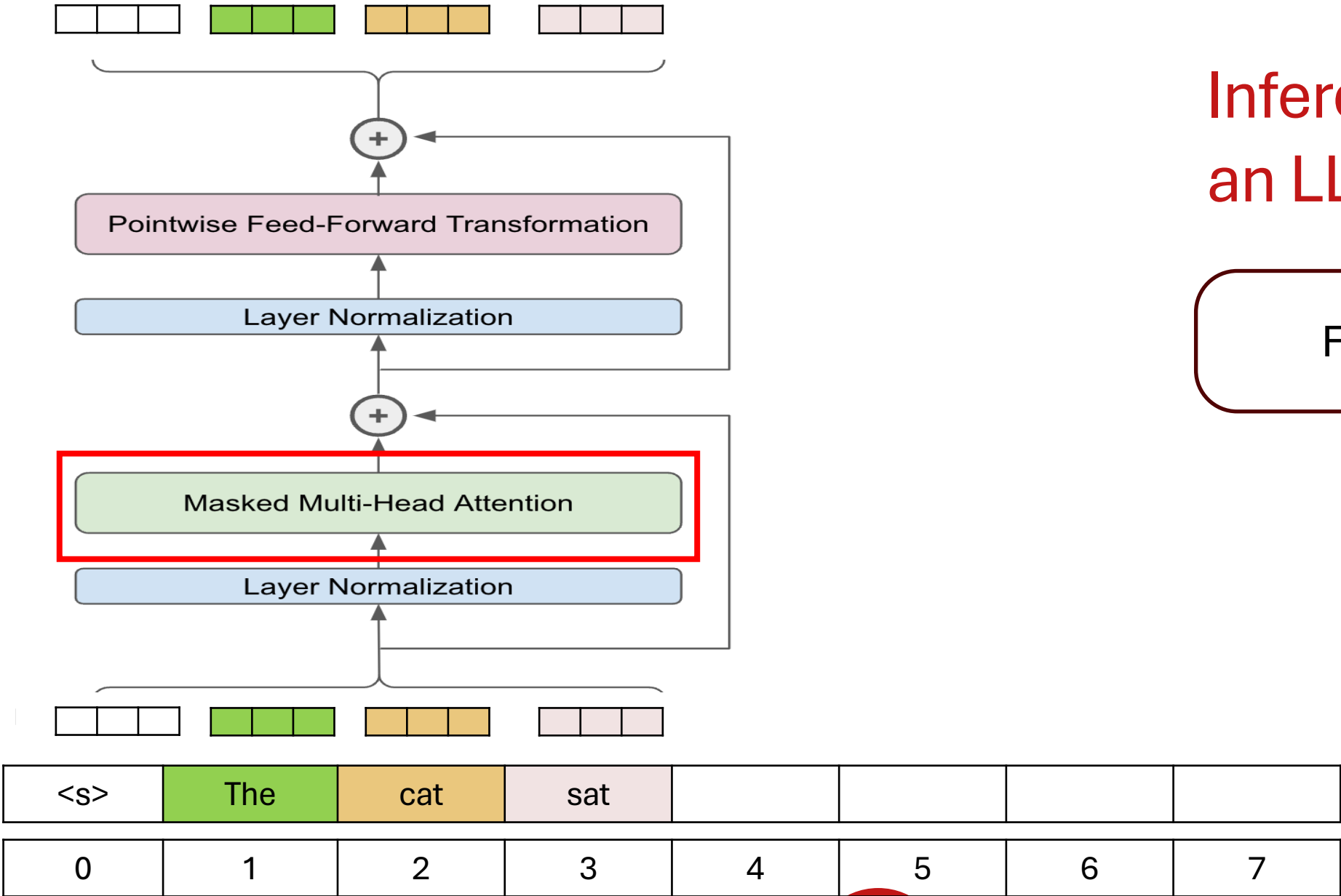
Inference through an LLM

Forward Pass #1



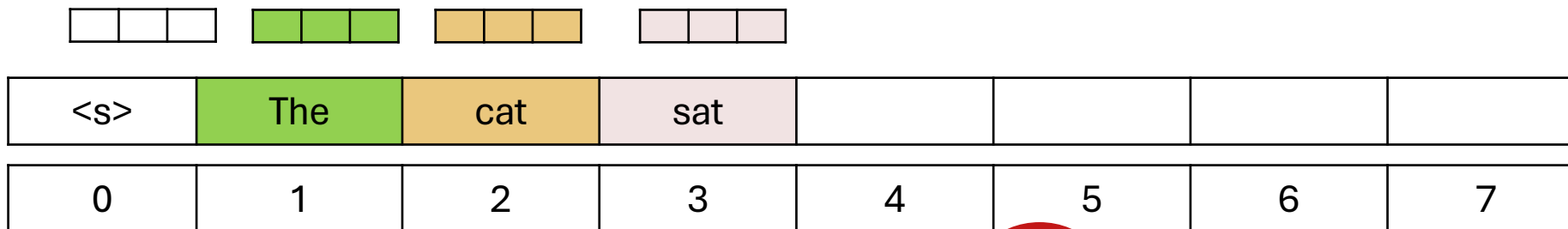
Inference through an LLM

Forward Pass #1



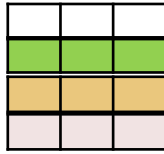
Inference through an LLM

Forward Pass #1



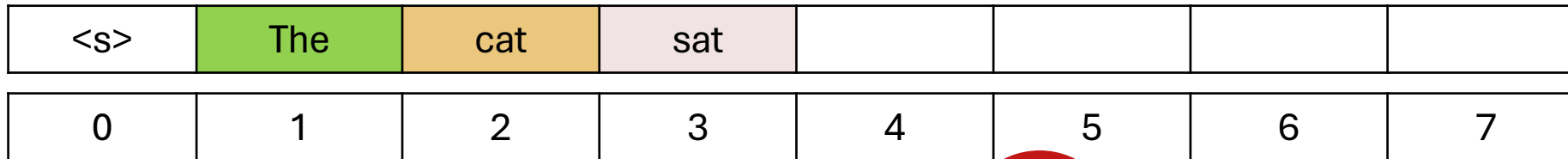
Content credits: <https://cameronwolfe.substack.com/p/decoder-only-transformers-the-workhorse>





Inference through an LLM

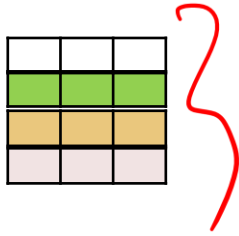
Forward Pass #1



Content credits: <https://cameronwolfe.substack.com/p/decoder-only-transformers-the-workhorse>



W_Q



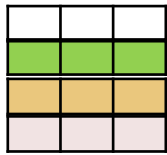
=



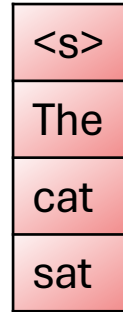
$Q: 4 \times d \text{ dim.}$

Inference through an LLM

W_K



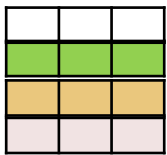
=



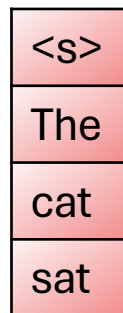
$K: 4 \times d \text{ dim.}$

Forward Pass #1

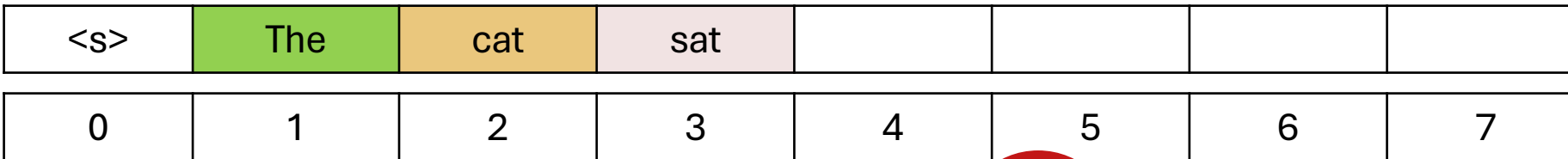
W_V



=



$V: 4 \times d \text{ dim.}$



Content credits: <https://cameronwolfe.substack.com/p/decoder-only-transformers-the-workhorse>



Inference through an LLM

Q : $4 \times d$ dim.

<s>
The
cat
sat

V : $4 \times d$ dim.

<s>
The
cat
sat

Forward Pass #1

<s>	The	cat	sat
-----	-----	-----	-----

K^T : $d \times 4$ dim.

<s>	The	cat	sat				
0	1	2	3	4	5	6	7



Inference through an LLM

Forward Pass #1

Q: $4 \times d$ dim.

<s>
The
cat
sat

A: 4×4 dim.

$$A = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$$

V: $4 \times d$ dim.

<s>
The
cat
sat

<s>	The	cat	sat
-----	-----	-----	-----

K^T: $d \times 4$ dim.

<s>	The	cat	sat				
0	1	2	3	4	5	6	7



Inference through an LLM

Q: $4 \times d$ dim.

A: 4×4 dim.

V: $4 \times d$ dim.

<s>
The
cat
sat

1			
0.2	0.8		
0.1	0.3	0.6	
0.01	0.19	0.3	0.5

<s>
The
cat
sat

<s>	The	cat	sat
-----	-----	-----	-----

K^T: $d \times 4$ dim.

<s>	The	cat	sat				
0	1	2	3	4	5	6	7

Forward Pass #1



Inference through an LLM

Forward Pass #1

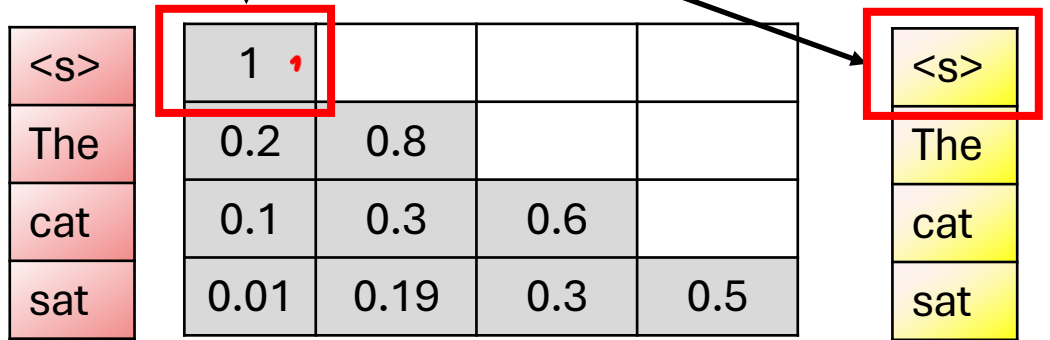
$$\frac{\exp(q_0 k_0^T) v_0}{S_0}$$



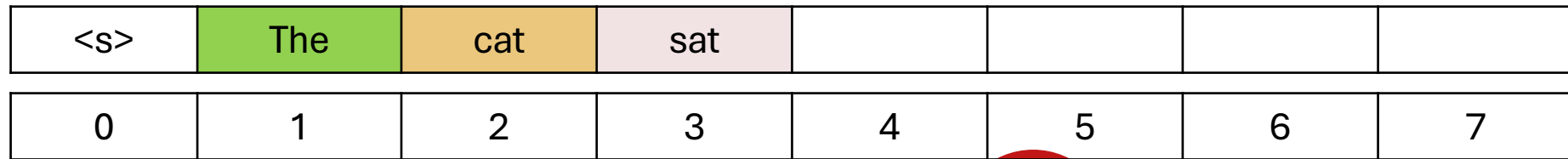
Q: $4 \times d$ dim.

A: 4×4 dim.

V: $4 \times d$ dim.



K^T : $d \times 4$ dim.



Inference through an LLM

Forward Pass #1

$$\frac{\exp(q_1 k_0^T) v_0}{S_1} + \frac{\exp(q_1 k_1^T) v_1}{S_1}$$

<s> The cat sat

Q: 4 x d dim.

A: 4 x 4 dim.

V: 4 x d dim.

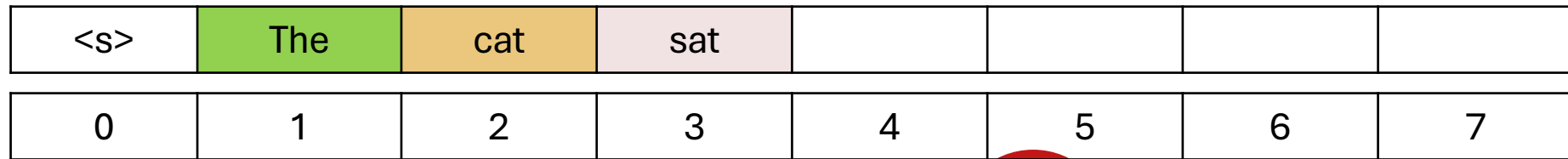
<s>
The
cat
sat

	1			
The	0.2	0.8		
cat	0.1	0.3	0.6	
sat	0.01	0.19	0.3	0.5

<s>
The
cat
sat

<s> The cat sat

K^T: d x 4 dim.



Inference through an LLM

Forward Pass #1

$$\frac{\exp(q_2 k_0^T) v_0}{S_2} + \frac{\exp(q_2 k_1^T) v_1}{S_2} + \frac{\exp(q_2 k_2^T) v_2}{S_2}$$

<s> The cat sat

Q: 4 x d dim.

A: 4 x 4 dim.

V: 4 x d dim.

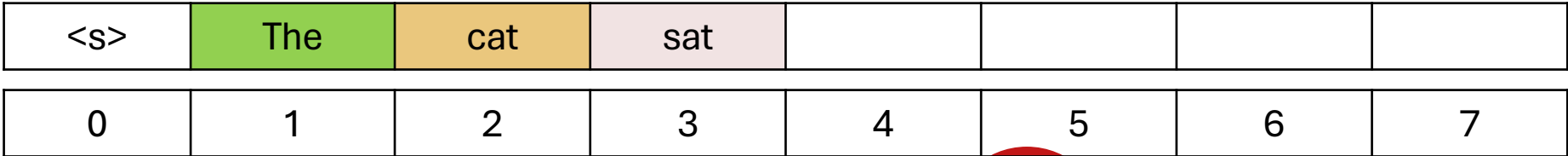
<s>
The
cat
sat

1			
0.2	0.8		
0.1	0.3	0.6	
0.01	0.19	0.3	0.5

<s>
The
cat
sat

<s> The cat sat

K^T: d x 4 dim.



Inference through an LLM

$$\frac{\exp(q_3 k_0^T) v_0}{S_3} + \frac{\exp(q_3 k_1^T) v_1}{S_3} + \frac{\exp(q_3 k_2^T) v_2}{S_3} + \frac{\exp(q_3 k_3^T) v_3}{S_3}$$

<s> The cat sat

Q: 4 x d dim.

A: 4 x 4 dim.

V: 4 x d dim.

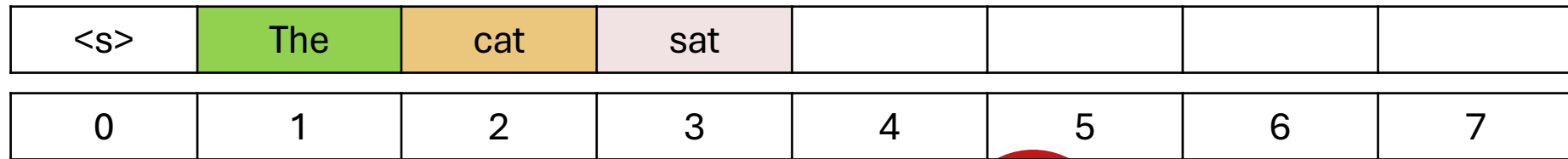
<s>
The
cat
sat

1			
0.2	0.8		
0.1	0.3	0.6	
0.01	0.19	0.3	0.5

<s>
The
cat
sat

<s> The cat sat

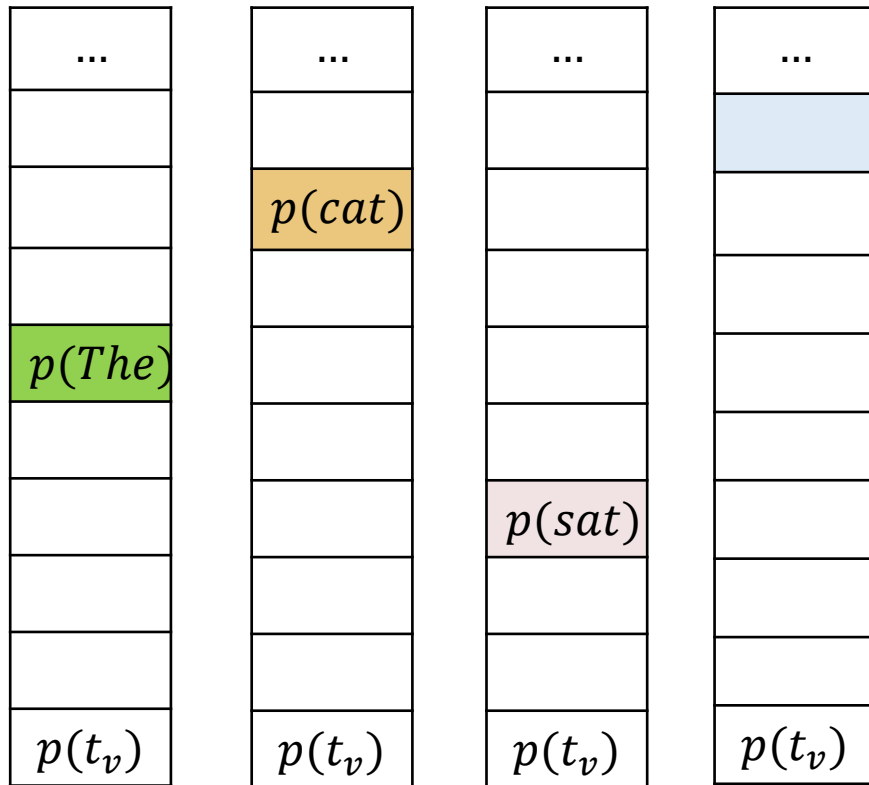
K^T: d x 4 dim.



Forward Pass #1

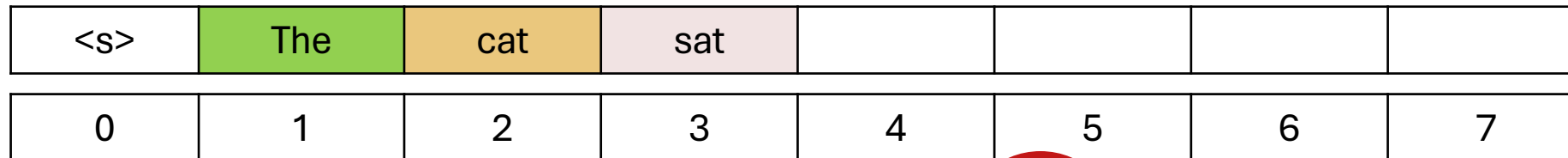


Inference through an LLM

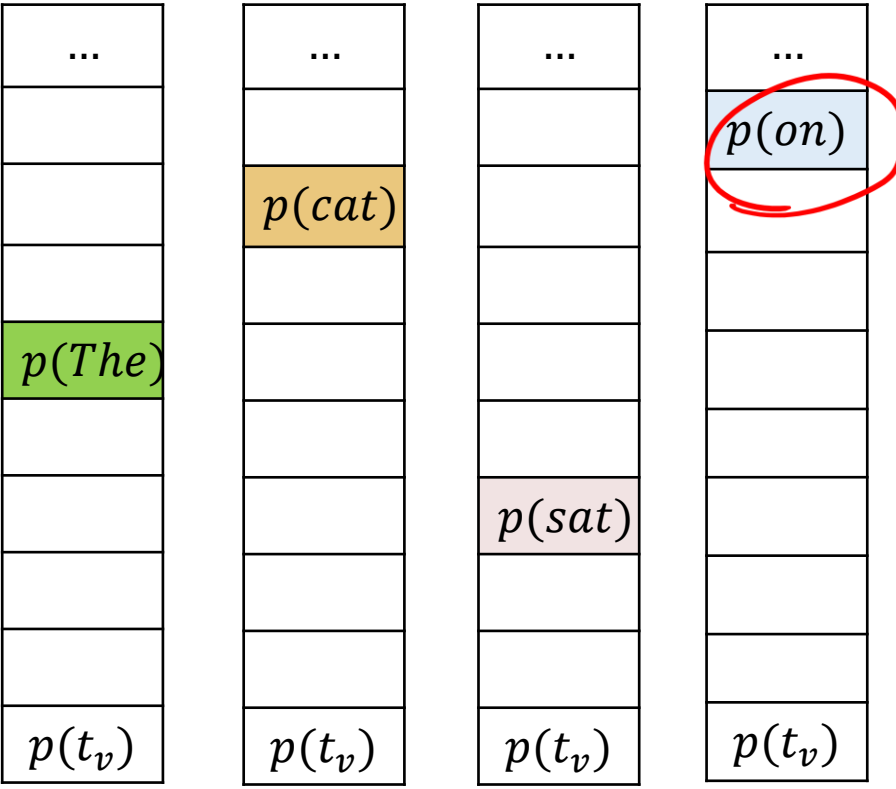


- Emb. of `sat` at the last layer
- Pass through classifier to get distribution over tokens
- Pick the token having max. probability at step 3

Transformer based LLM (θ)

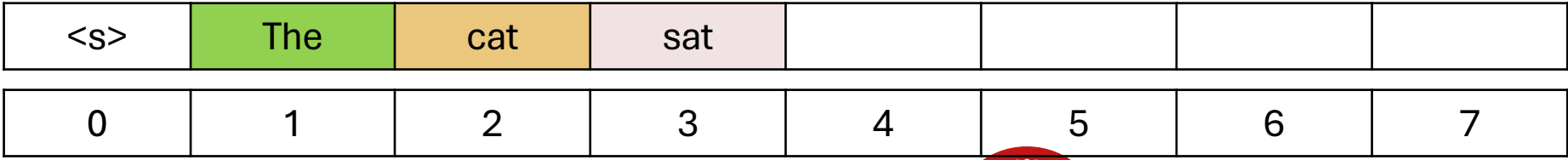


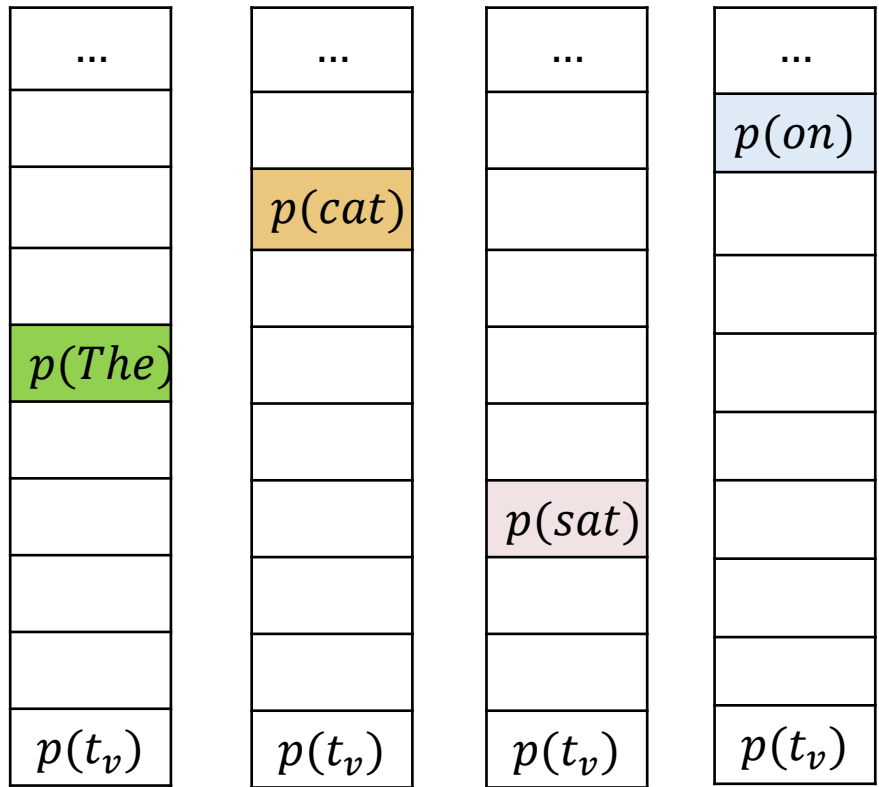
Inference through an LLM



- Emb. of `sat` at the last layer
- Pass through classifier to get distribution over tokens
- Pick the token having max. probability at step 3

Transformer based LLM (θ)

