

# Large Language Models

## Advanced Attention Mechanisms - I

ELL881 · AIL821

Sourish Dasgupta

Assistant Professor, DA-IICT, Gandhinagar

<https://www.daiict.ac.in/faculty/sourish-dasgupta>

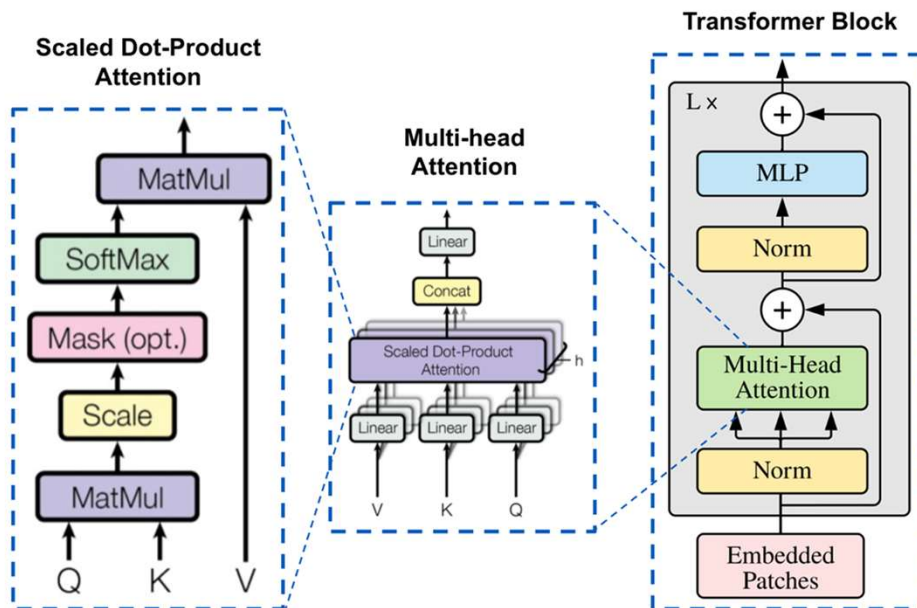


Semester 1, 2024-2025

Year: 2017, NeurIPS



# Self Attention



similarity

contextual embedding

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

$$\begin{matrix} \text{Q} \\ (6, 4096) \end{matrix} \times \begin{matrix} \text{K}^T \\ (4096, 6) \end{matrix} = \frac{\sqrt{4096}}{\sqrt{4096}}$$

	THE	CAT	IS	ON	A	CHAIR
THE	0.268	0.119	0.134	0.148	0.179	0.152
CAT	0.124	0.278	0.201	0.128	0.154	0.115
IS	0.147	0.132	0.262	0.097	0.218	0.145
ON	0.210	0.128	0.206	0.212	0.119	0.125
A	0.146	0.158	0.152	0.143	0.227	0.174
CHAIR	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

*Stylized handwritten text:*

Stylized (loves. John) % ? %  
 John % loves % IL made in NY %  
 % % % %  
 v : values.



Year: 2017, NeurIPS



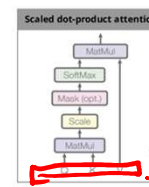
# Causal (Forward Masked) Attention

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

$$\begin{matrix} \boxed{Q} & \times & \boxed{K^T} \\ (6, 4096) & & (4096, 6) \end{matrix} = \frac{\quad}{\sqrt{4096}}$$

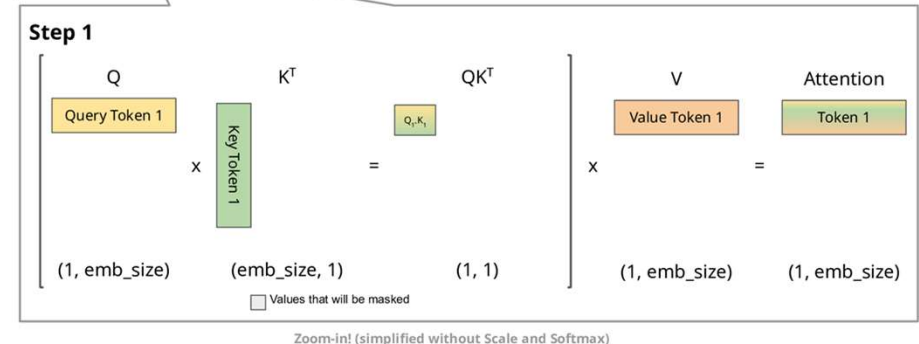
	THE	CAT	IS	ON	A	CHAIR
THE	0.268	—∞	—∞	—∞	—∞	—∞
CAT	0.124	0.278	—∞	—∞	—∞	—∞
IS	0.147	0.132	0.262	—∞	—∞	—∞
ON	0.210	0.128	0.206	0.212	—∞	—∞
A	0.146	0.158	0.152	0.143	0.227	—∞
CHAIR	0.195	0.114	0.203	0.103	0.157	0.229

(6, 6)