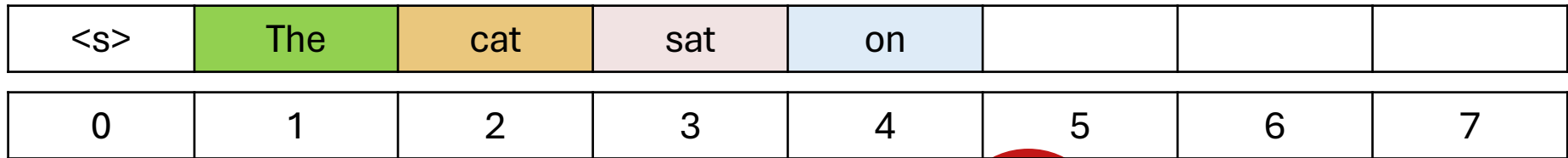




Inference through an LLM

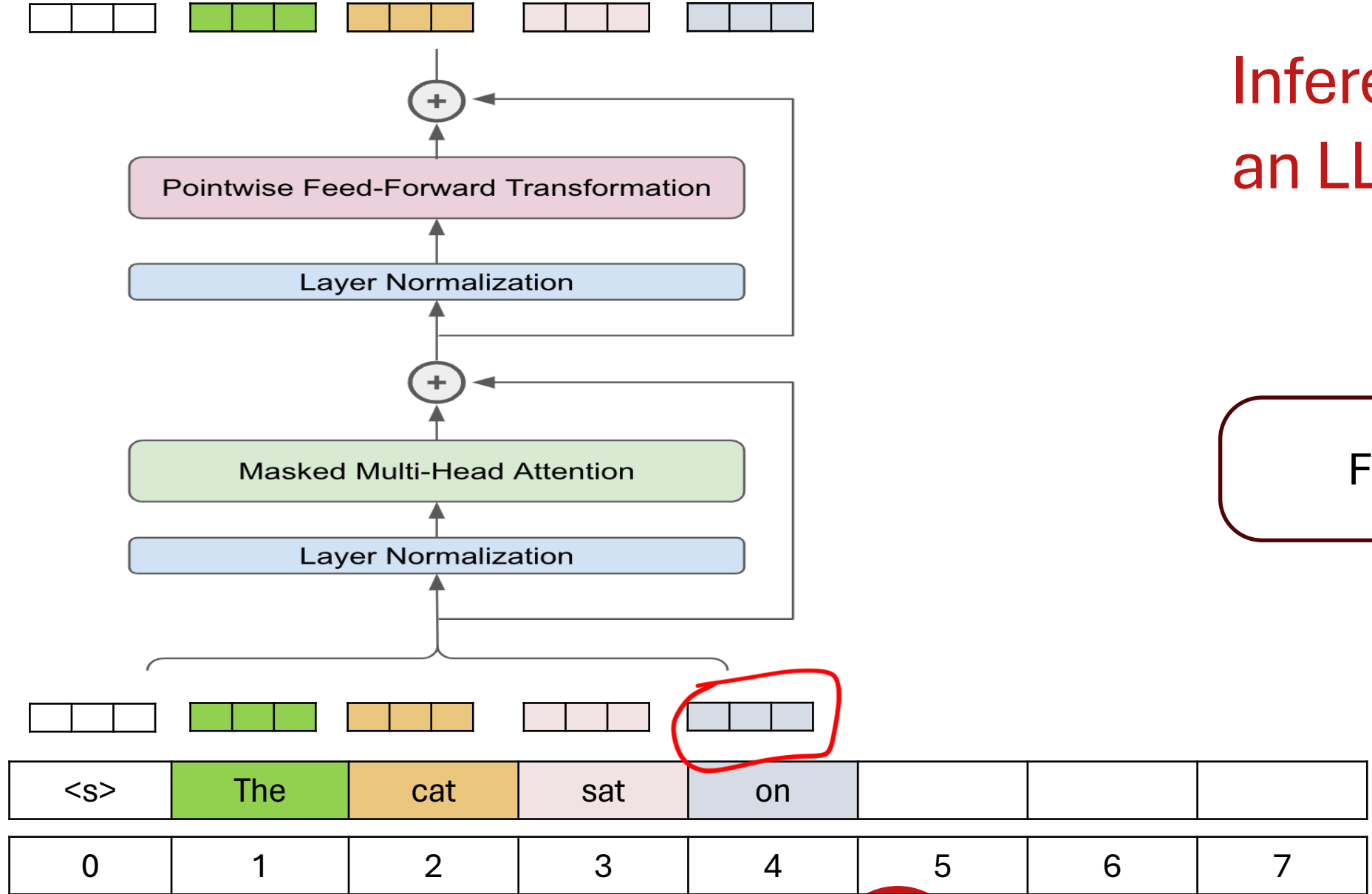
Fill at step 4

Transformer based LLM (θ)



Inference through an LLM

Forward Pass #2

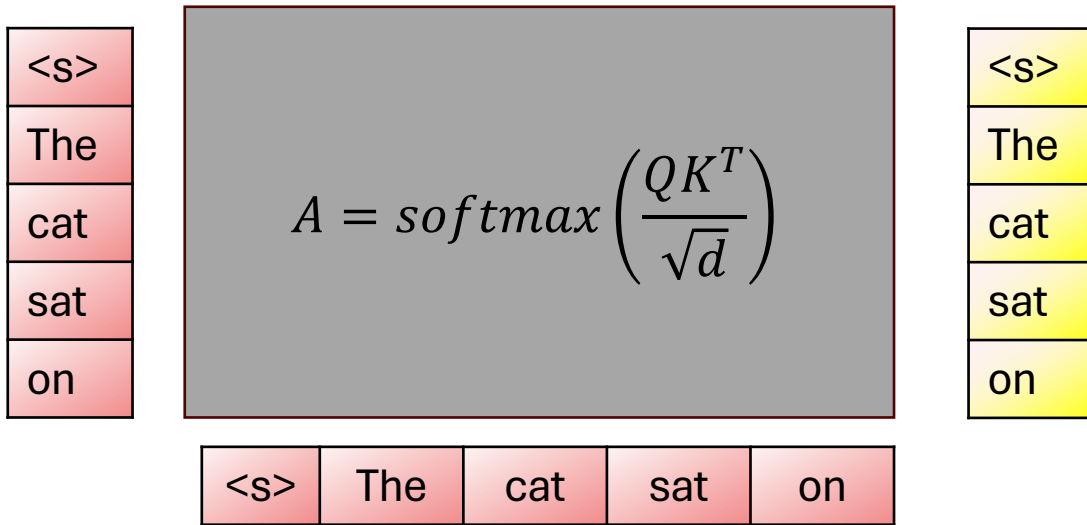


Content credits: <https://cameronwolfe.substack.com/p/decoder-only-transformers-the-workhorse>



Inference through an LLM

Q: $5 \times d$ dim. **A:** 5×5 dim. **V:** $5 \times d$ dim.



K^T: $d \times 5$ dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7

- A lot of computation already done in Fwd. pass #1

Forward Pass #2



Inference through an LLM

Q: $5 \times d$ dim.

<s>
The
cat
sat
on

A: 5×5 dim.

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	

V: $5 \times d$ dim.

<s>
The
cat
sat
on

- Attention matrix already computed in #1

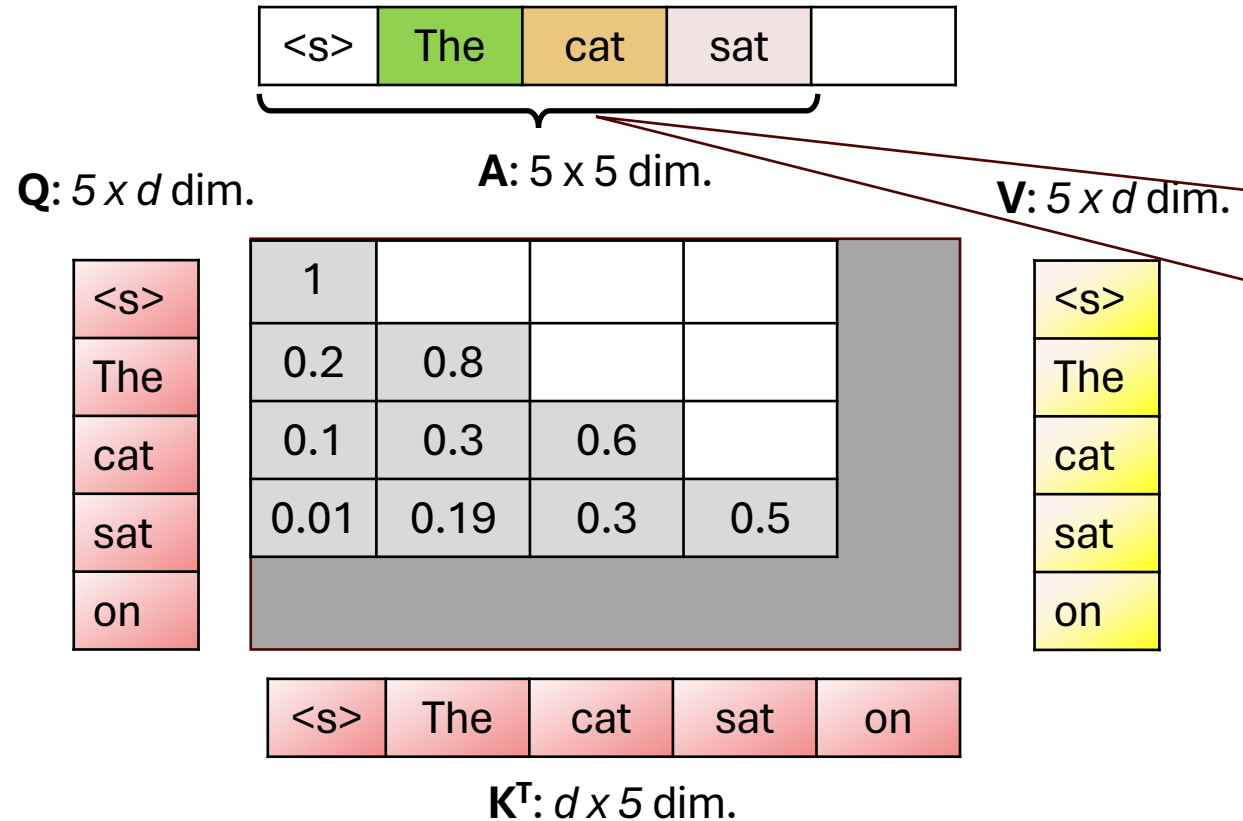
<s>	The	cat	sat	on
-----	-----	-----	-----	----

K^T: $d \times 5$ dim.

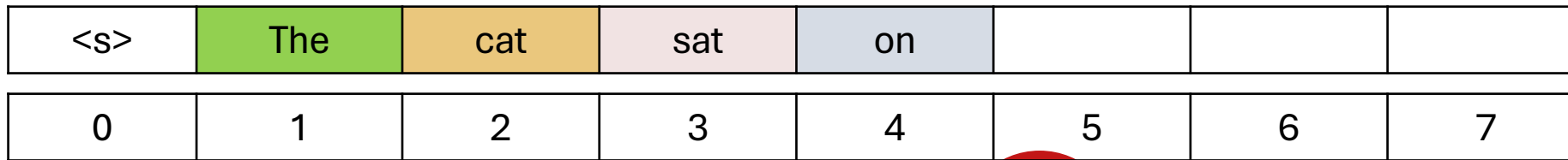
<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM



- Attention matrix already computed in #1
- Output embed. already computed in #1



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1			
0.2	0.8		
0.1	0.3	0.6	
0.01	0.19	0.3	0.5

<s>
The
cat
sat
on

<s> The cat sat on

K^T : $d \times 5$ dim.

- Attention matrix already computed in #1
- Output embed. already computed in #1
- Keys and Values already computed in #1

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1			
0.2	0.8		
0.1	0.3	0.6	
0.01	0.19	0.3	0.5

<s>
The
cat
sat
on

<s> The cat sat on

K^T : $d \times 5$ dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7

- Attention matrix already computed in #1
- Output embed. already computed in #1
- Keys and Values already computed in #1
- Queries not required in #2



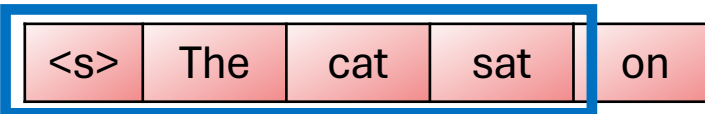
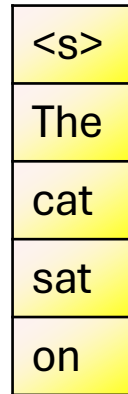
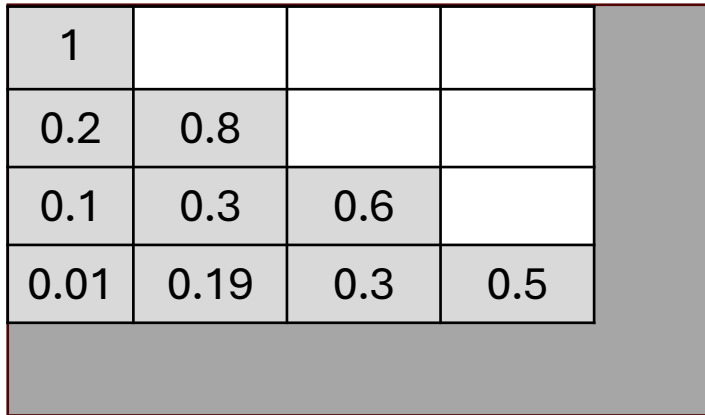
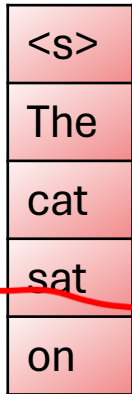
Inference through an LLM



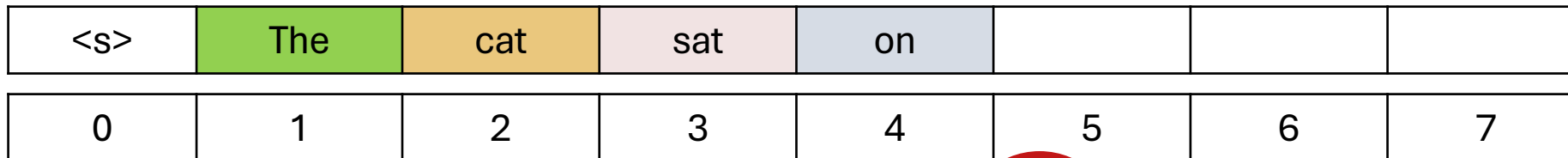
Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.



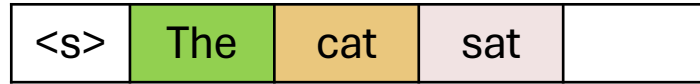
K^T : $d \times 5$ dim.



- Cache the already computed matrices



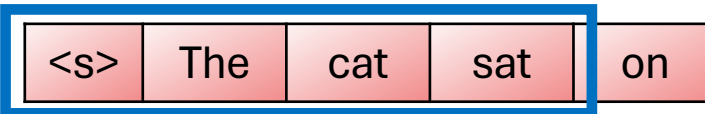
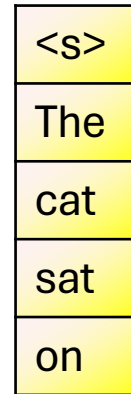
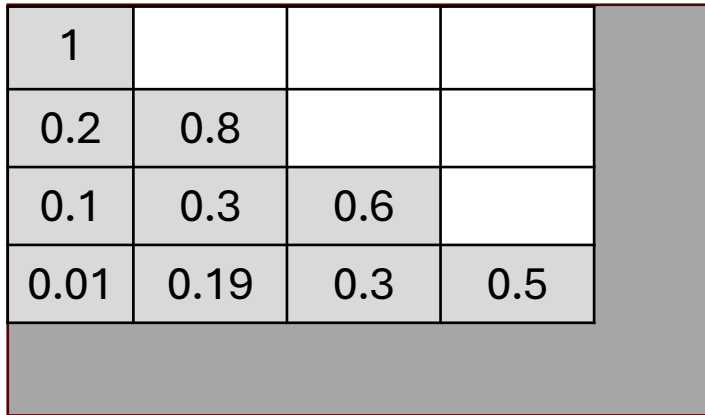
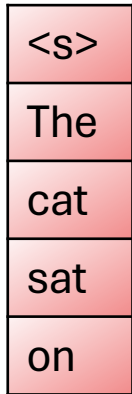
Inference through an LLM



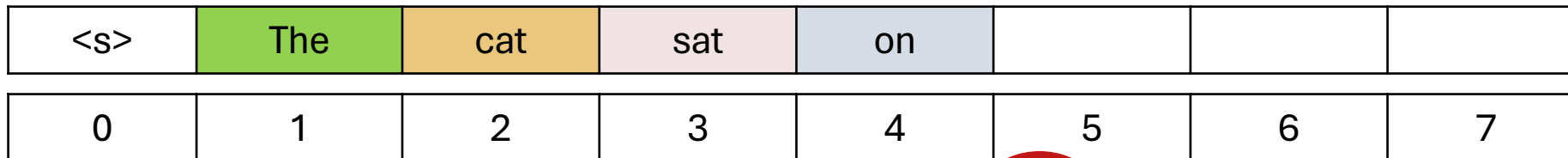
Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.



K^T : $d \times 5$ dim.



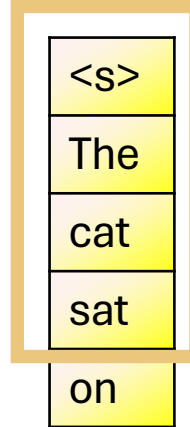
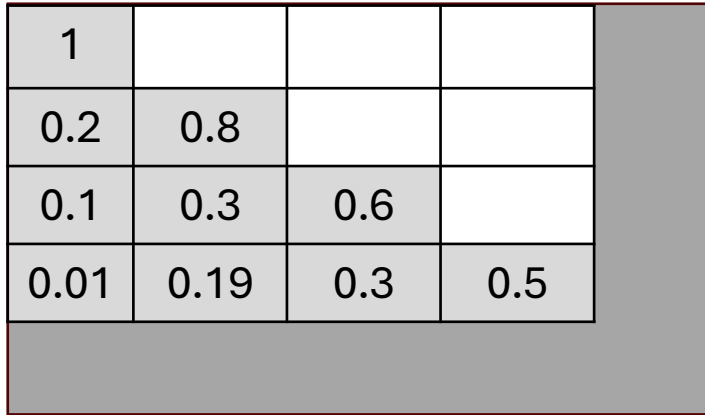
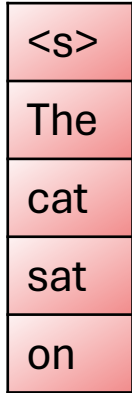
Inference through an LLM



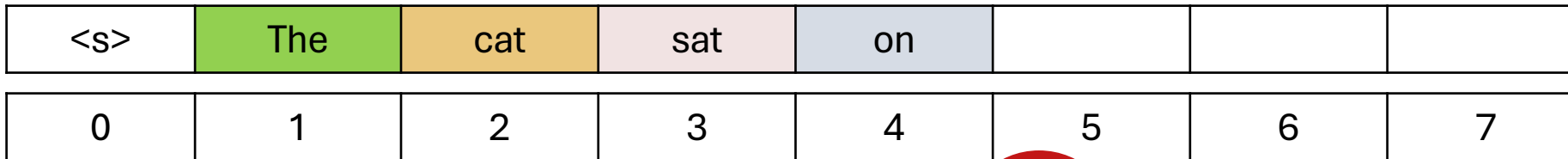
Q: $5 \times d$ dim.

A: 5×5 dim.

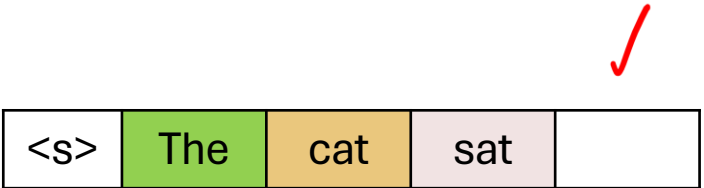
V: $5 \times d$ dim.



K^T : $d \times 5$ dim.



Inference through an LLM



Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	

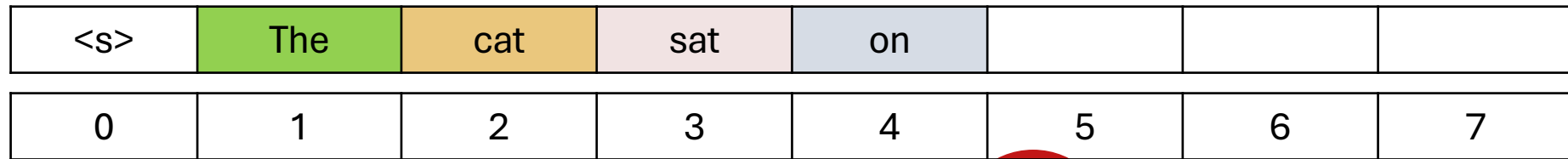
<s>
The
cat
sat
on

<s>	The	cat	sat	on
-----	-----	-----	-----	----

K^T : $d \times 5$ dim.

K cache
<s>
The
cat
sat

V cache
<s>
The
cat
sat



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

K cache

<s>
The
cat
sat

V cache

<s>
The
cat
sat

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	

<s>
The
cat
sat
on

<s> The cat sat on

K^T : $d \times 5$ dim.

Compute Query vector for token on

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	

<s>
The
cat
sat
on

K cache

<s>
The
cat
sat

V cache

<s>
The
cat
sat

<s> The cat sat on

K^T : $d \times 5$ dim.

Read Key vector for <s> token from cache

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
$q_4 k_0^T$				

<s>
The
cat
sat
on

K cache

<s>
The
cat
sat

V cache

<s>
The
cat
sat

<s> The cat sat on

$K^T: d \times 5$ dim.

Dot product to compute attention score b/w query "on" and key "<s>"

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

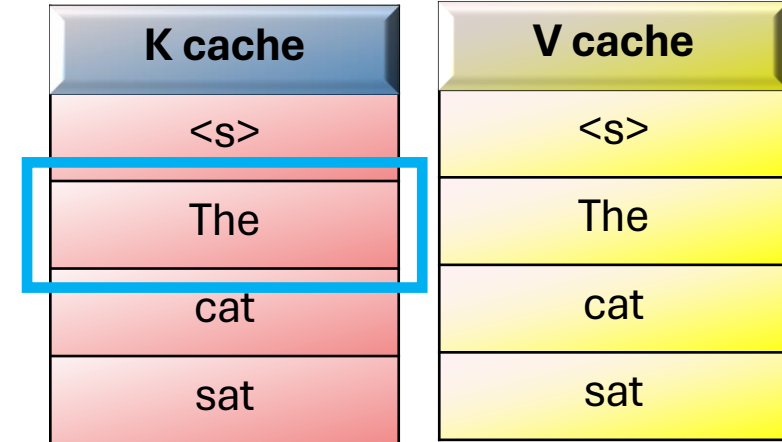
A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
$q_4 k_0^T$	$q_4 k_1^T$			

<s>
The
cat
sat
on



<s> The cat sat on

K^T : $d \times 5$ dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

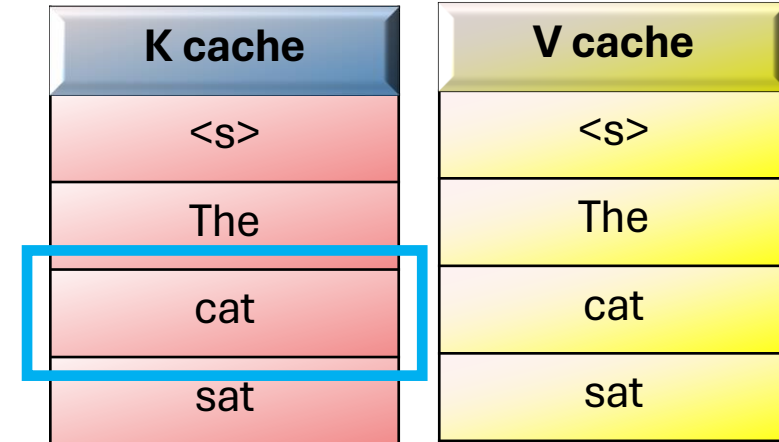
A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
$q_4 k_0^T$	$q_4 k_1^T$	$q_4 k_2^T$		

<s>
The
cat
sat
on



<s> The cat sat on

K^T : $d \times 5$ dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

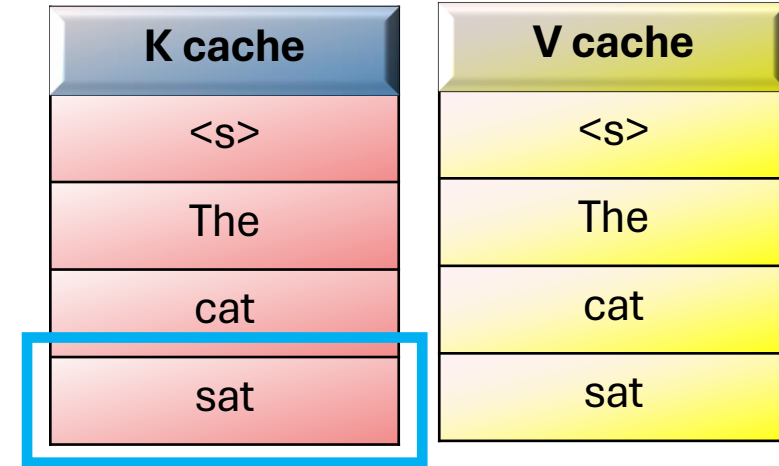
A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
$q_4 k_0^T$	$q_4 k_1^T$	$q_4 k_2^T$	$q_4 k_3^T$	

<s>
The
cat
sat
on



<s> The cat sat on

K^T : $d \times 5$ dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
$q_4 k_0^T$	$q_4 k_1^T$	$q_4 k_2^T$	$q_4 k_3^T$	$q_4 k_4^T$

<s>
The
cat
sat
on

K cache

<s>
The
cat
sat

V cache

<s>
The
cat
sat

<s> The cat sat on

K^T : $d \times 5$ dim.

Compute the Key emb. of token "on"

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
$q_4 k_0^T$	$q_4 k_1^T$	$q_4 k_2^T$	$q_4 k_3^T$	$q_4 k_4^T$

<s>
The
cat
sat
on

K cache

<s>
The
cat
sat
on

V cache

<s>
The
cat
sat

<s> The cat sat on

$K^T: d \times 5$ dim.

Add the Key emb. of token "on" to the cache

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

V: $5 \times d$ dim.

<s>	1				
The	0.2	0.8			
cat	0.1	0.3	0.6		
sat	0.01	0.19	0.3	0.5	
on	0.03	0.07	0.1	0.3	0.4

<s> The cat sat on

K^T : $d \times 5$ dim.

<s>
The
cat
sat
on

K cache
<s>
The
cat
sat
on

V cache
<s>
The
cat
sat

Convert attn. scores to probability

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

<s> The cat sat

Q: $5 \times d$ dim.

A: 5×5 dim.

<s>	1				
The	0.2	0.8			
cat	0.1	0.3	0.6		
sat	0.01	0.19	0.3	0.5	
on	0.03	0.07	0.1	0.3	0.4

<s> The cat sat on

K^T : $d \times 5$ dim.

V: $5 \times d$ dim.

<s>
The
cat
sat
on

K cache

<s>
The
cat
sat
on

V cache

<s>
The
cat
sat

Load Value vectors from V cache

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4}$$

<s>	The	cat	sat	
-----	-----	-----	-----	--

Q: 5 x d dim.

A: 5 x 5 dim.

V: 5 x d dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
0.03	0.07	0.1	0.3	0.4

<s>
The
cat
sat
on

K cache
<s>
The
cat
sat
on

V cache
<s>
me
cat
sat

<s>	The	cat	sat	on
-----	-----	-----	-----	----

K^T: d x 5 dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4} + \frac{(q_4 k_1^T) v_1}{S_4}$$

<s> The cat sat

Q: 5 x d dim.

A: 5 x 5 dim.

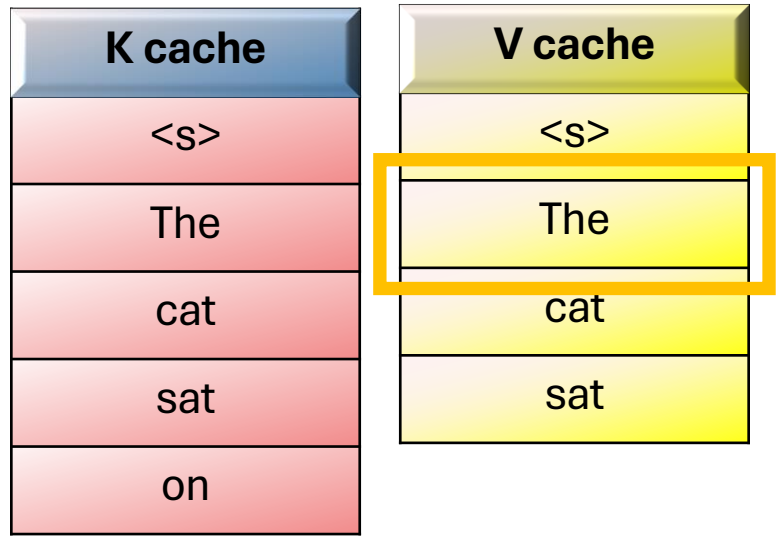
V: 5 x d dim.

<s>	1				
The	0.2	0.8			
cat	0.1	0.3	0.6		
sat	0.01	0.19	0.3	0.5	
on	0.03	0.07	0.1	0.3	0.4

<s> The cat sat on

K^T: d x 5 dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4} + \frac{(q_4 k_1^T) v_1}{S_4} + \frac{(q_4 k_2^T) v_2}{S_4}$$

<s> The cat sat

Q: 5 x d dim.

A: 5 x 5 dim.

V: 5 x d dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
0.03	0.07	0.1	0.3	0.4

<s>
The
cat
sat
on

K cache
<s>
The
cat
sat
on

V cache
<s>
The
cat
sat

<s> The cat sat on

K^T: d x 5 dim.

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4} + \frac{(q_4 k_1^T) v_1}{S_4} + \frac{(q_4 k_2^T) v_2}{S_4} + \frac{(q_4 k_3^T) v_3}{S_4}$$

<s> The cat sat

Q: 5 x d dim.

A: 5 x 5 dim.

V: 5 x d dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
0.03	0.07	0.1	0.3	0.4

<s>
The
cat
sat
on

<s> The cat sat on

K^T: d x 5 dim.

K cache
<s>
The
cat
sat
on

V cache
<s>
The
cat
sat

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4} + \frac{(q_4 k_1^T) v_1}{S_4} + \frac{(q_4 k_2^T) v_2}{S_4} + \frac{(q_4 k_3^T) v_3}{S_4}$$

<s> The cat sat

Q: 5 x d dim.

A: 5 x 5 dim.

V: 5 x d dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
0.03	0.07	0.1	0.3	0.4

<s>
The
cat
sat
on

K cache

<s>
The
cat
sat
on

V cache

<s>
The
cat
sat

<s> The cat sat on

K^T: d x 5 dim.

Compute V emb. of on

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4} + \frac{(q_4 k_1^T) v_1}{S_4} + \frac{(q_4 k_2^T) v_2}{S_4} + \frac{(q_4 k_3^T) v_3}{S_4}$$

<s> The cat sat

Q: 5 x d dim.

A: 5 x 5 dim.

V: 5 x d dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
0.03	0.07	0.1	0.3	0.4

<s>
The
cat
sat
on

K cache
<s>
The
cat
sat
on

V cache
<s>
The
cat
sat
on

<s> The cat sat on

K^T: d x 5 dim.

Add V emb. of on to V-cache

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7



Inference through an LLM

$$\frac{(q_4 k_0^T) v_0}{S_4} + \frac{(q_4 k_1^T) v_1}{S_4} + \frac{(q_4 k_2^T) v_2}{S_4} + \frac{(q_4 k_3^T) v_3}{S_4} + \frac{(q_4 k_4^T) v_4}{S_4}$$

<s> The cat sat on

Q: 5 x d dim.

A: 5 x 5 dim.

V: 5 x d dim.

<s>
The
cat
sat
on

1				
0.2	0.8			
0.1	0.3	0.6		
0.01	0.19	0.3	0.5	
0.03	0.07	0.1	0.3	0.4

<s>
The
cat
sat
on

K cache
<s>
The
cat
sat
on

V cache
<s>
The
cat
sat
on

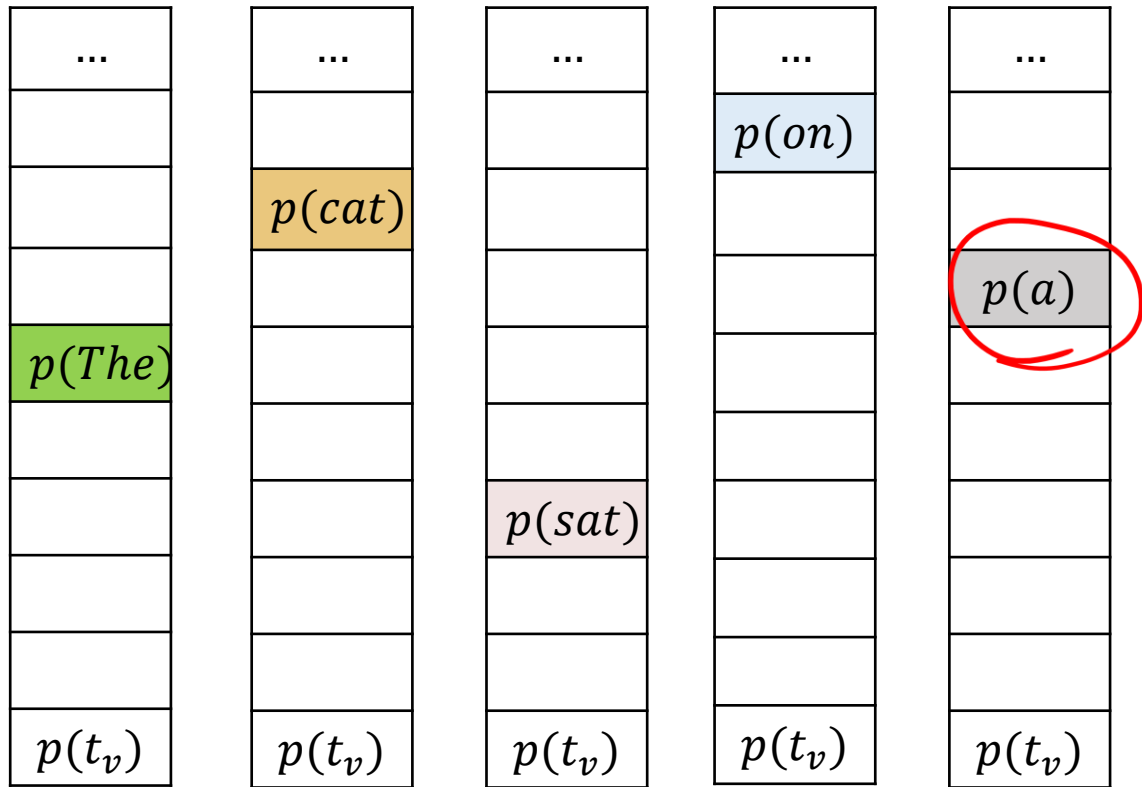
<s> The cat sat on

K^T: d x 5 dim.

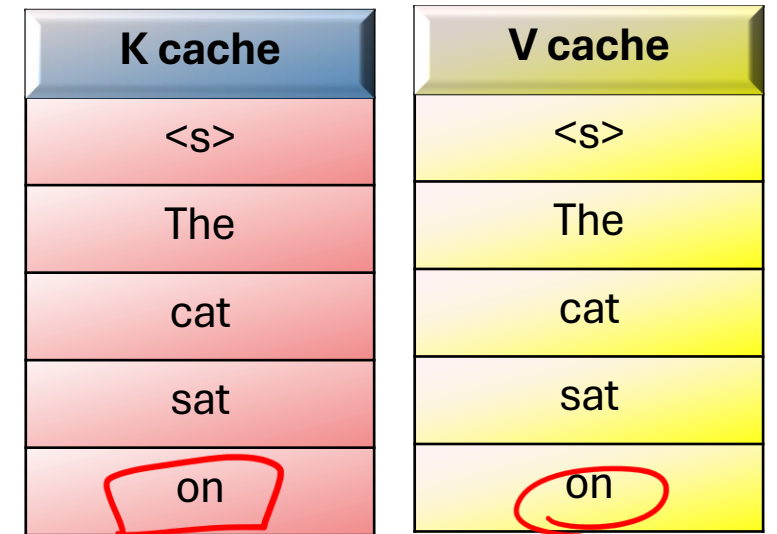
We get output emb. of on

<s>	The	cat	sat	on			
0	1	2	3	4	5	6	7

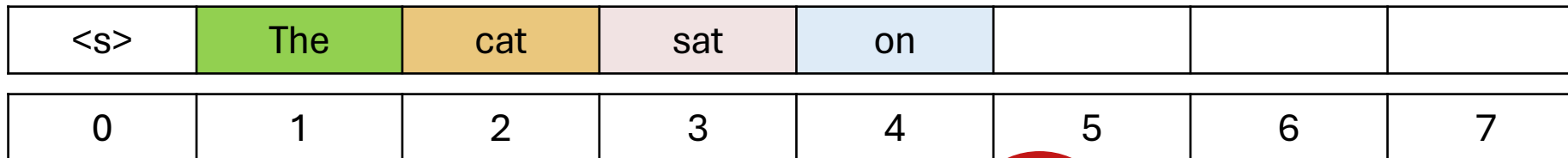




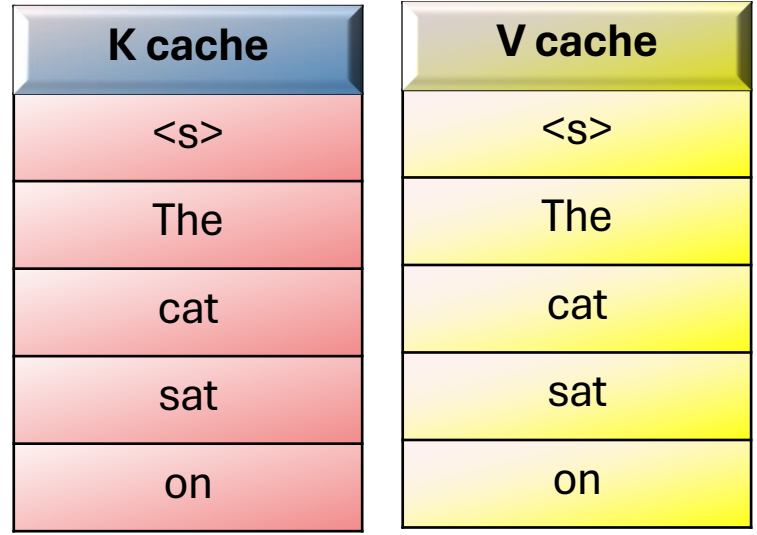
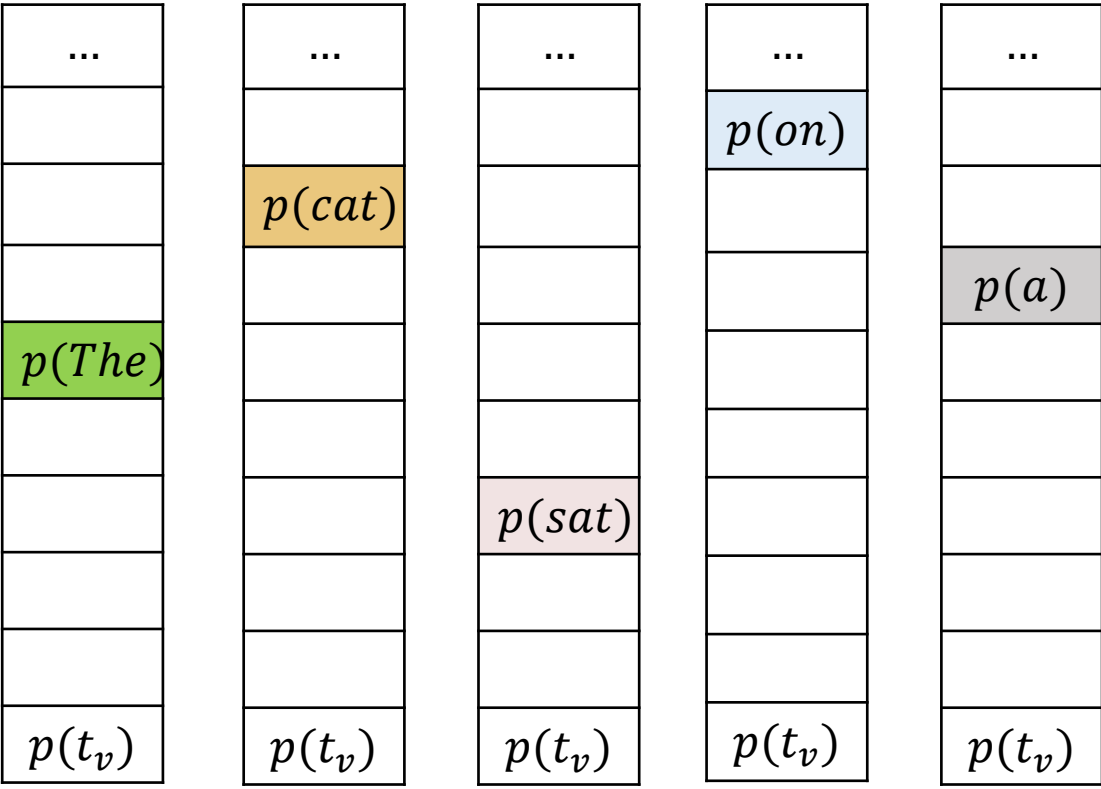
Inference through an LLM



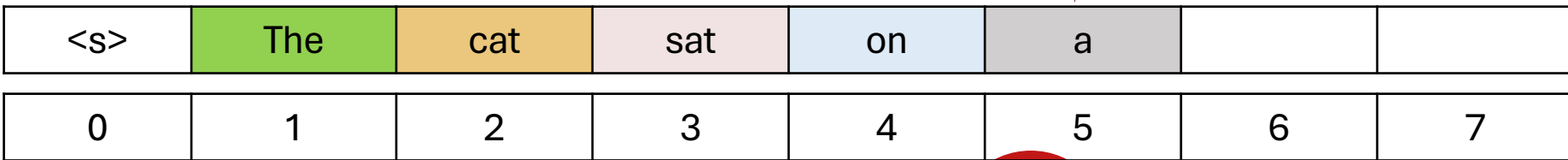
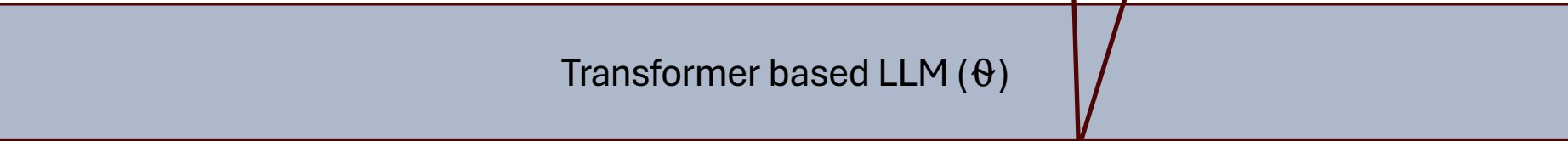
Transformer based LLM (θ)

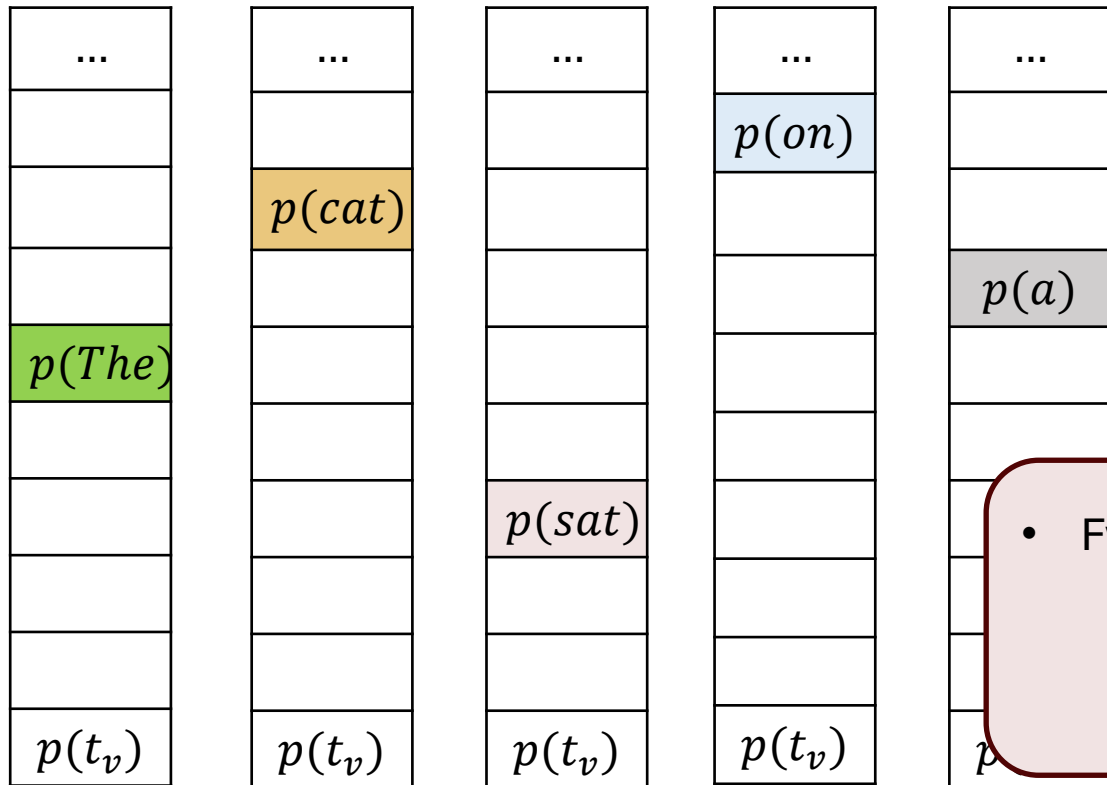


Inference through an LLM



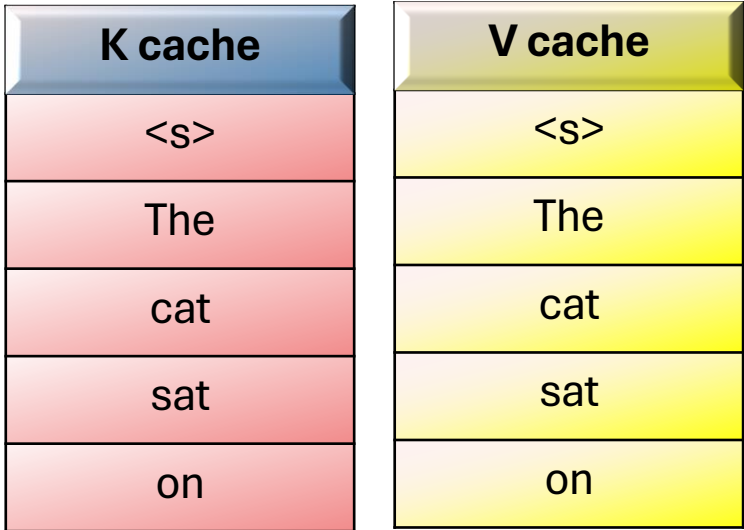
Fill at step 5



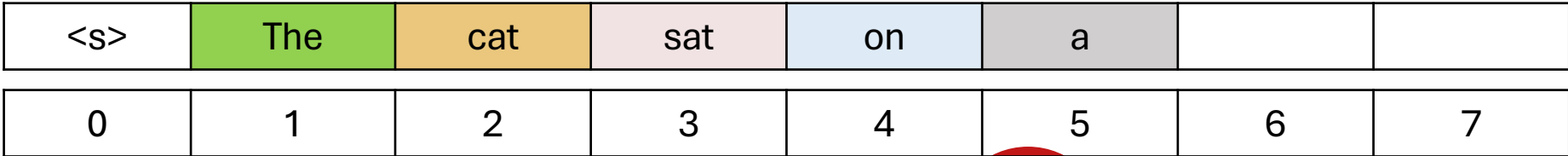


- Fwd. pass again (#3)

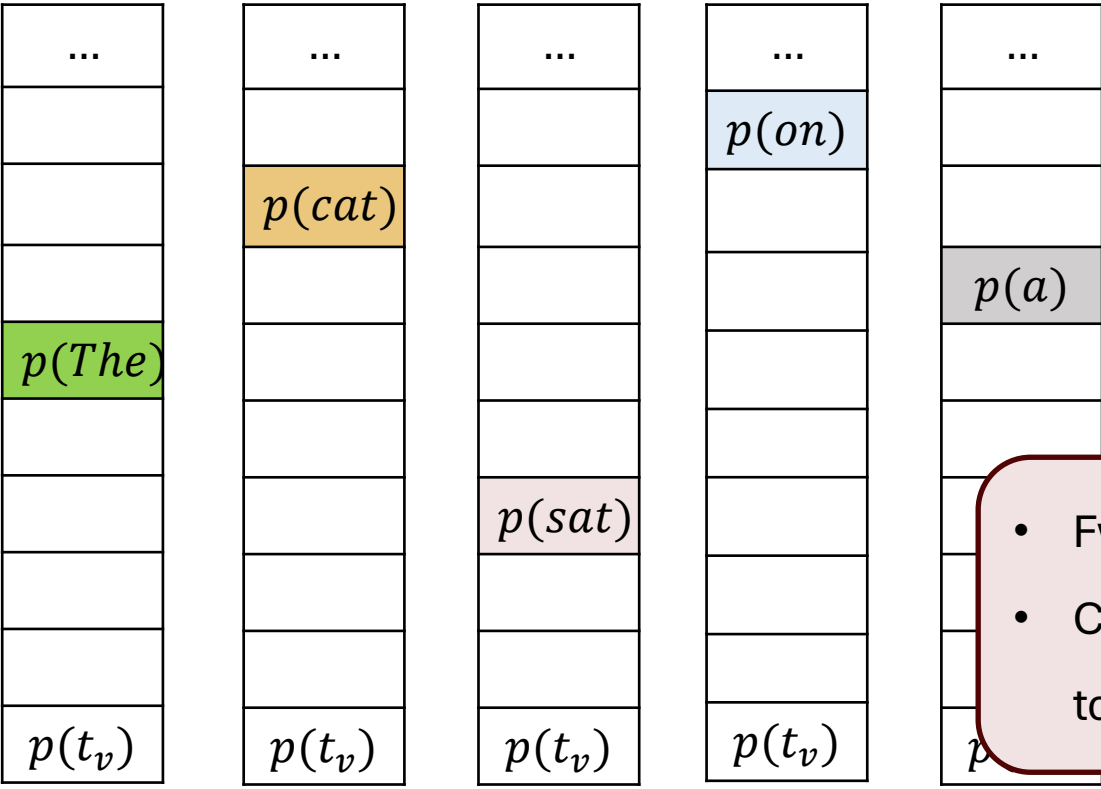
Inference through an LLM



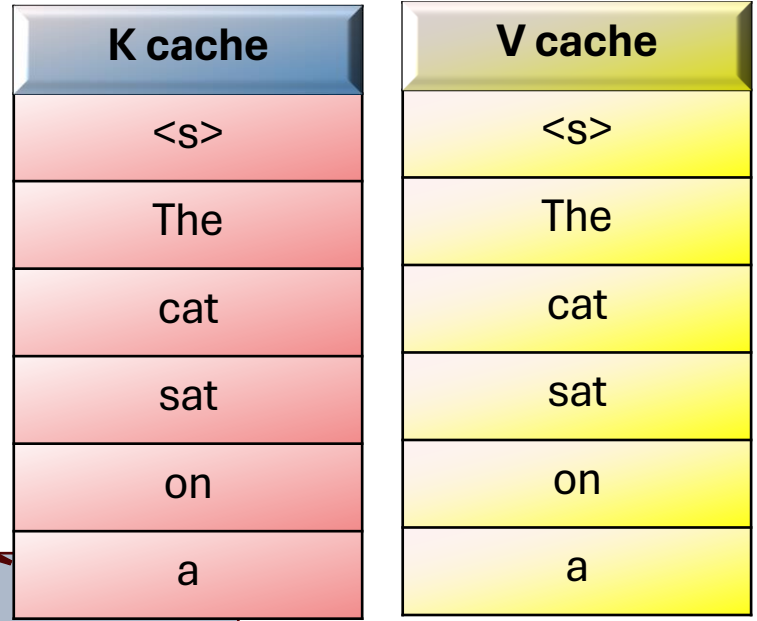
Transformer based LLM (θ)



Inference through an LLM



- Fwd. pass again (#3)
- Cache K, V emb. of token **a**

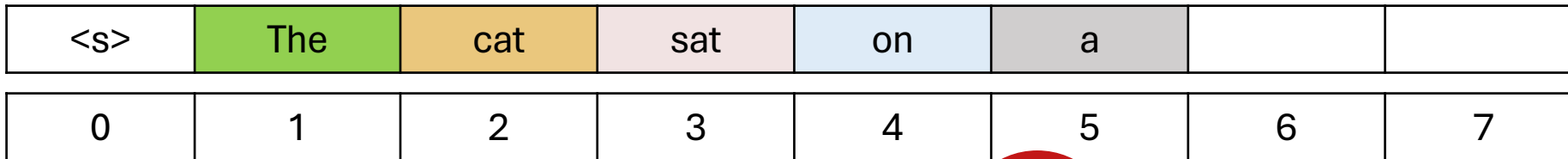
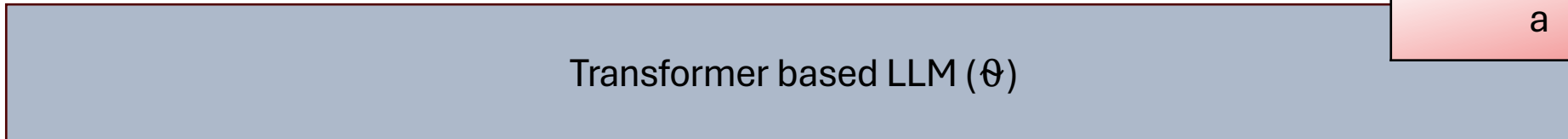
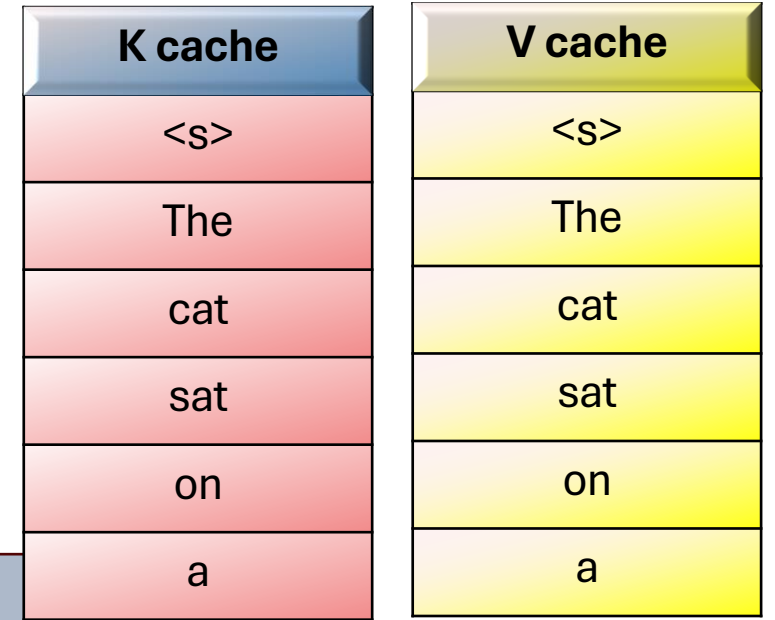
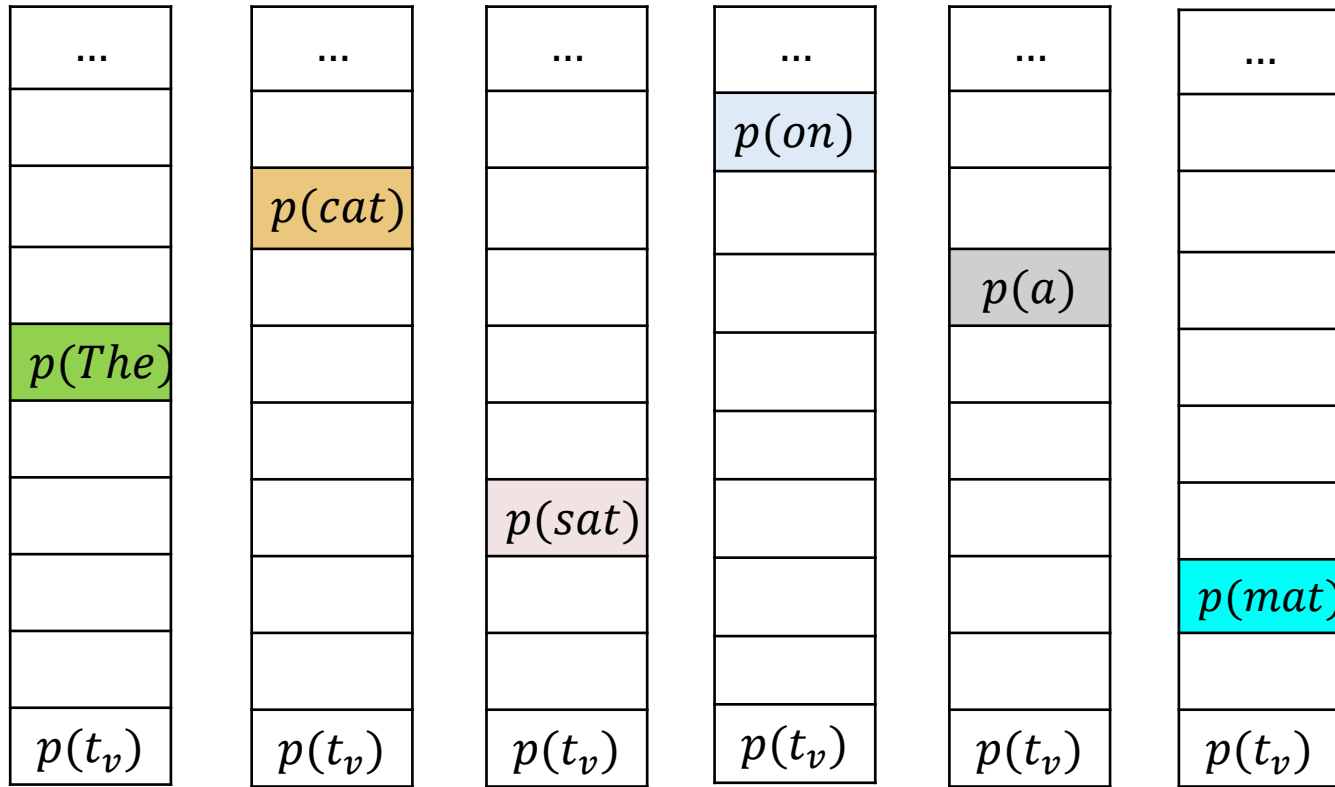


Transformer based LLM (θ)

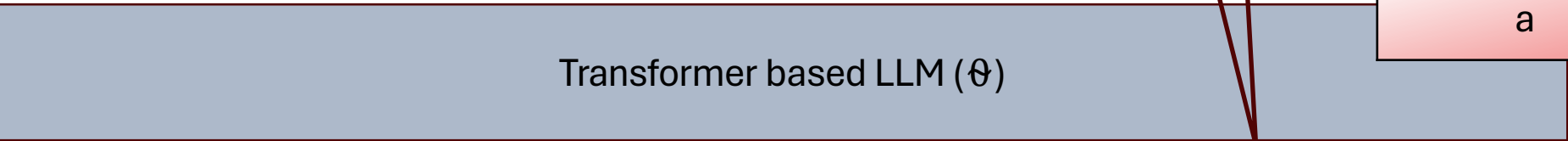
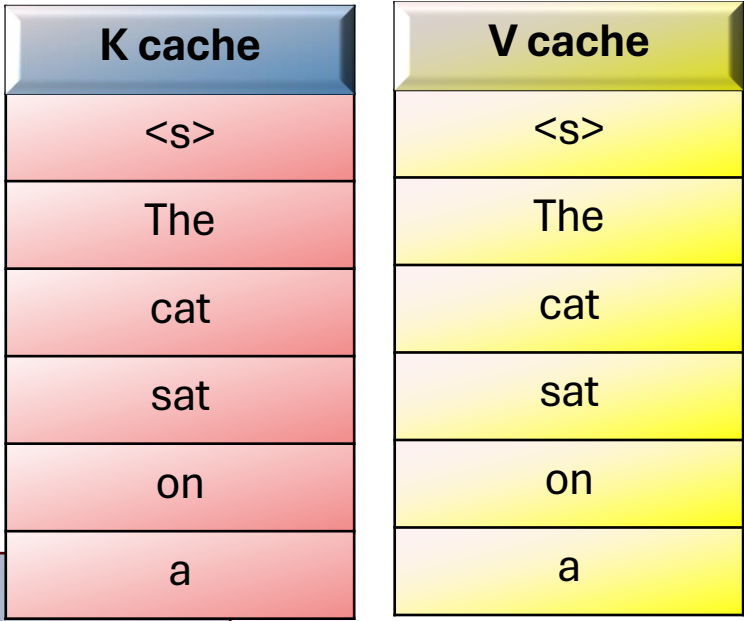
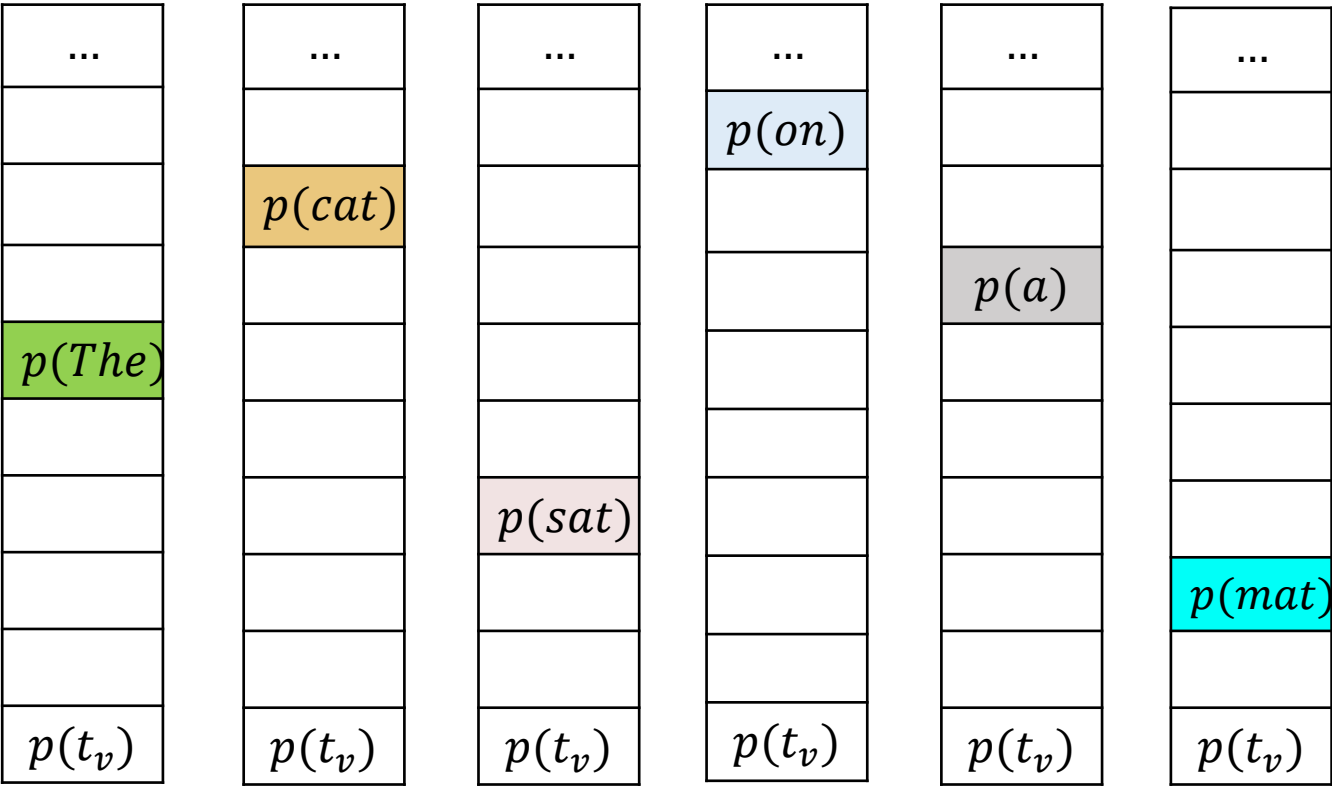
<s>	The	cat	sat	on	a		
0	1	2	3	4	5	6	7



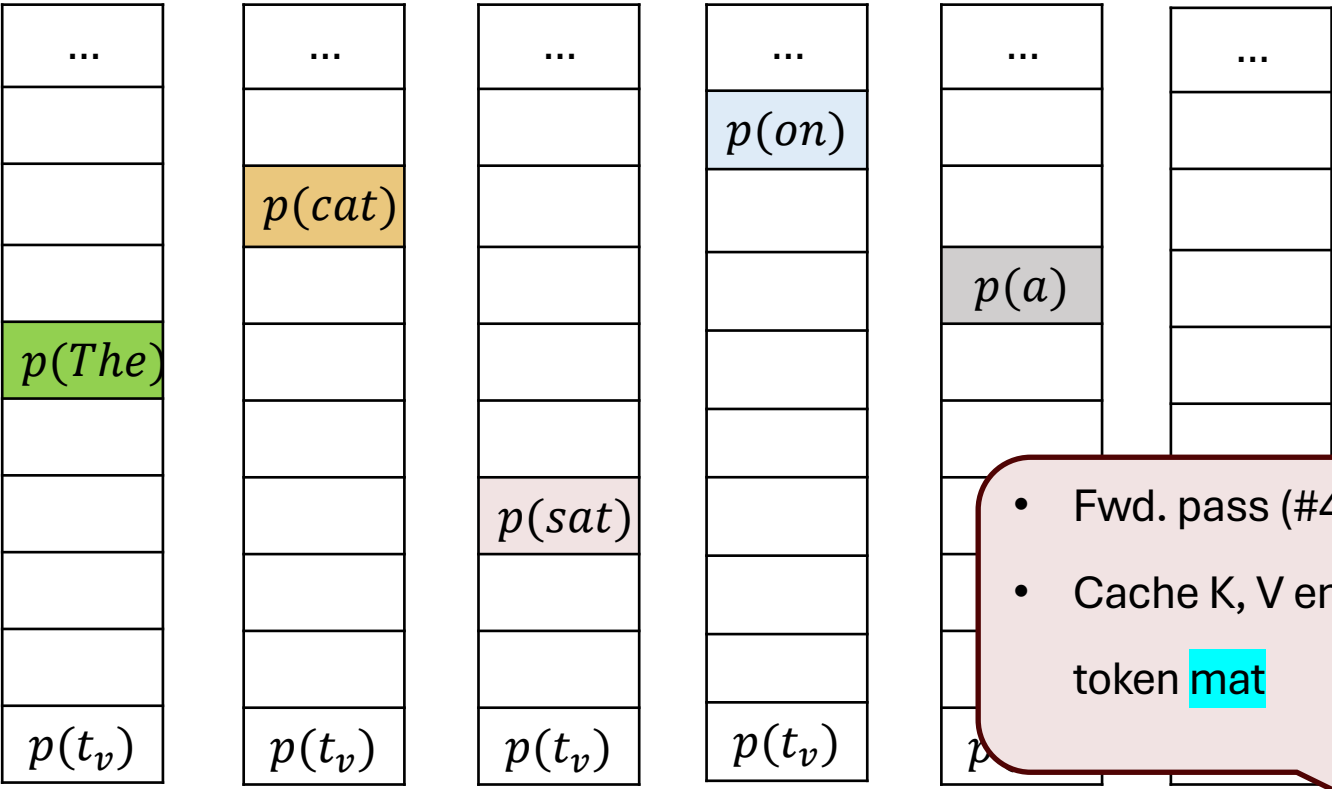
Inference through an LLM



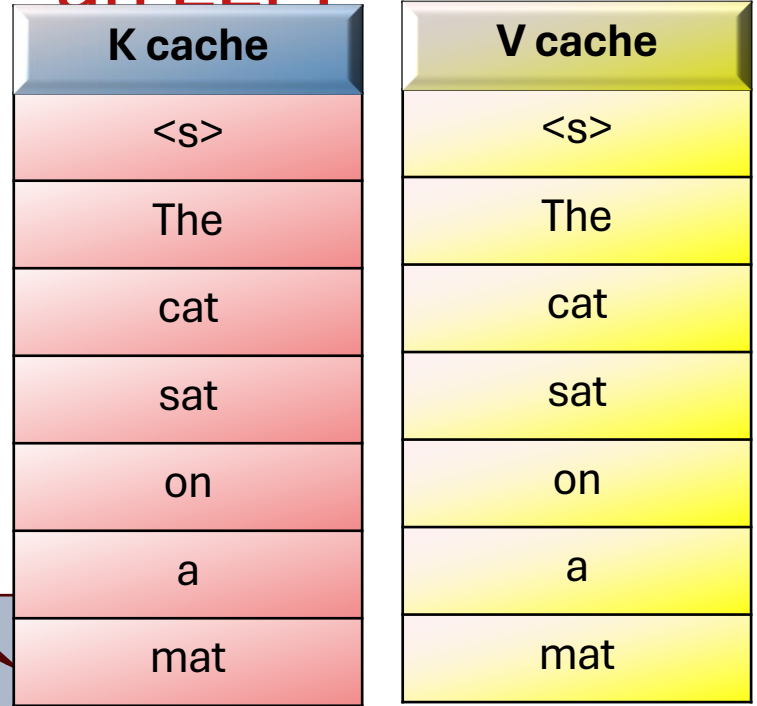
Inference through an LLM



Inference through an LLM



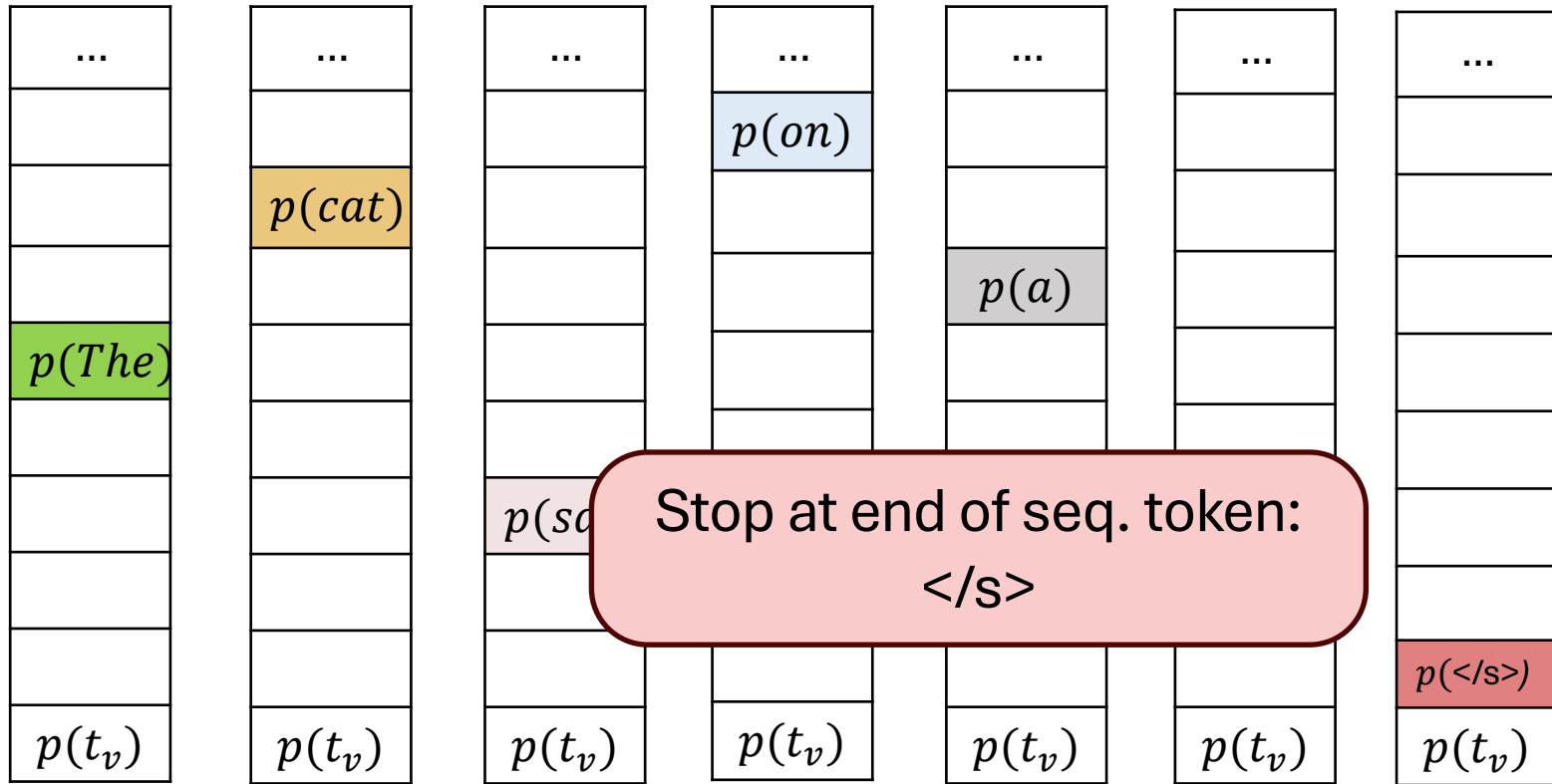
- Fwd. pass (#4)
- Cache K, V emb. of token **mat**



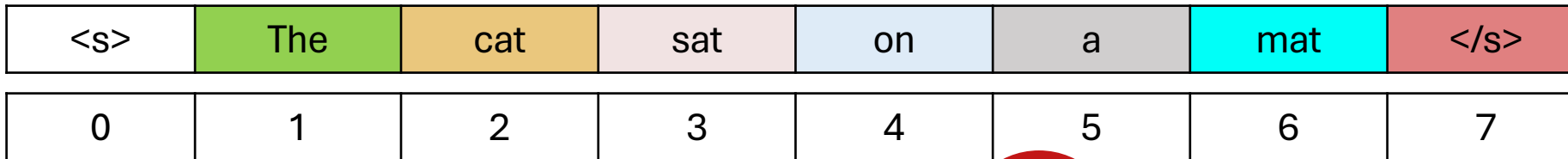
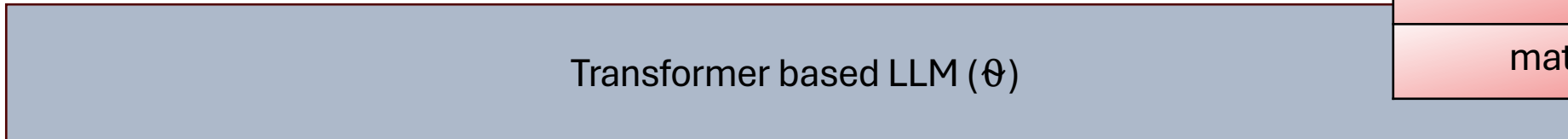
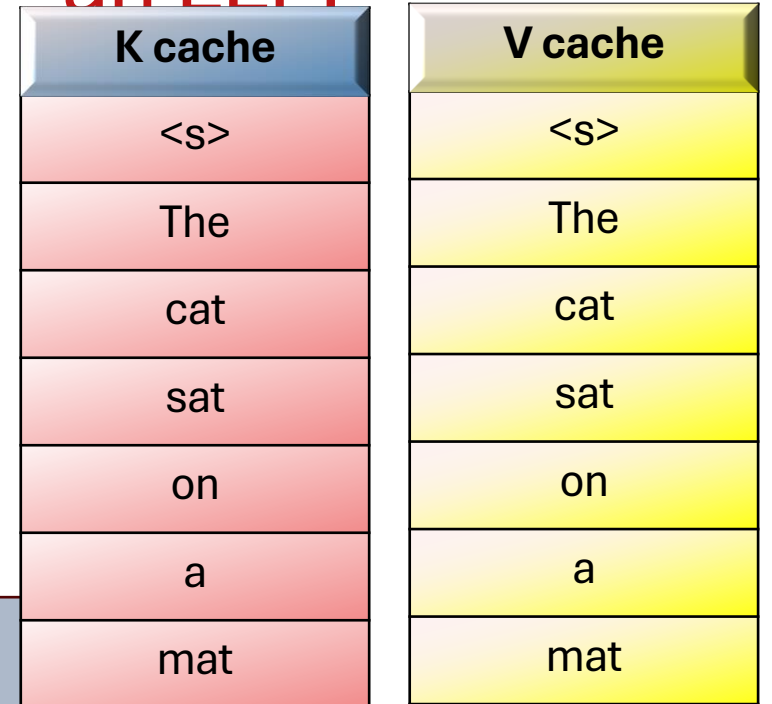
Transformer based LLM (θ)

<s>	The	cat	sat	on	a	mat	
0	1	2	3	4	5	6	7



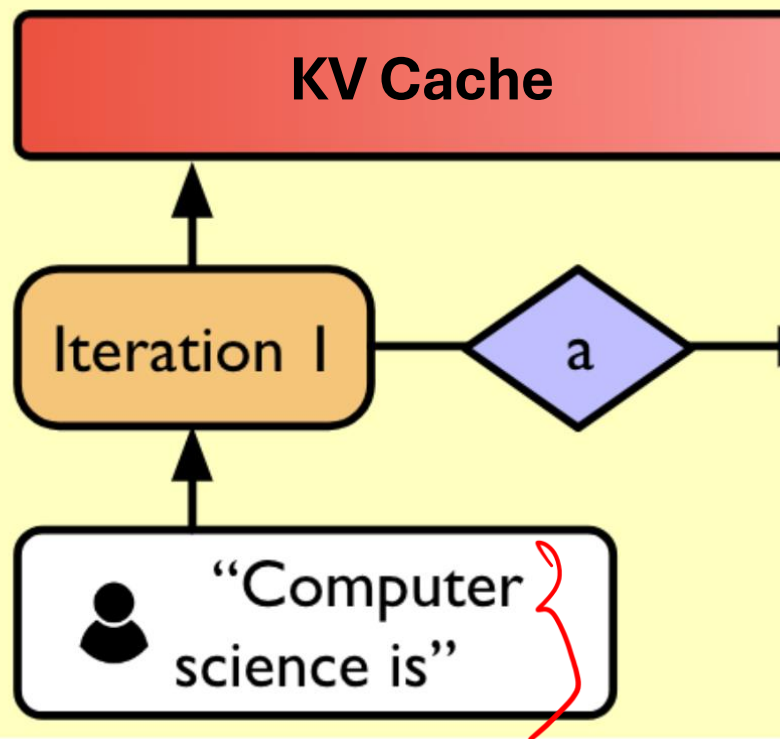


Inference through an LLM



Two stages of LLM inference

Prefill Phase



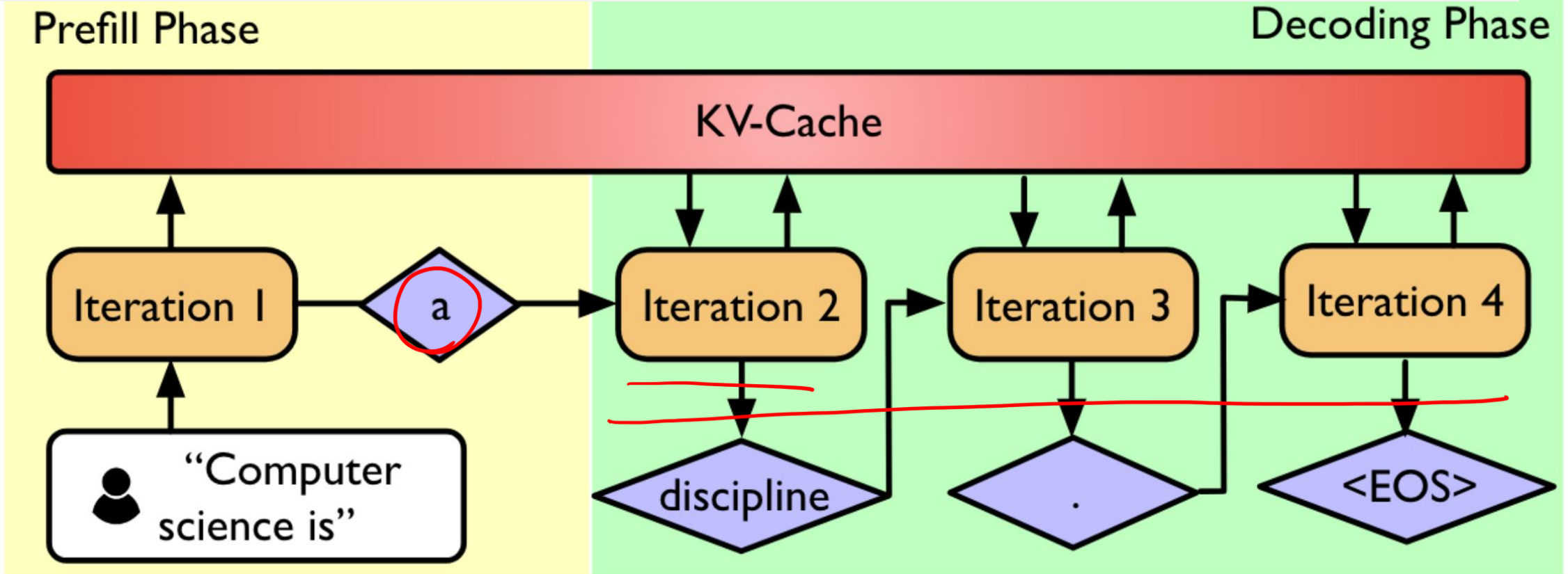
- 1st forward pass (**Pre-fill step**) **Highly parallel**
 - ❖ The entire prompt is embedded and encoded – High latency
 - ❖ Multi-head attention computes the keys and values (KV)
 - ❖ Large matrix multiplication, high usage of the hardware accelerator

Content credits: Li et al, 2024 LLM Inference Serving: Survey of Recent Advances



Remaining forward passes (Output generation): **sequential**

- The answer is generated **one token** at a time – Low latency per step
- Each generated token is **appended** to the previous input
- The process is repeated until the **stopping criteria** is met (max. length or EOS)
- Low usage of the hardware accelerator



Content credits: Li et al, 2024 LLM Inference Serving: Survey of Recent Advances and Opportunities



Inference through an LLM

- 1st forward pass (**Pre-fill step**) **Highly parallel**
 - The entire prompt is embedded and encoded – High latency
 - Multi-head attention computes the keys and values (KV)
 - Large matrix multiplication, high usage of the hardware accelerator
- Remaining forward passes (**Output generation**): **sequential**
 - The answer is generated **one token** at a time – Low latency per step
 - Each generated token is **appended** to the previous input
 - The process is repeated until the **stopping criteria** is met (max. length or EOS)
 - Low usage of the hardware accelerator

Content credits: <https://www.slideshare.net/slideshow/julien-simon-deep-dive-optimizing-llm-inference-69d3/270921961>



Memory Usage of KV cache

$$2 * \textit{precision} * N_{\textit{layers}} * d_{\textit{model}} * \textit{seqlen} * \textit{batch}$$

2 : Two matrices for K and V

precision : bytes per parameter (e.g. 4 for fp32)

$N_{\textit{layers}}$: layers in the model

$d_{\textit{model}}$: dimension of embeddings

seqlen : length of context in tokens

batch : batch size



Memory Usage of KV cache: Example OPT-13B

$$2 * precision * N_{layers} * d_{model} * seq_{len} * batch$$

- 2 : Two matrices for K and V
- $precision$: bytes per parameter (e.g. 4 for fp32)
- N_{layers} : layers in the model
- d_{model} : dimension of embeddings
- seq_{len} : length of context in tokens
- $batch$: batch size

2 (KV)
2 bytes (fp16)
40 layers
5120 dim.
2048 tokens
10



Memory Usage of KV cache: Example OPT-13B

$$2 * precision * N_{layers} * d_{model} * seq_{len} * batch$$

KV Cache: 17 GB

Model Size: $2 * 13 = 26$ GB

On a 40GB A100

- 65% (26GB) used by model parameters
- ~30% (12 GB) available for KV cache
- Expected throughput ~ 8 batch size of 2048 tokens

2 (KV)
2 bytes (fp16)
40 layers
5120 dim.
2048 tokens
10



Memory Management of KV Cache

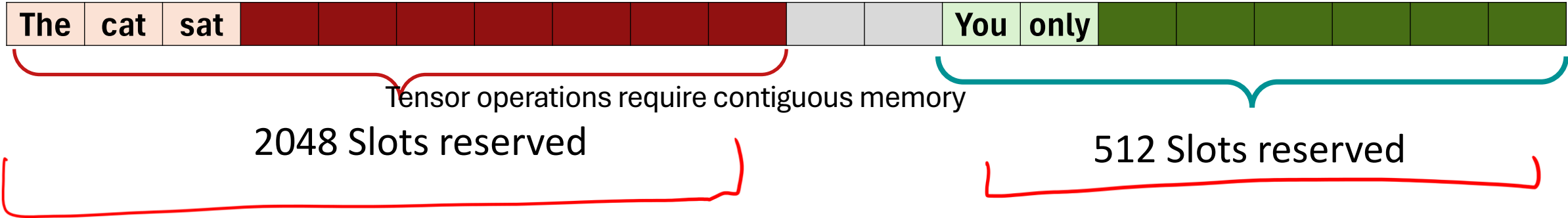
Prompt A: *“The cat sat”*
Max Tokens: 2048

Prompt B: *“You only”*
Max Tokens: 512



Memory Management of KV Cache

Tensor operations require contiguous memory

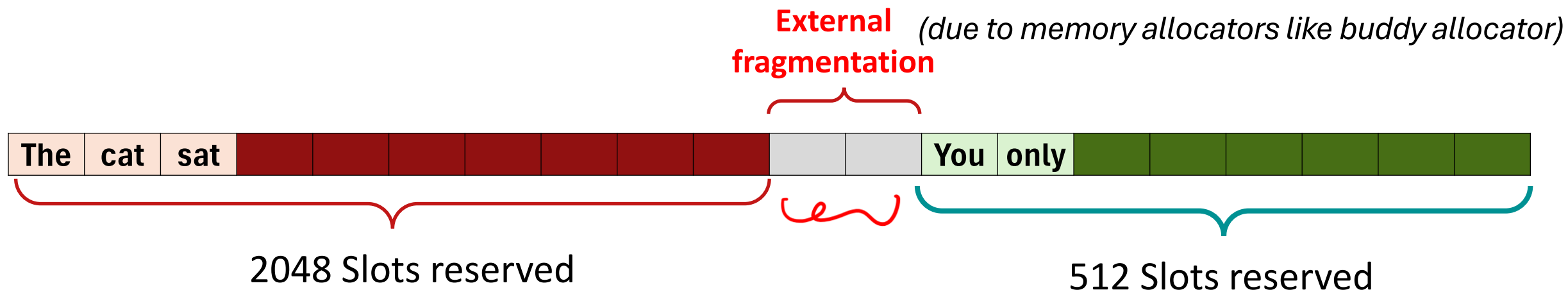


Prompt A: ***"The cat sat"***
Max Tokens: ***2048***

Prompt B: ***"You only"***
Max Tokens: ***512***



Memory Management of KV Cache



Prompt A: ***"The cat sat"***

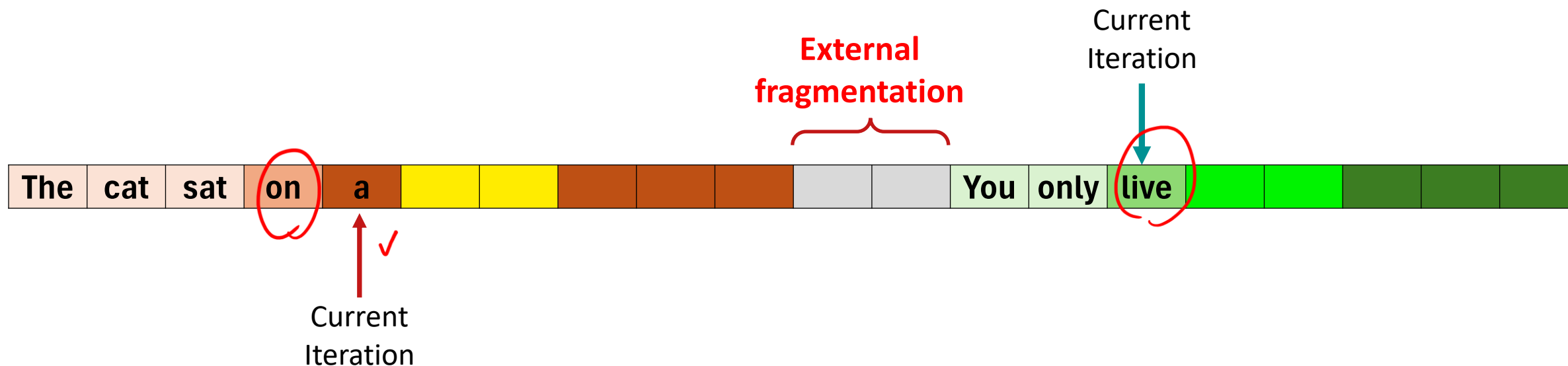
Max Tokens: ***2048***

Prompt B: ***"You only"***

Max Tokens: ***512***



Memory Management of KV Cache

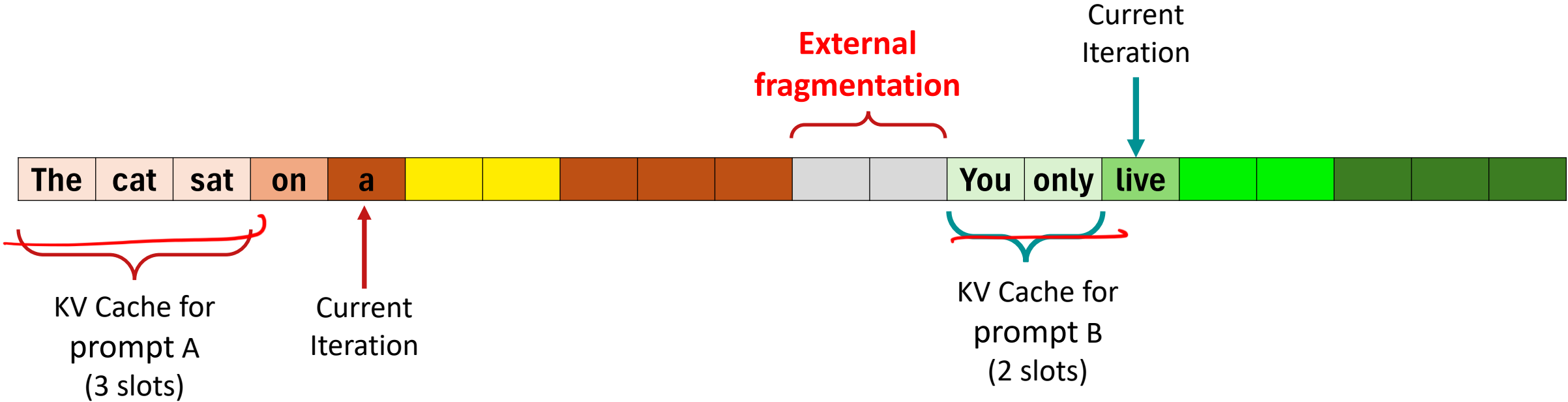


Prompt A: *"The cat sat"*
Max Tokens: 2048

Prompt B: *"You only"*
Max Tokens: 512



Memory Management of KV Cache

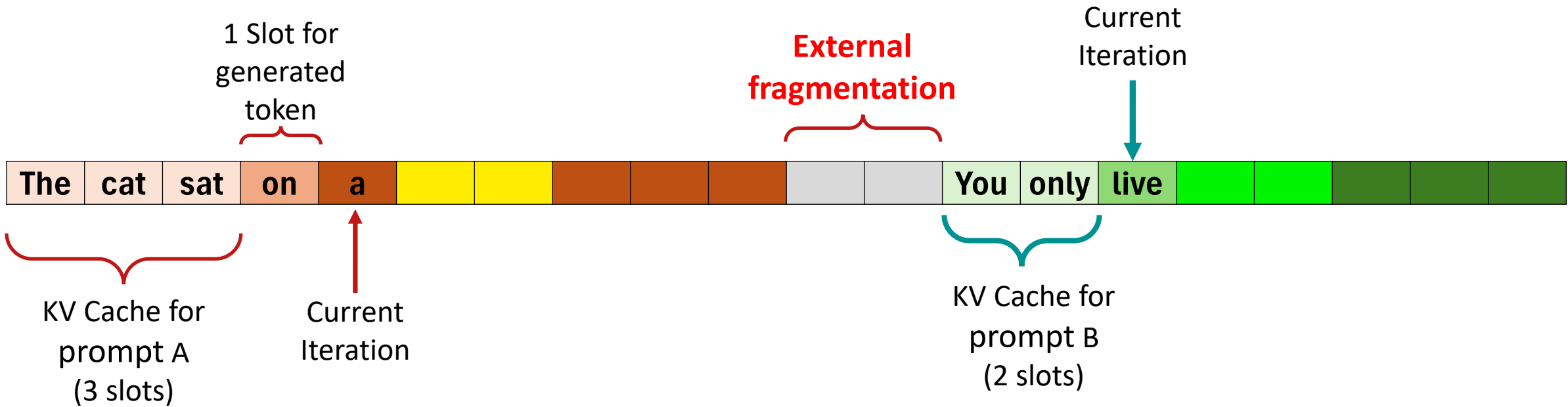


Prompt A: *"The cat sat"*
Max Tokens: 2048

Prompt B: *"You only"*
Max Tokens: 512



Memory Management of KV Cache

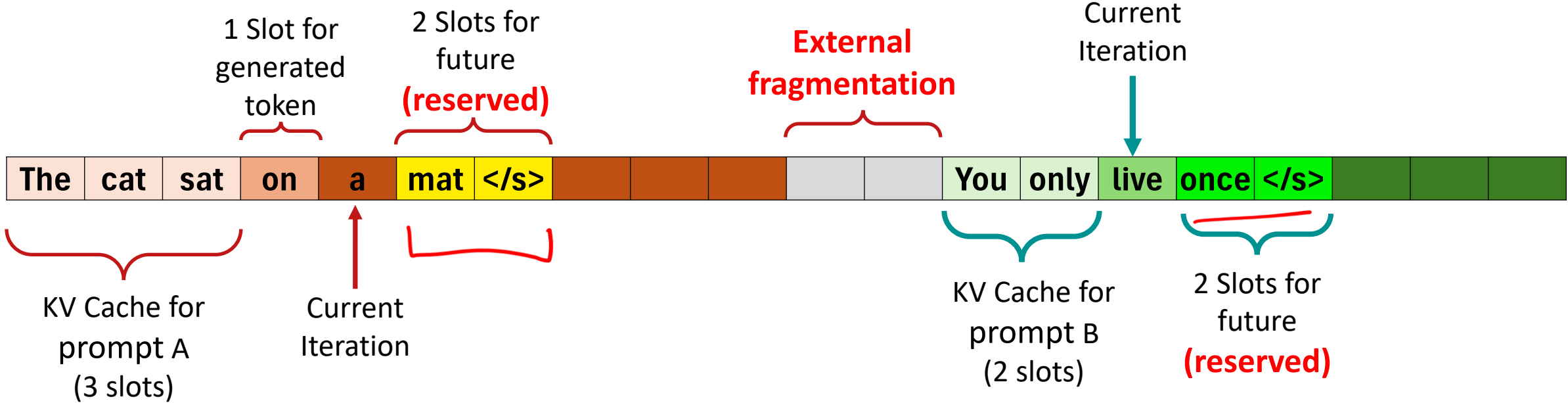


Prompt A: *"The cat sat"*
Max Tokens: 2048

Prompt B: *"You only"*
Max Tokens: 512



Memory Management of KV Cache

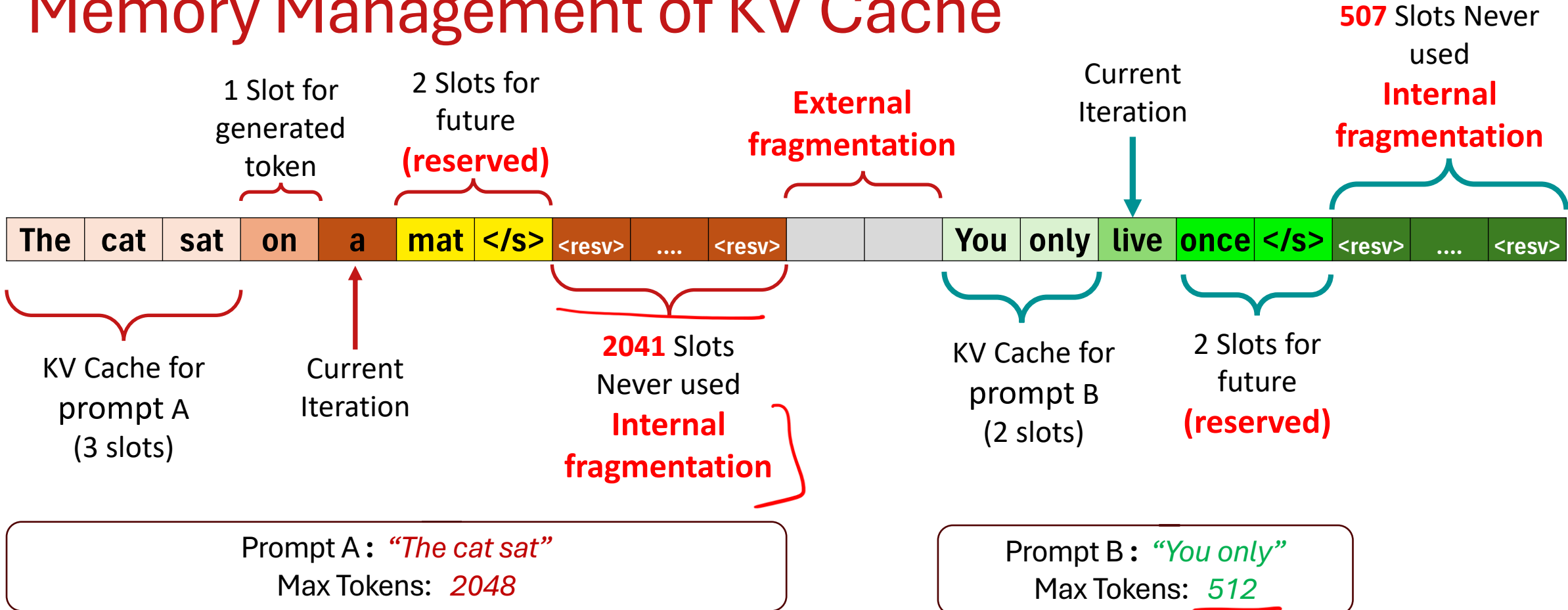


Prompt A: *"The cat sat"*
 Max Tokens: 2048

Prompt B: *"You only live once"*
 Max Tokens: 512



Memory Management of KV Cache



Memory Management of KV Cache

Chunk Pre-allocation scheme

- KV cache stored in contiguous memory
- Chunks of memory allocated statically, based on max. tokens.
- Actual input or eventual output length ignored while allocating memory



Memory Management of KV Cache

Chunk Pre-allocation scheme

- KV cache stored in contiguous memory
- Chunks of memory allocated statically, based on max. tokens.
- Actual input or eventual output length ignored while allocating memory

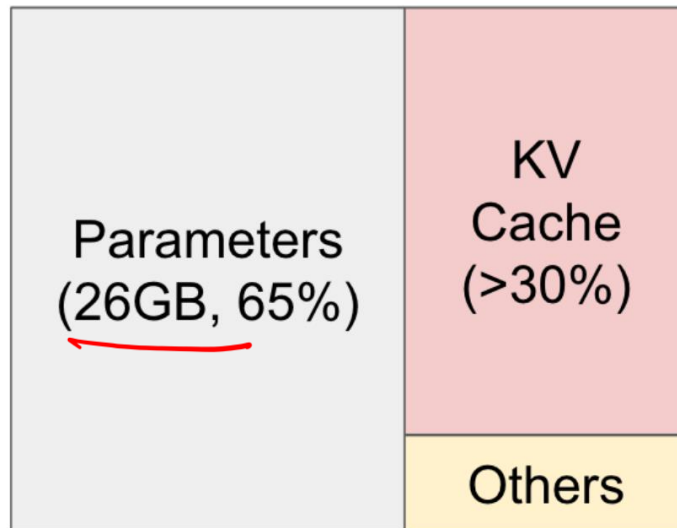
Results in 3 types of memory wastes –

- **Reserved slots** for future tokens
- **Internal fragmentation** due to over-provisioning for maximum sequence lengths
- **External fragmentation** from the memory allocator.



Memory Layout for 13B-OPT model on A100 (40GB)

20.4-38.2% utilized

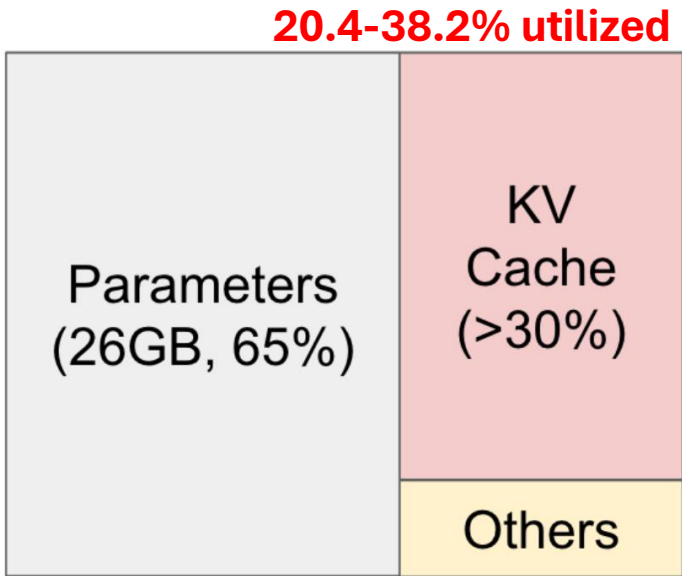


NVIDIA A100 40GB

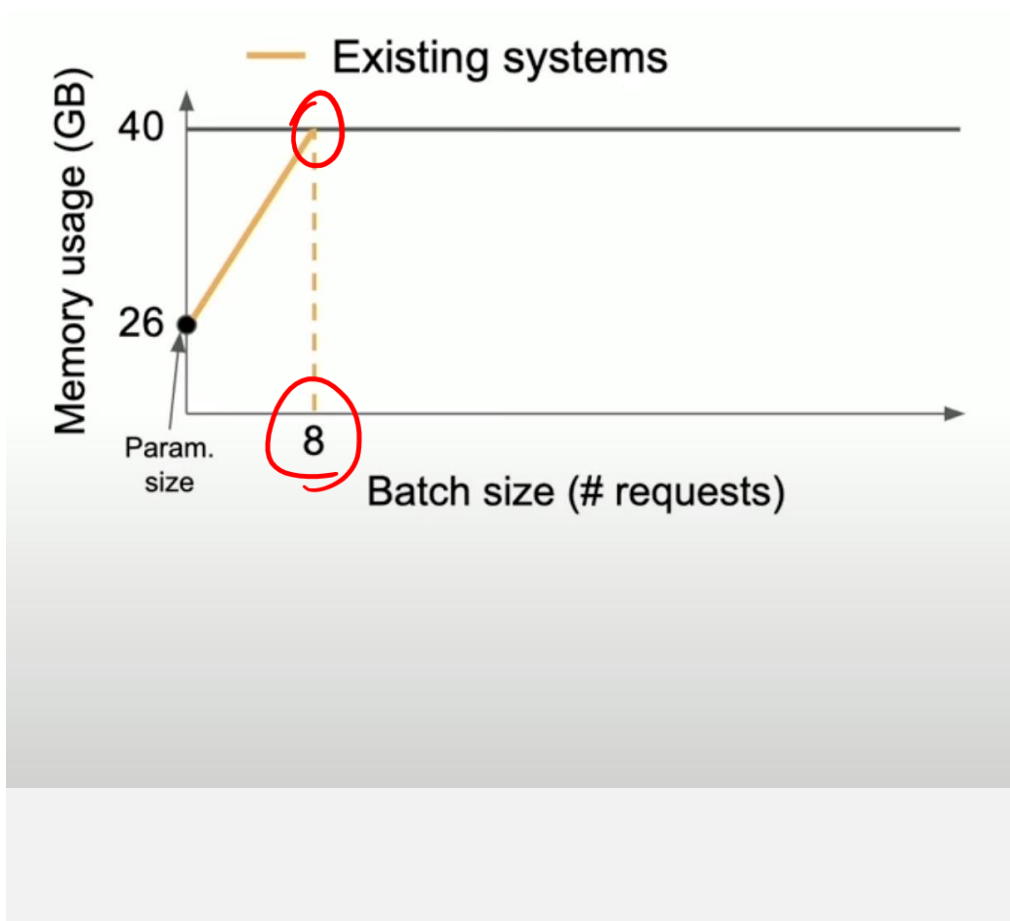
Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale



Memory Layout for 13B-OPT model on A100 (40GB)



NVIDIA A100 40GB



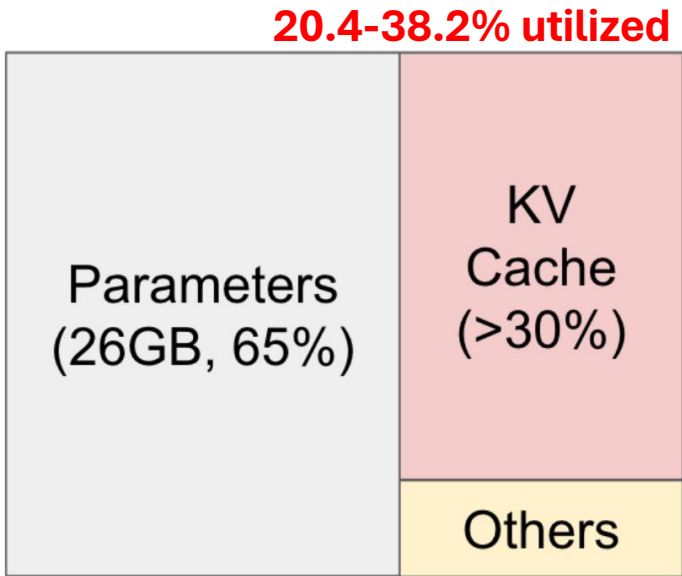
Existing systems

- max batch size - 8

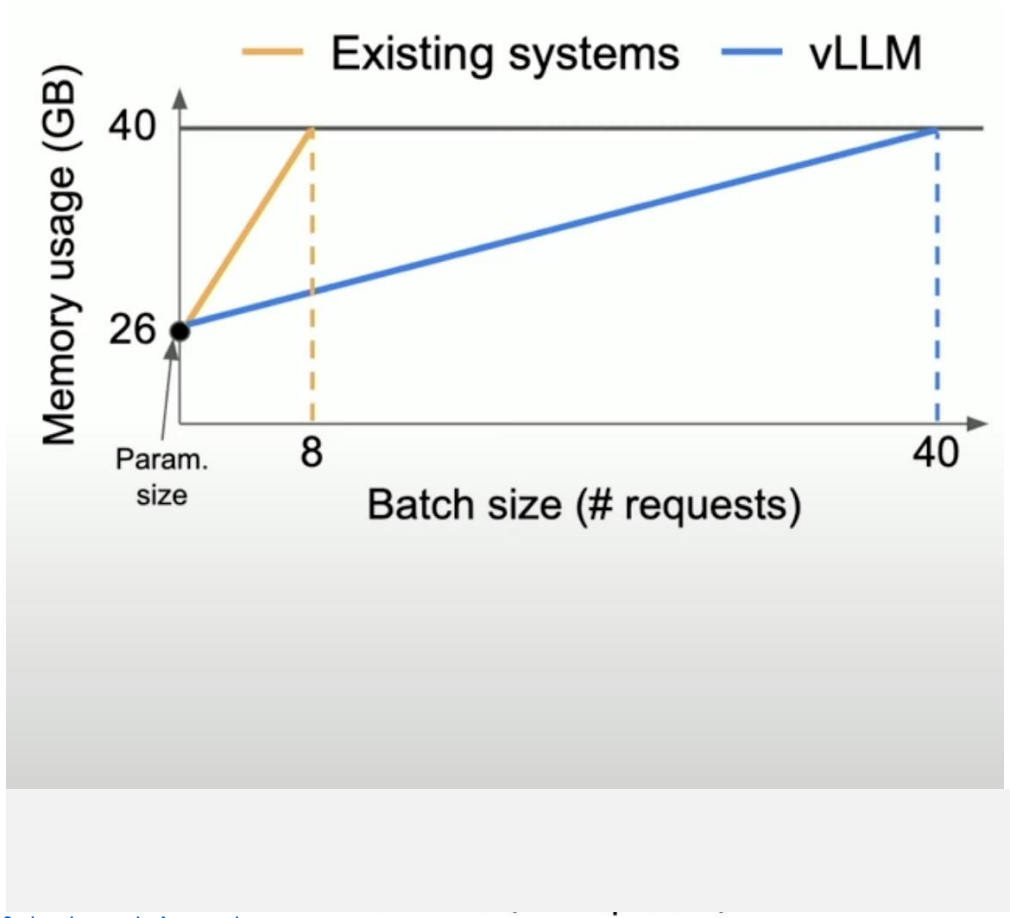
Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale



Memory Layout for 13B-OPT model on A100 (40GB)



NVIDIA A100 40GB



Existing systems

- max batch size - 8

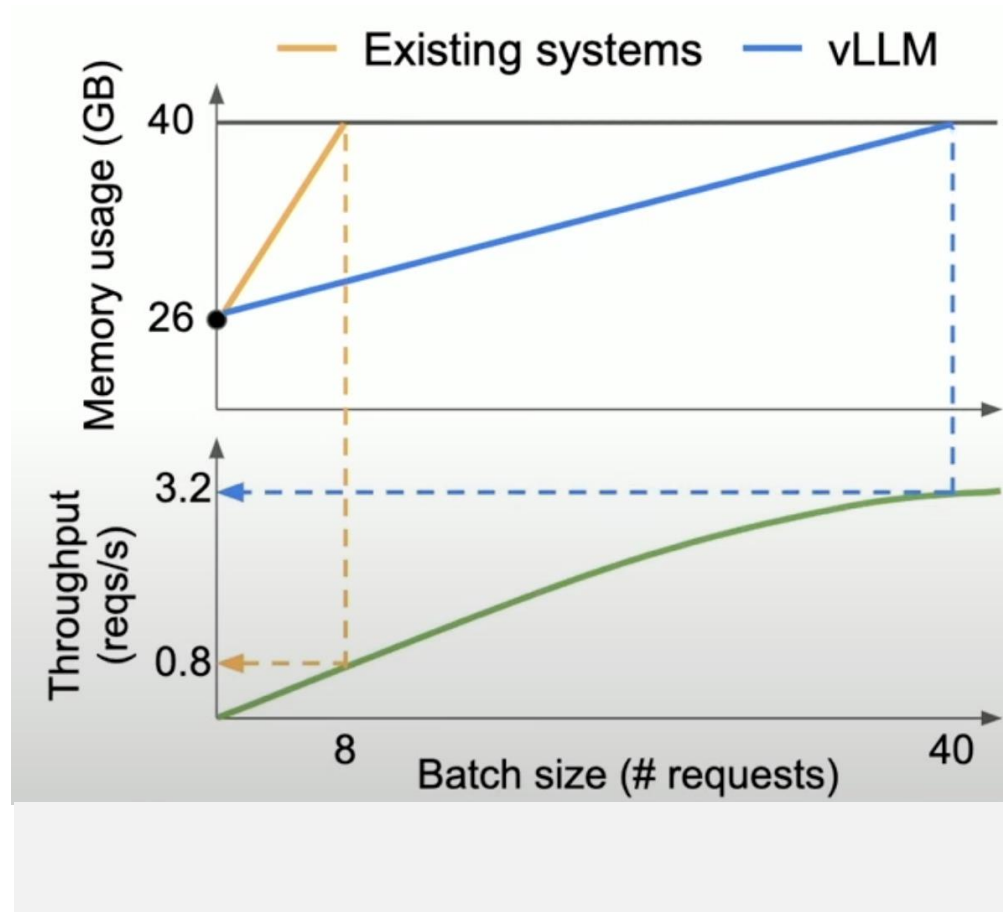
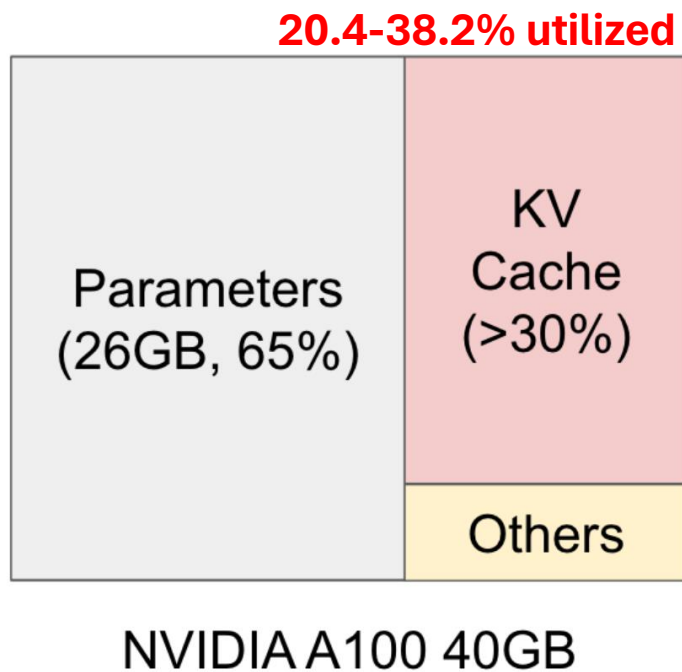
vLLM (paged attention)

- Max batch size ~ 40

Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale



Memory Layout for 13B-OPT model on A100 (40GB)



Existing systems

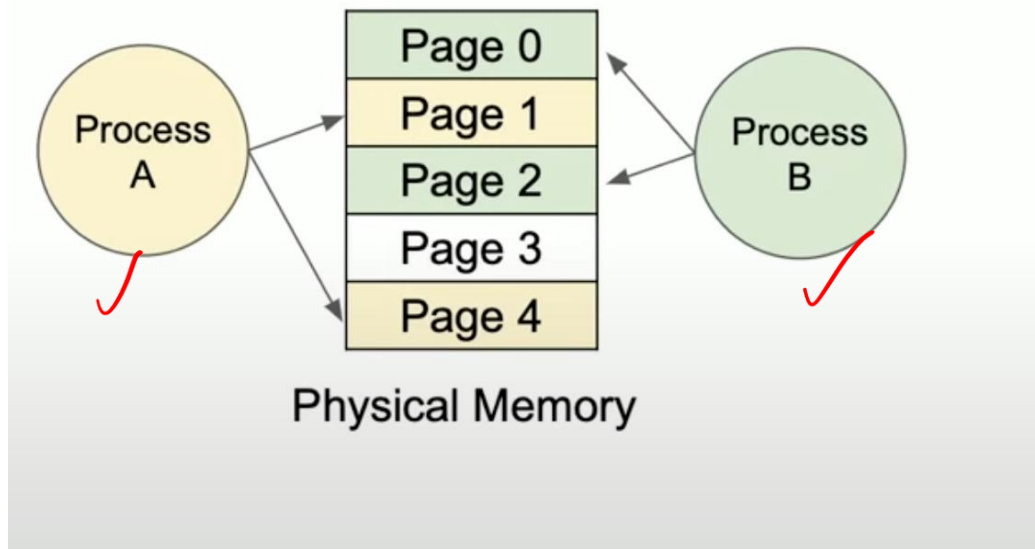
- max batch size - 8
- ~ 0.8 requests / sec

vLLM (paged attention)

- Max batch size ~ 38
- ~ 3.2 requests per sec

vLLM: Efficient KV cache management

Inspired by **Virtual memory** and paging



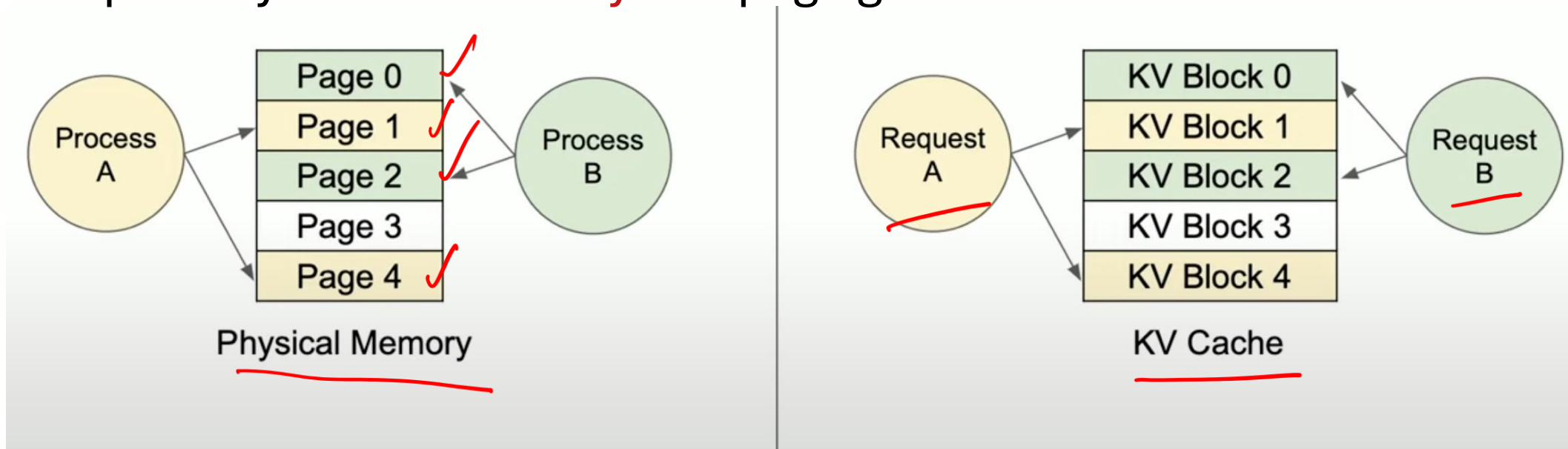
Memory management in OS

Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale



vLLM: Efficient KV cache management

Inspired by **Virtual memory** and paging

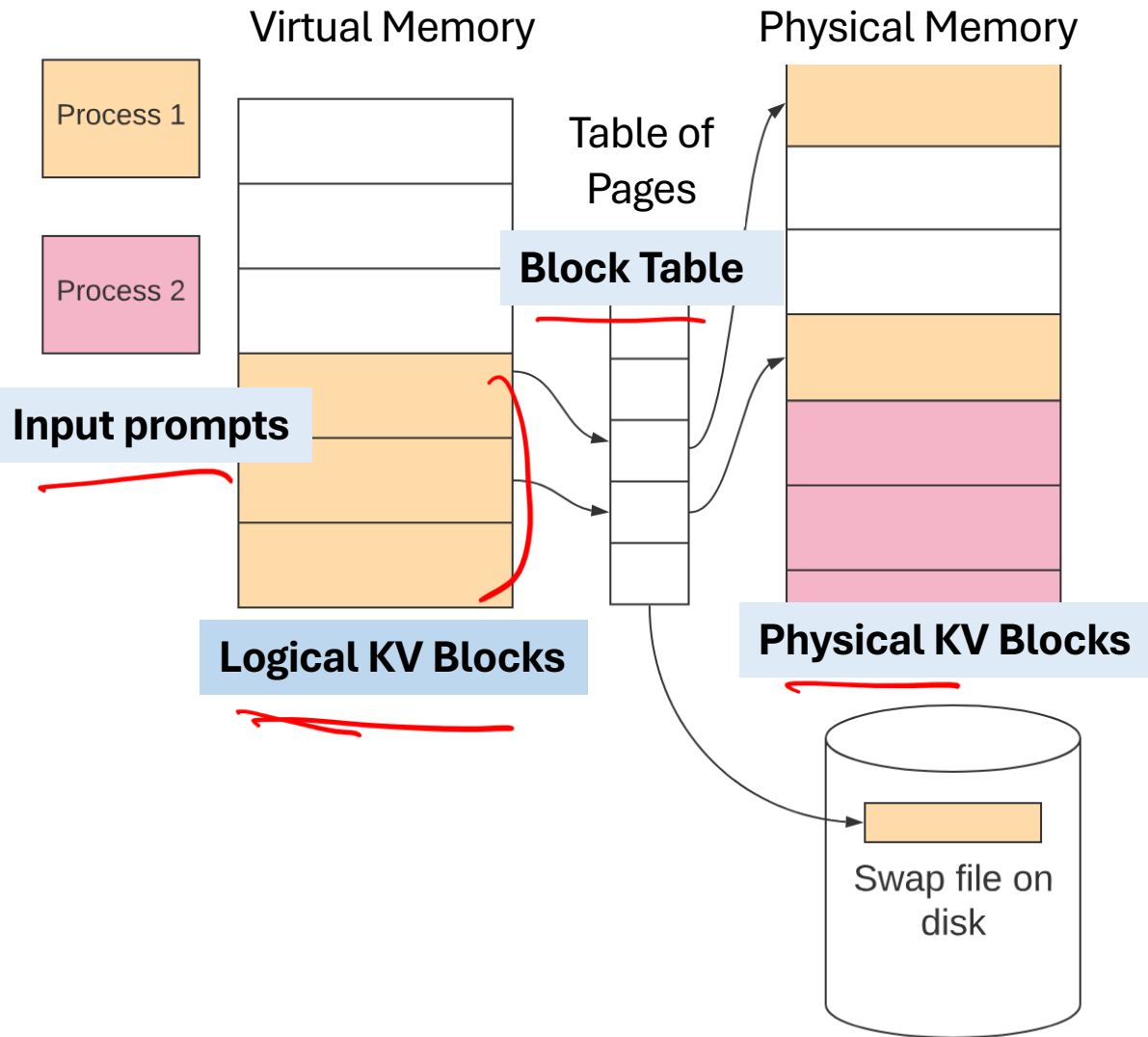


Memory management in OS

Memory management in vLLM

Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale





Efficient KV cache management

Inspired by **Virtual memory** and paging

- ❑ Processes as **incoming requests** (input to the model)
- ❑ Virtual Memory to **Logical KV Blocks**
- ❑ Physical Memory to **Physical KV Blocks**
- ❑ Page table to **Block Table**

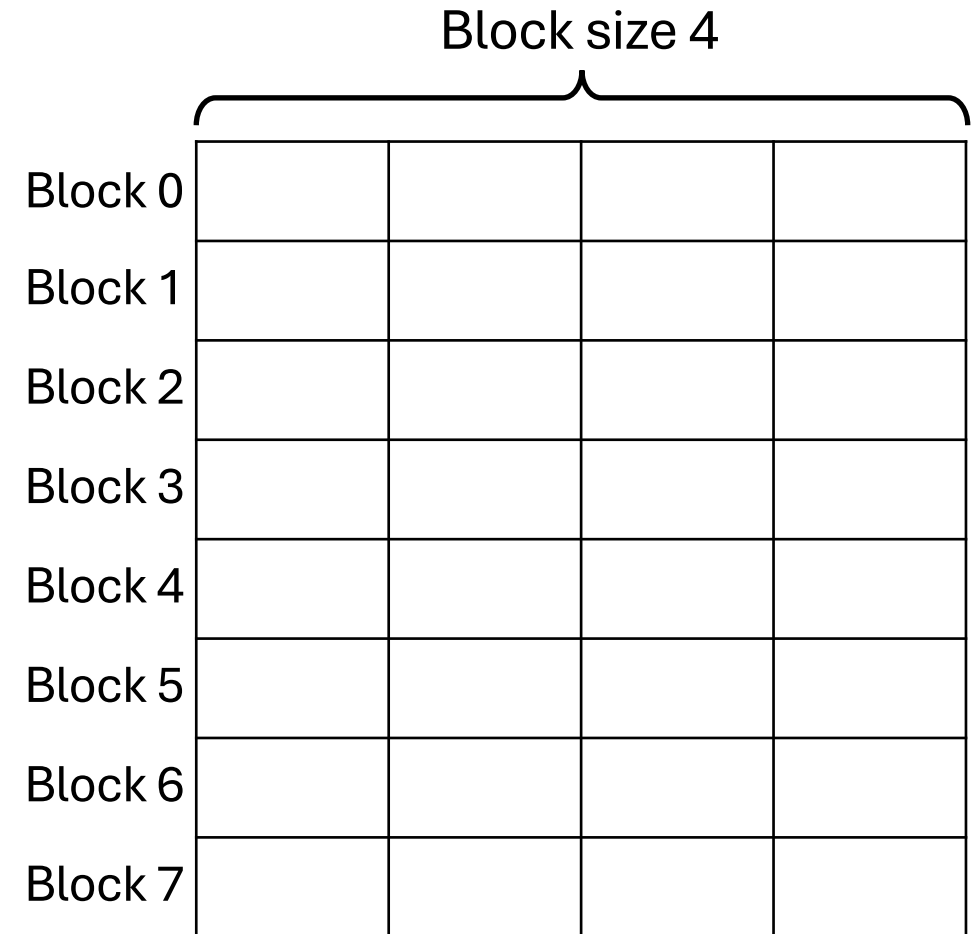
KV Blocks

KV Cache

Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale



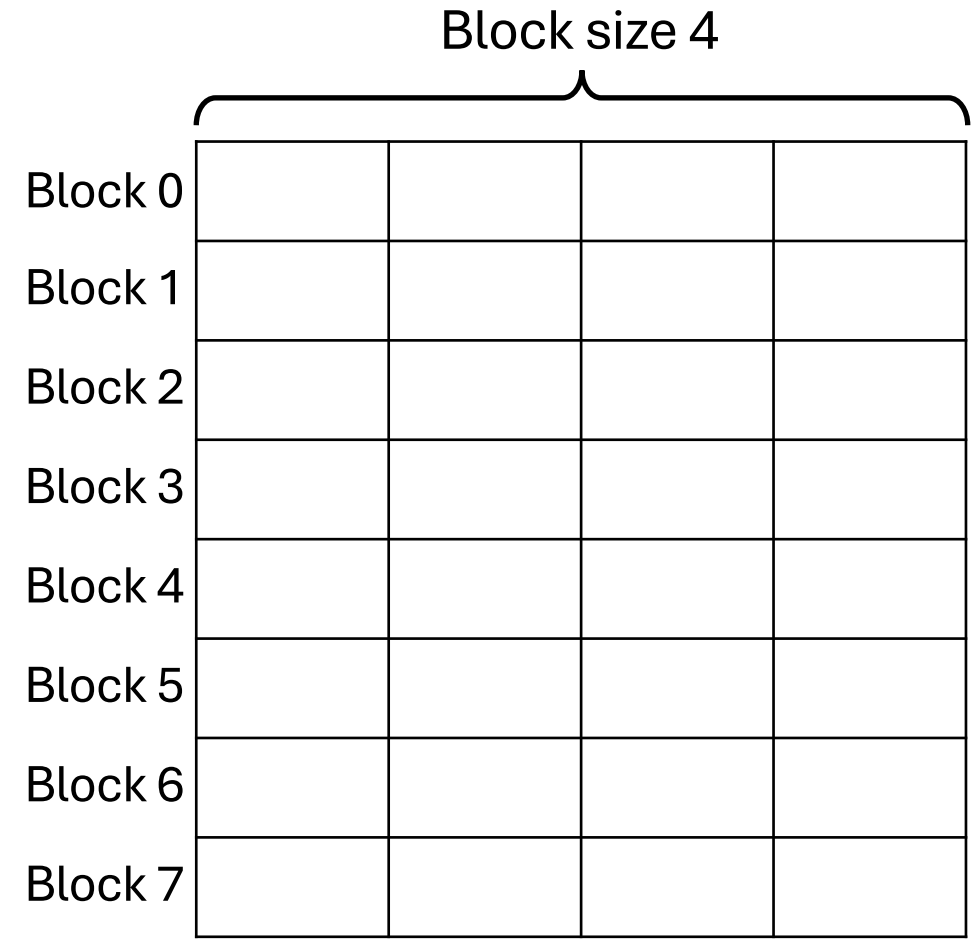
KV Blocks



Content credits: https://www.youtube.com/watch?v=5ZlavKF_98U&t=1646s&ab_channel=Anyscale



KV Blocks



Physical KV Blocks



Physical vs Logical KV Blocks

Block 0				
1				
2				
3				

Logical KV Blocks

Block size 4

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

Physical KV Blocks



Physical vs Logical KV Blocks

Block 0				
1				
2				
3				

Logical KV Blocks

Phys. Block	# Filled

Block Table

Block size 4

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

Physical KV Blocks



Physical vs Logical KV Blocks

Prompt: “Today we are learning about LLMs and”

Block 0	Today	we	are	learning
1	about	LLMs	and	
2				
3				

Logical KV Blocks

Phys. Block	# Filled

Block Table

Block size 4

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

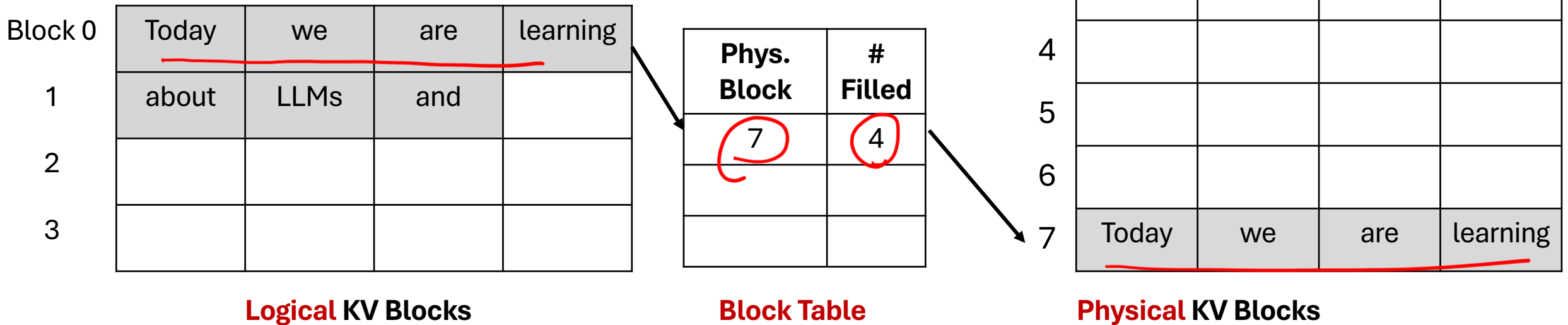
Physical KV Blocks



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"

Block size 4



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"

Block size 4

Block 0

0	Today	we	are	learning
1	about	LLMs	and	
2				
3				

Logical KV Blocks

Phys. Block	# Filled
7	4
1	3

Block Table

Block 0

0				
1	about	LLMs	and	
2				
3				
4				
5				
6				
7	Today	we	are	learning

Physical KV Blocks



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"
 Completion: "*memory*"

Block size 4

Block 0

0	Today	we	are	learning
1	about	LLMs	and	memory
2				
3				

Logical KV Blocks

Phys. Block	# Filled
7	4
1	4

Block Table

Block 0

0				
1	about	LLMs	and	
2				
3				
4				
5				
6				
7	Today	we	are	learning

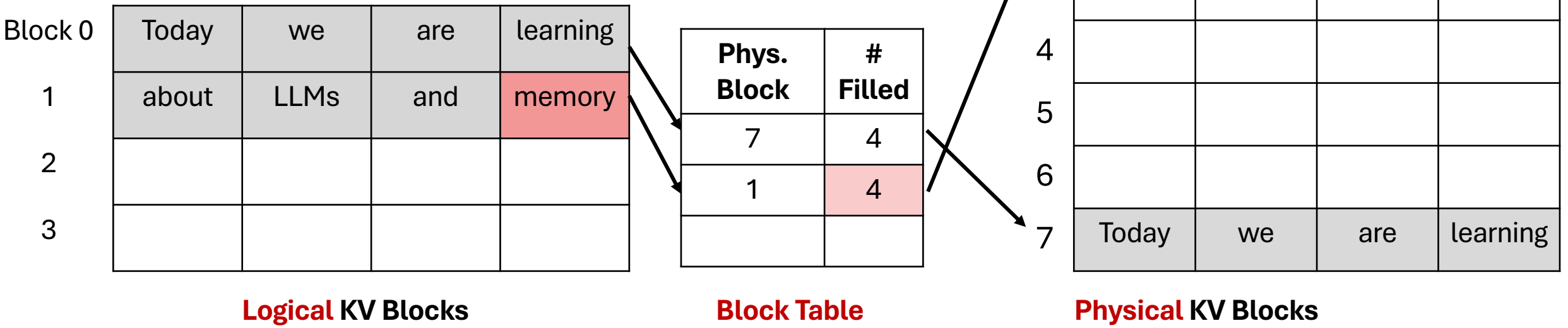
Physical KV Blocks



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"
 Completion: "*memory*"

Block size 4



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"
 Completion: "*memory on*"

Block size 4

Block 0

0	Today	we	are	learning
1	about	LLMs	and	memory
2	on			
3				

Logical KV Blocks

Phys. Block	# Filled
7	4
1	4

Block Table

Block 0

0				
1	about	LLMs	and	memory
2				
3				
4				
5				
6				
7	Today	we	are	learning

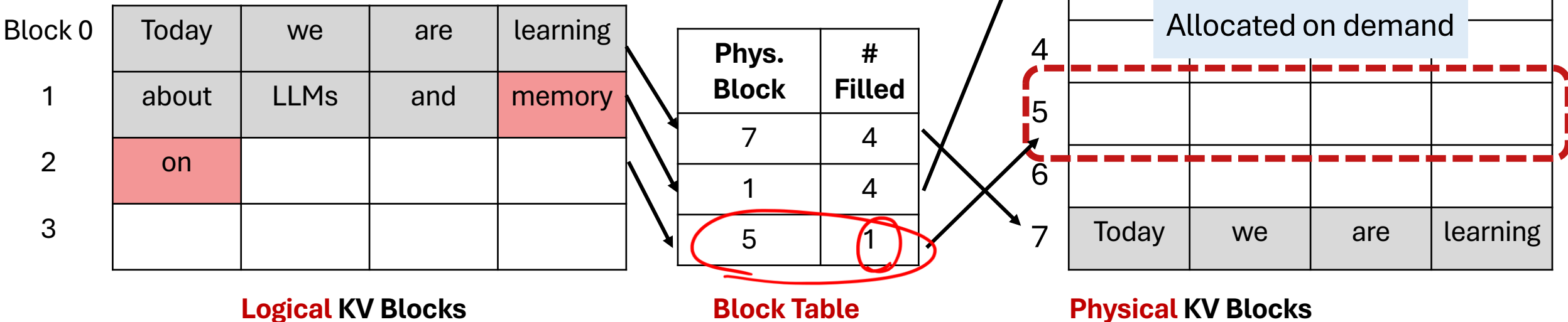
Physical KV Blocks



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"
Completion: "*memory on*"

Block size 4



Content credits: <https://youtu.be/yVXtLTcdO1Q?si=XO2Dk-VYOShUMH1u>



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"
 Completion: "*memory on*"

Block size 4

Block 0
1
2
3

Block 0	Today	we	are	learning
1	about	LLMs	and	memory
2	on			
3				

Logical KV Blocks

Phys. Block	# Filled
7	4
1	4
5	1

Block Table

Block 0
1
2
3
4
5
6
7

Block size 4				
1	about	LLMs	and	memory
2				
3				
4				
5	on			
6				
7	Today	we	are	learning

Allocated on demand

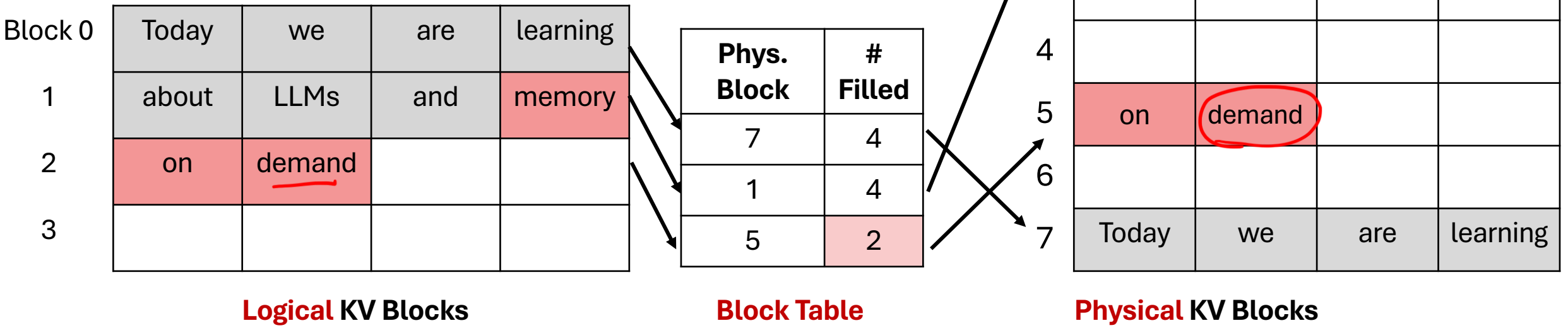
Physical KV Blocks



Physical vs Logical KV Blocks

Prompt: "Today we are learning about LLMs and"
 Completion: "*memory on demand*"

Block size 4

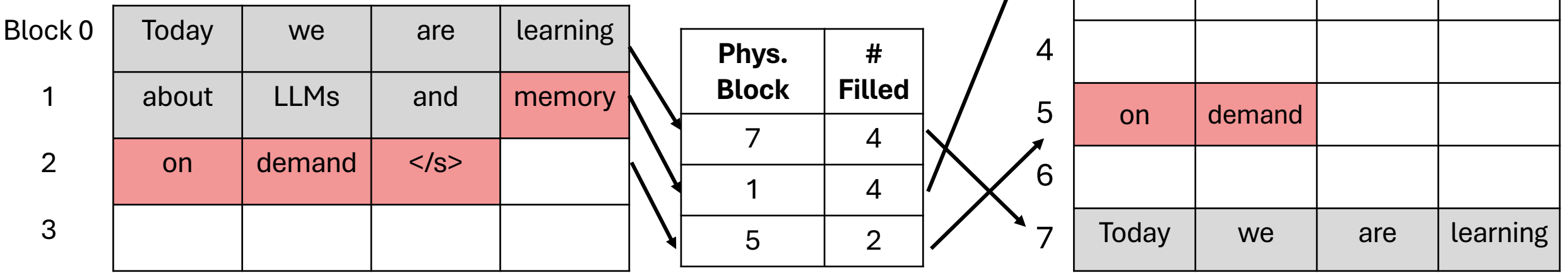


Content credits: <https://youtu.be/yVXtLTcdO1Q?si=XO2Dk-VYOShUMH1u>



Physical vs Logical KV Blocks

Block size 4



Logical KV Blocks

Block Table

Physical KV Blocks

Prompt A: “Today we are learning about LLMs and”

Completion: “*memory on demand </s>*”

Content credits: <https://youtu.be/yVXtLTcdO1Q?si=XO2Dk-VYOShUMH1u>



Physical vs Logical KV Blocks

Block size 4

Block 0
1
2
3

Block 0	Today	we	are	learning
1	about	LLMs	and	memory
2	on	demand	</s>	
3				

Logical KV Blocks

Phys. Block	# Filled
7	4
1	4
5	2

Block Table

Block 0
1
2
3
4
5
6
7

about	LLMs	and	memory
on	demand		
Today	we	are	learning

Physical KV Blocks

Internal fragmentation

Prompt A: "Today we are learning about LLMs and"
Completion: "memory on demand </s>"

Content credits: <https://youtu.be/yVXtLTcdO1Q?si=XO2Dk-VYOShUMH1u>



0	Today	we	are	learning
1	about	LLMs	and	
2				
3				

Logical KV Blocks - B

0	Today	we	are	learning
1	about	LLMs	and	memory
2	on	demand	</s>	
3				

Logical KV Blocks - A

Phys. Block	# Filled
7	4
1	4
5	2

Block Table -A

Block size 4

0				
1	about	LLMs	and	memory
2				
3				
4				
5	on	demand		
6				
7	Today	we	are	learning

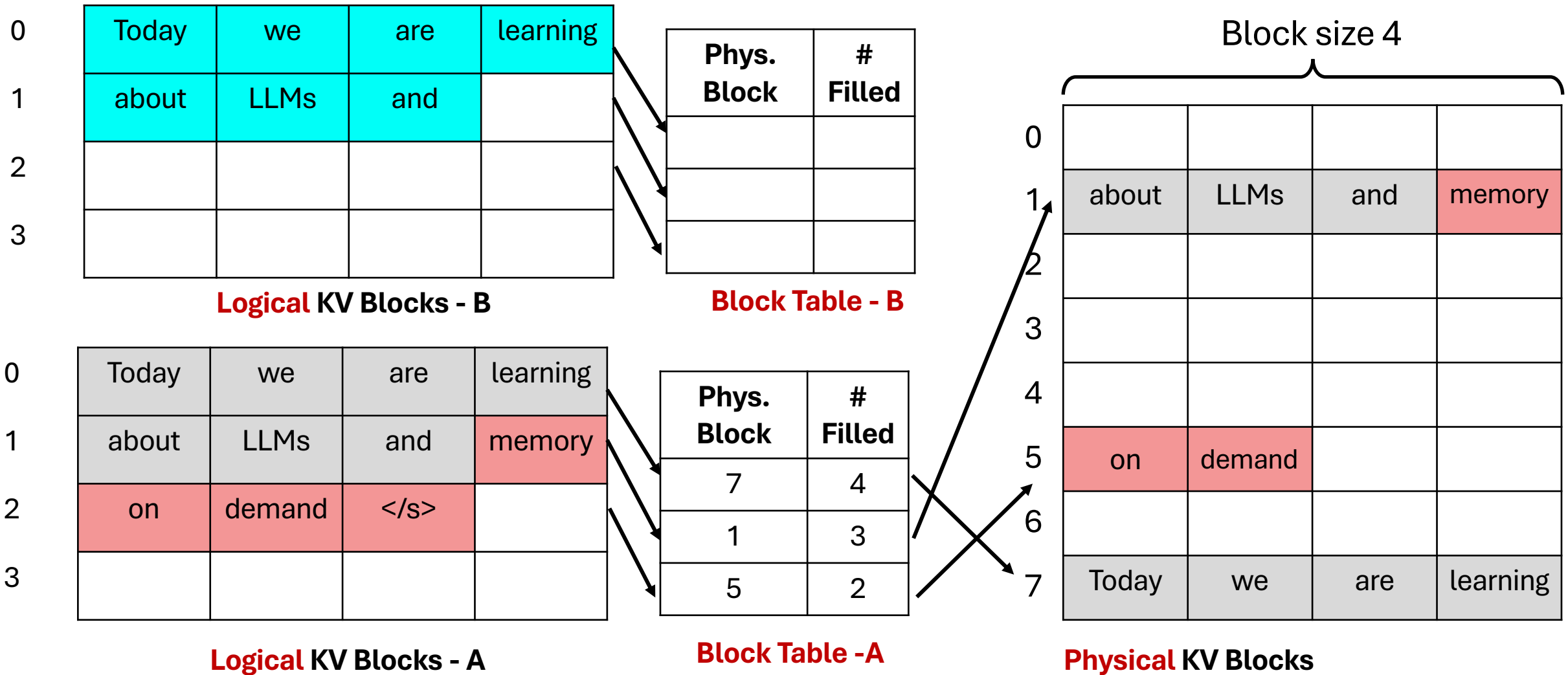
Physical KV Blocks

Prompt A: “Today we are learning about LLMs and”
Completion: “*memory on demand</s>*”

Prompt B: “Today we are learning about LLMs and”
Completion:

Content credits: <https://youtu.be/yVXtLTcdO1Q?si=XO2Dk-VYOShUMH1u>

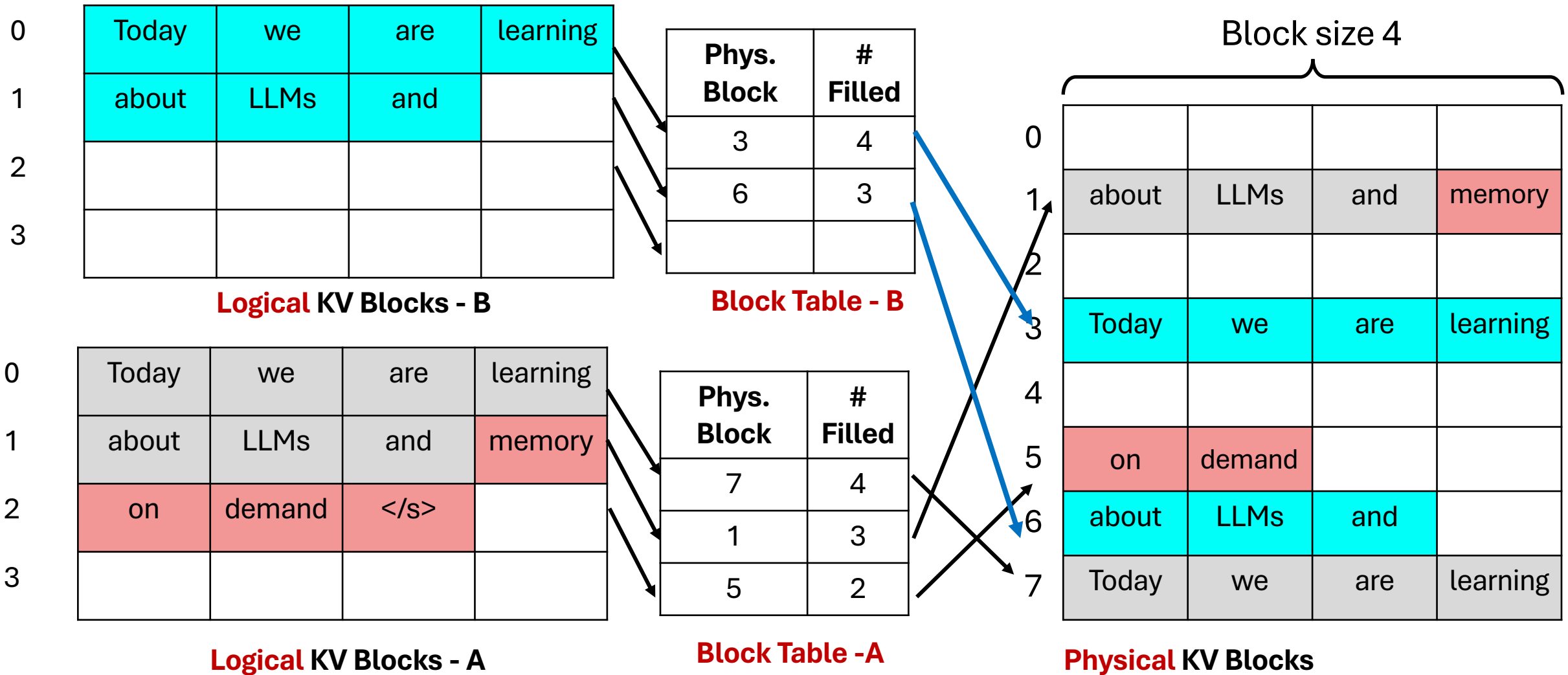




Prompt A: “Today we are learning about LLMs and”
Completion: “*memory on demand </s>*”

Prompt B: “Today we are learning about LLMs and”
Completion:

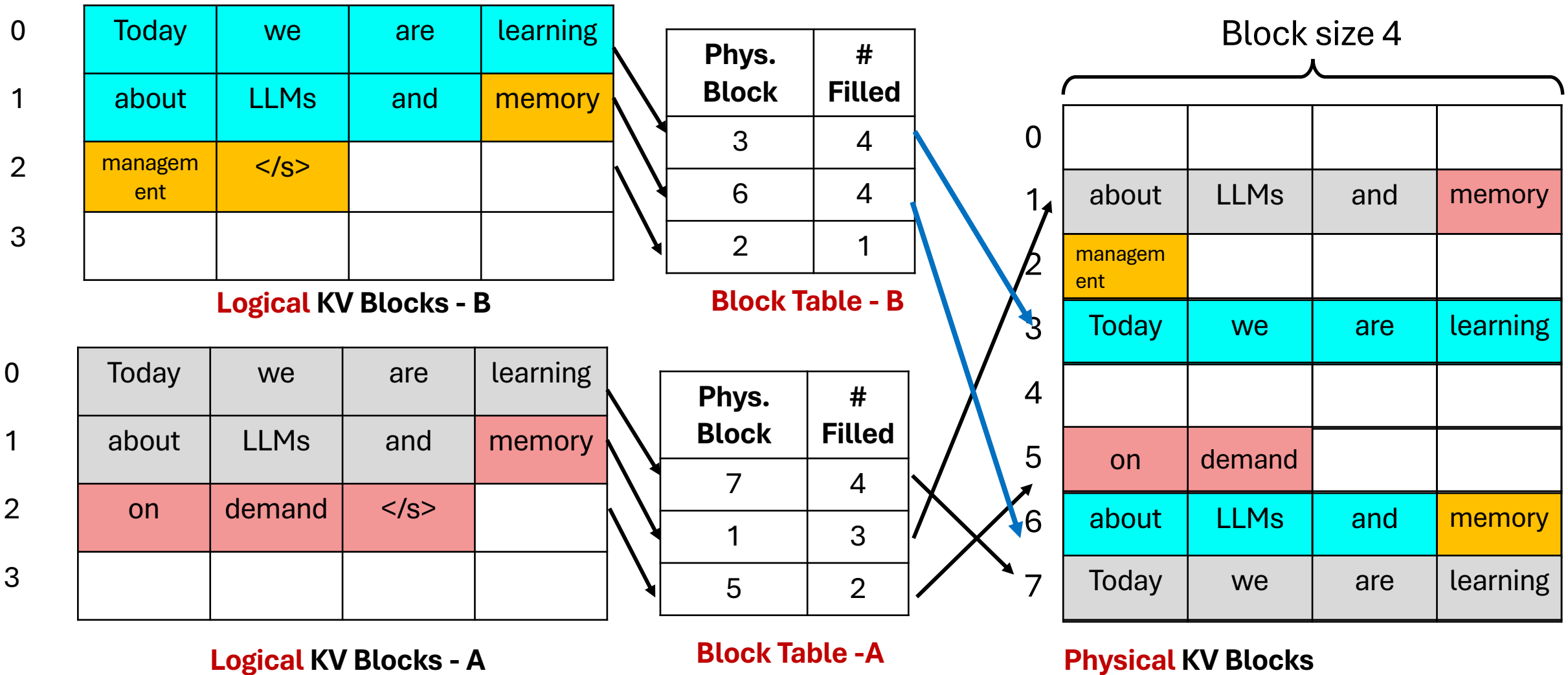




Prompt A: “Today we are learning about LLMs and”
Completion: “*memory on demand </s>*”

Prompt B: “Today we are learning about LLMs and”
Completion: “memory”



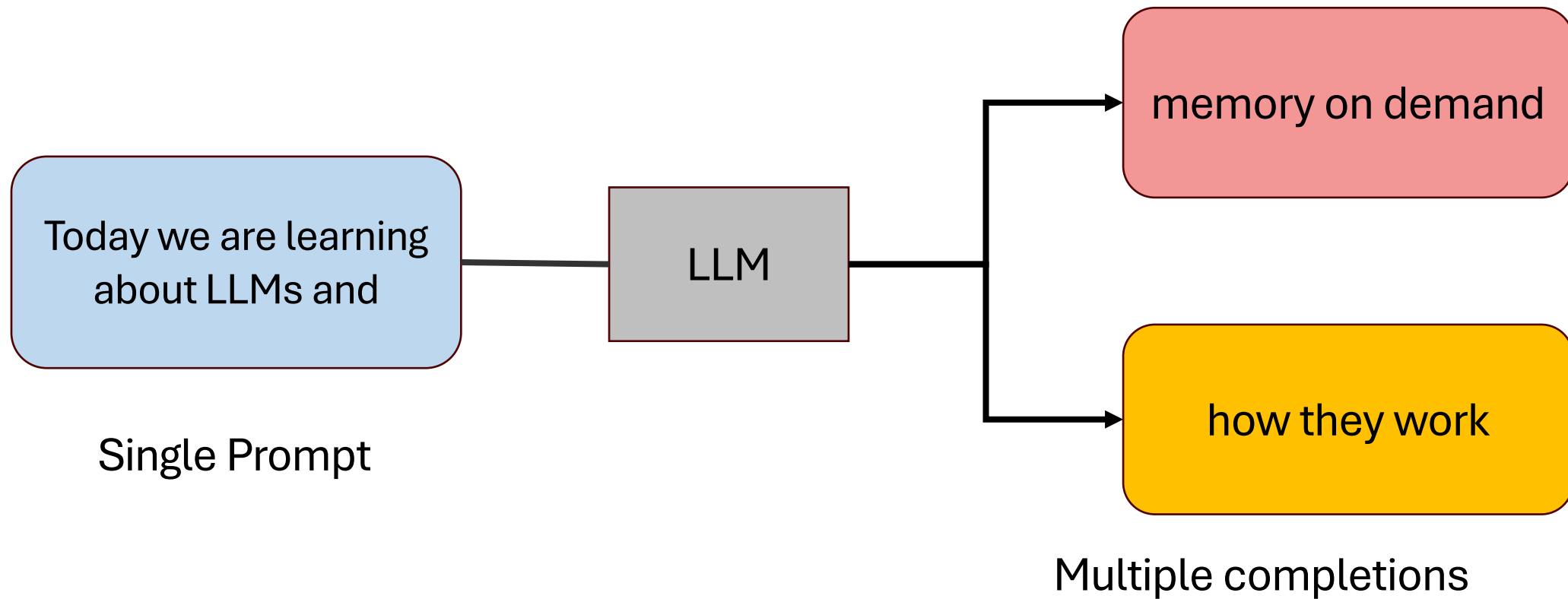


Prompt A: “Today we are learning about LLMs and”
Completion: “*memory on demand </s>*”

Prompt B: “Today we are learning about LLMs and”
Completion: “*memory management </s>*”



Dynamic block mapping enables sharing



Sharing KV blocks in parallel sampling

B.T. for KV A

Phys. Block	# Filled
5	4
7	3

B.T. for KV B

Phys. Block	# Filled
7 5	4
5 7	3

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
1				
2				
3				
4				
5	Today	we	are	learn ng
6				
7	about	LLMs	and	memory

Physical KV Blocks

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
7	3

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	

Logical KV Blocks - A

0				
1				
2				
3				
4				
5	Today	we	are	learn ng
	about	LLMs	and	

Physical KV Blocks

Ref count: 2

Today	we	are	learning
about	LLMs	and	

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
7	3

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
1				
2				
3				
4				
5	Today	we	are	learn ng
	about	LLMs	and	

Physical KV Blocks

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B

Ref counter 2



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
7	3

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
1				
2				
3				
4				
5	Today	we	are	learni ng
6	about	LLMs	and	

Physical KV Blocks

Ref count: 2 → 1

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
7	3

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
	copy-on-write			
2	about	LLMs	and	
3				...
4				
5	Today	we	are	learn ng
6				
7	about	LLMs	and	

Physical KV Blocks

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
7	3

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
1	copy-on-write			
2	about	LLMs	and	memory
3				
4				
5	Today	we	are	learning
6				
7	about	LLMs	and	

Ref count: 1

Physical KV Blocks

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
2	4

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
1	copy-on-write			
2	about	LLMs	and	memory
3				...
4				
5	Today	we	are	learning
6	Ref count: 1			
7	about	LLMs	and	

Physical KV Blocks

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
2	4

Today	we	are	learning
about	LLMs	and	memory

Logical KV Blocks - A

0				
copy-on-write				
2	about	LLMs	and	memory
3				..
4				
5	Today	we	are	learning
6				
Ref count: 1				
	about	LLMs	and	how

Physical KV Blocks

Phys. Block	# Filled
7	4
5	3

Today	we	are	learning
about	LLMs	and	how

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
2	4
0	2

Phys. Block	# Filled
5	4
3	4
4	2

Today	we	are	learning
about	LLMs	and	memory
on	demand		

Logical KV Blocks - A

0	on	demand		
1				
2	about	LLMs	and	memory
3				...
4	they	work		
5	Today	we	are	learning
6				
7	about	LLMs	and	how

Physical KV Blocks

Today	we	are	learning
about	LLMs	and	how
they	work		

Logical KV Blocks - B



Sharing KV blocks in parallel sampling

Phys. Block	# Filled
5	4
2	4
0	2

Phys. Block	# Filled
7	4
5	4
4	2

Today	we	are	learning
about	LLMs	and	memory
on	demand		

Logical KV Blocks - A

0	on	demand		
1				
2	about	LLMs	and	memory
3				.
4	they	work		
5	Today	we	are	learning
6				
7	about	LLMs	and	how

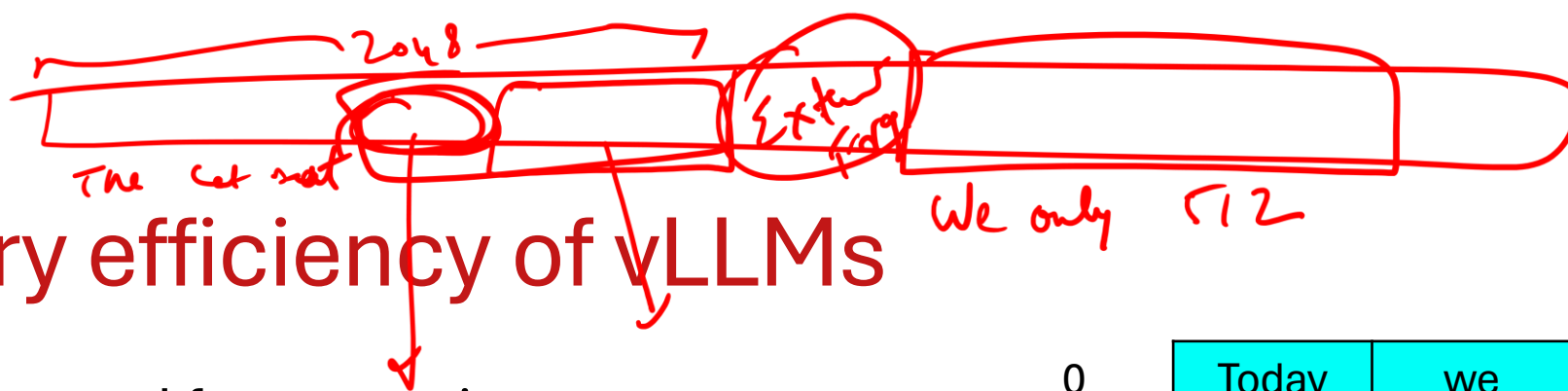
Ref count: 2

Physical KV Blocks

Today	we	are	learning
about	LLMs	and	how
they	work		

Logical KV Blocks - B





Memory efficiency of vLLMs

✓ Minimal internal fragmentation

- Only happens at the last block of a sequence

- **# wasted tokens / seq < block size**

- Sequence: $O(100)$ or $O(1000)$ tokens

- Block size: 16 or 32 tokens

✓ No external fragmentation

✓ On average, wasted space < 4% of KV cache

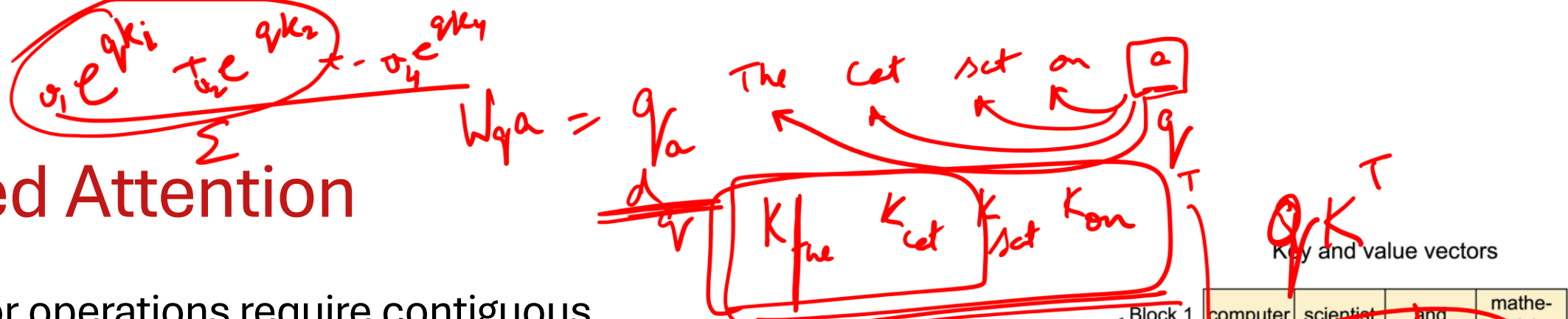
✓ 3-5x improved memory utilization!

0	Today	we	are	learning
1	about	LLMs	and	memory
2	managem ent			
3				

Internal fragmentation

Paged Attention

- Tensor operations require contiguous memory
- How to compute attention softmax across fragmented memory?
- Paged Attention!



$$\text{softmax}([A_1, A_2]) = [\alpha \text{softmax}(A_1), \beta \text{softmax}(A_2)]$$

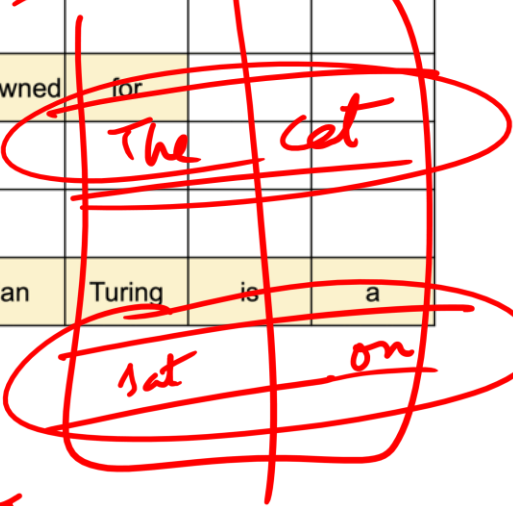
$$\text{softmax}([A_1, A_2]) \begin{bmatrix} V_1 \\ V_2 \end{bmatrix} = \alpha \text{softmax}(A_1) * V_1 + \beta \text{softmax}(A_2) * V_2$$

$$\alpha = \frac{\sum_i q_i r_i}{\sum_i q_i r_i}$$

Handwritten notes: $\sum_i p_i (V_i)$ and $(K_{the cat sat on})$.

Key and value vectors

Block 1	computer	scientist	and	mathe- mician
Block 2	renowned	for		
Block 0	Alan	Turing	is	a



How vLLM & Paged Attention results in efficient inference?

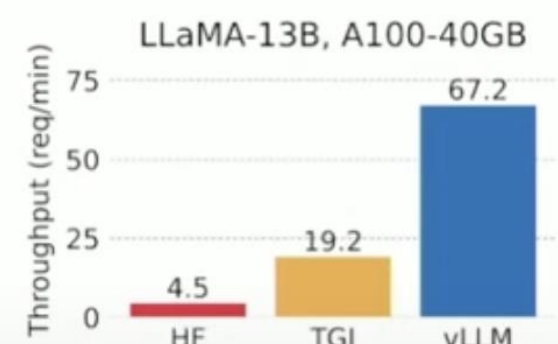
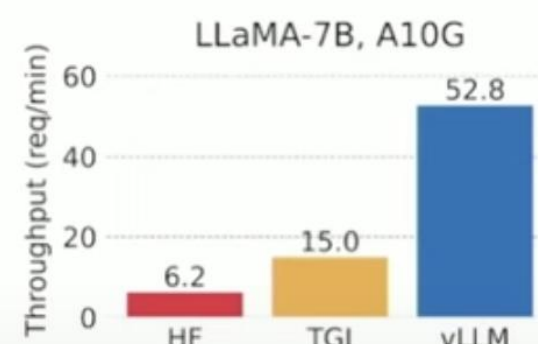
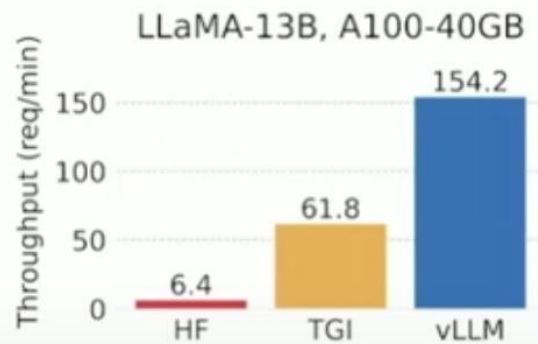
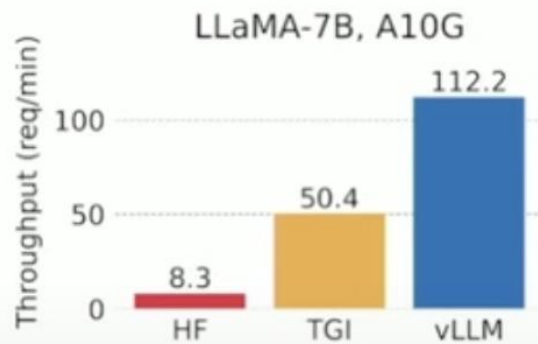
Reduce memory fragmentation with paging

Further reduce memory usage with sharing



Comparison with HuggingFace and TGI (2023)

- Up to **24x** higher throughput than HuggingFace (HF)
- Up to **3.5x** higher throughput than Text Generation Inference (TGI)

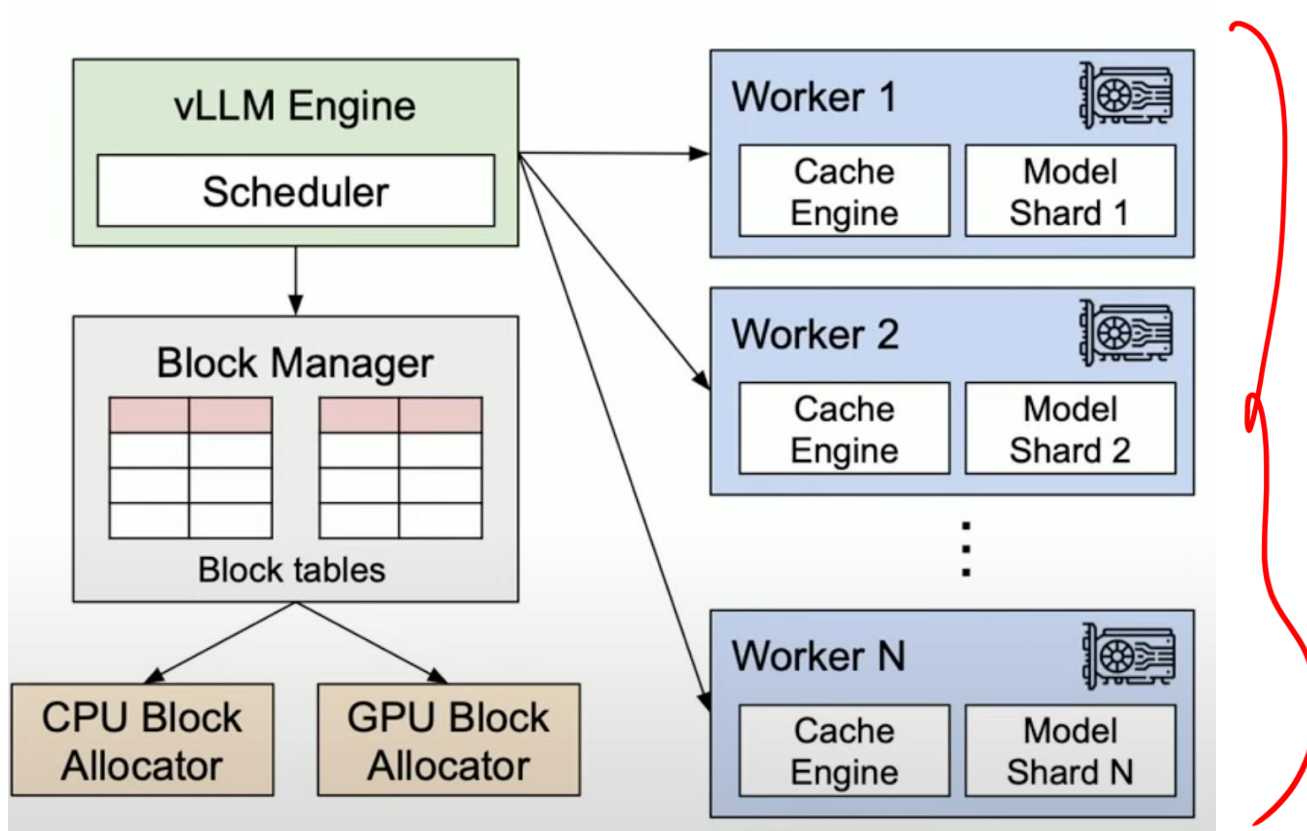


Serving throughput when each request asks for 1 output completion.

Serving throughput when each request asks for 3 output completions.



System Architecture and Implementation



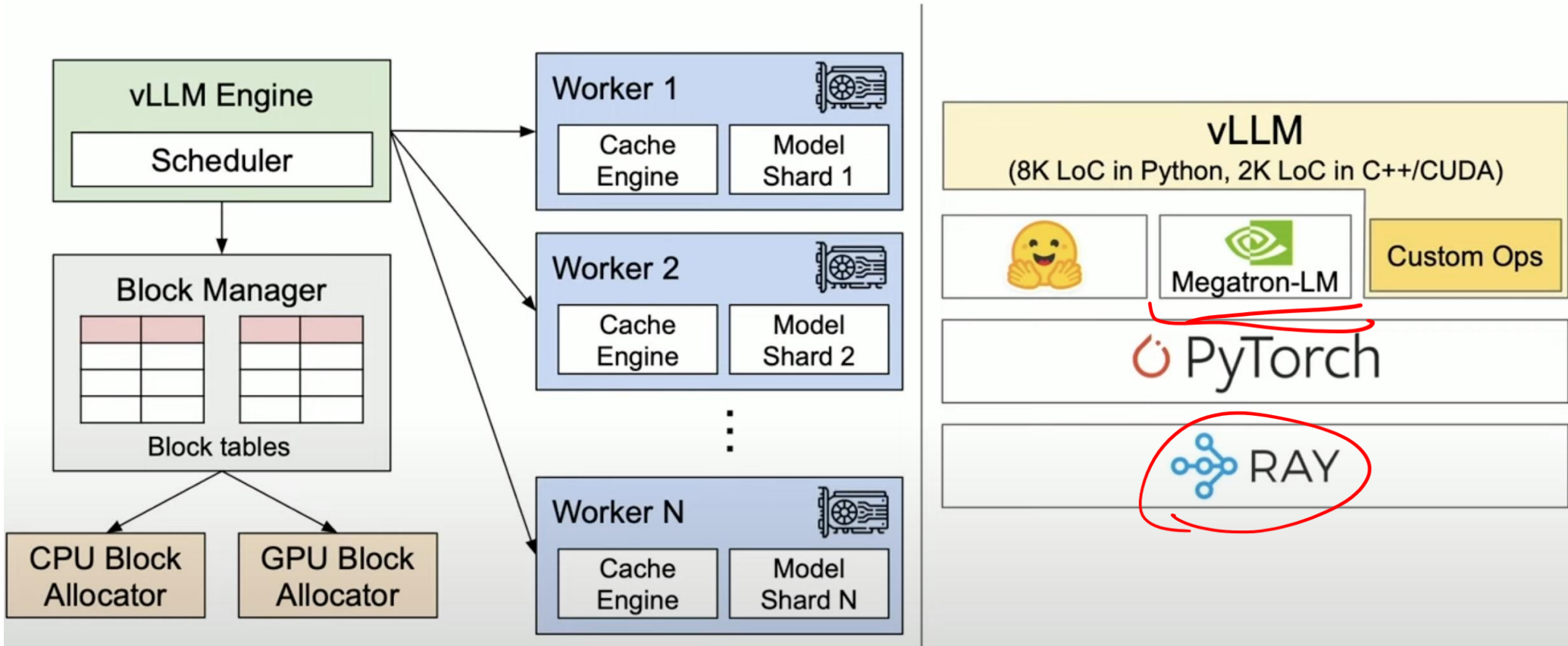
End to end llm serving engine

3 components –

- A frontend
- A distributed model executor
- A scheduler

A centralized engine that manages block table





Till now...

- **KV caching** – avoids re-computation of Keys and Value matrices
- **Paged Attention and vLLM** - efficient memory management
- Can we speed up attention computation?
- **Flash Attention?**

Redundant computation ✓
GPU v cache
✓ Efficient memory management
Speedy computation
Speculative decoding