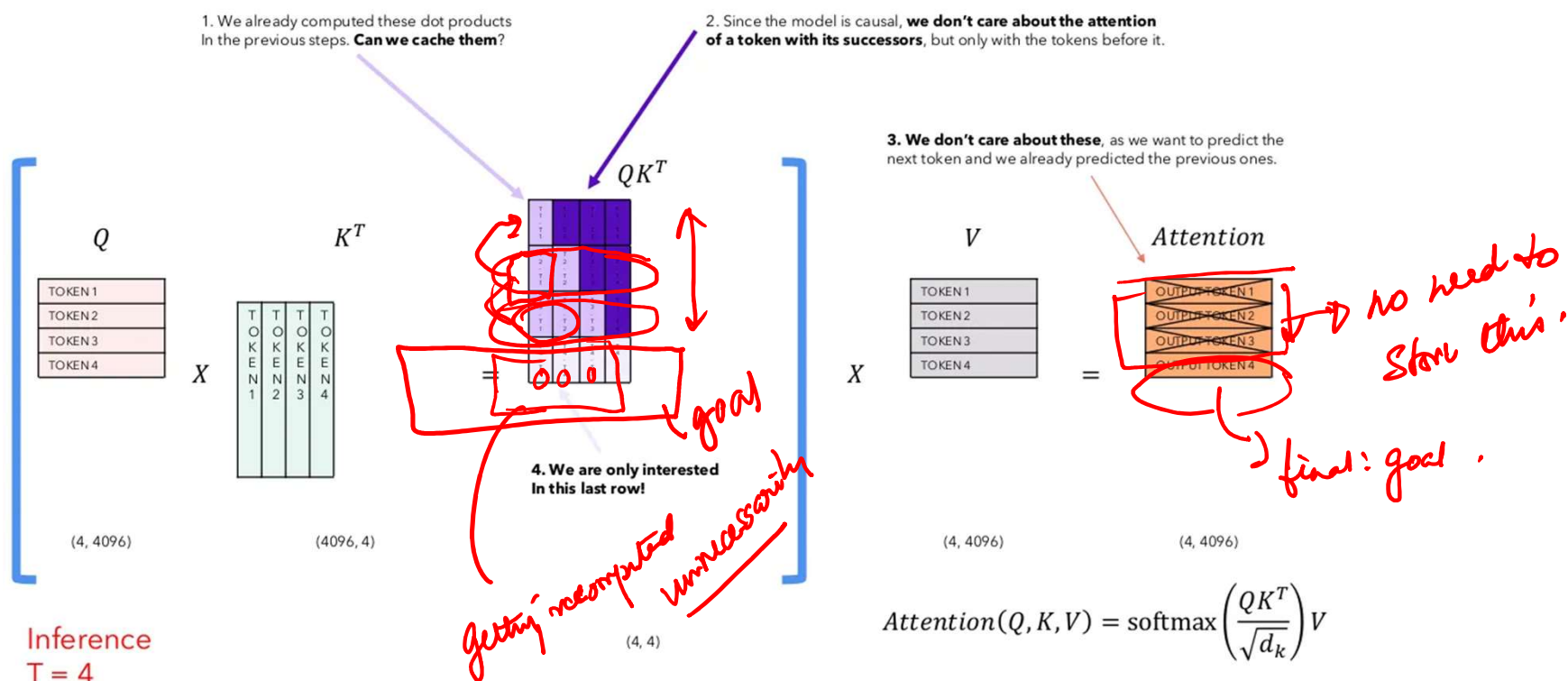


*N: Sequence length*  
*d: dimension*  
*Similarity (dot products)*  
*softmax*  
*x V*

Time:  $O(N^2 \cdot d) + O(N^2) + O(N^2 \cdot d) = O(N^2 \cdot d)$   
 Space:  $O(3 \cdot N \cdot d) + O(N^2) + O(N \cdot d) = O(N^2 + N \cdot d)$   
 $\approx (3 \times N \times d + N^2) \times \text{Size of a float}$

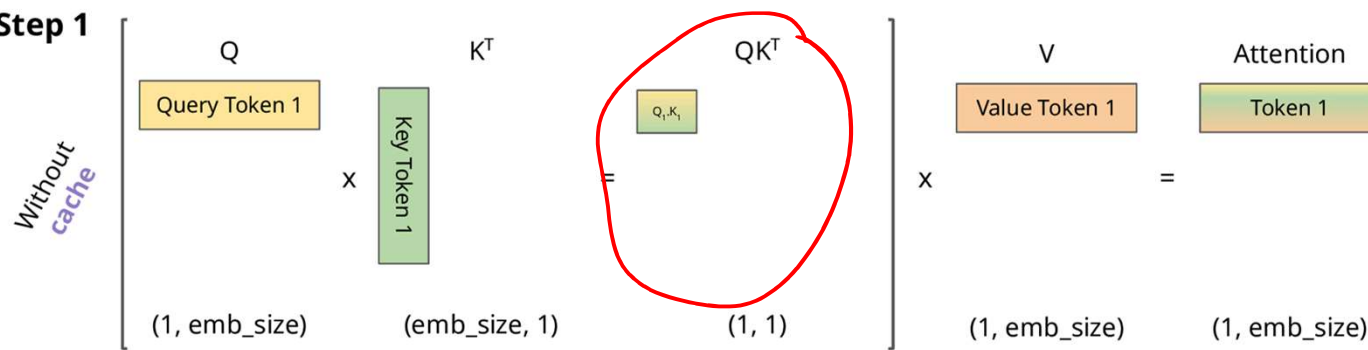


# Why do we need to do better?



# KV Cache based (Forward Masked) Attention

Step 1

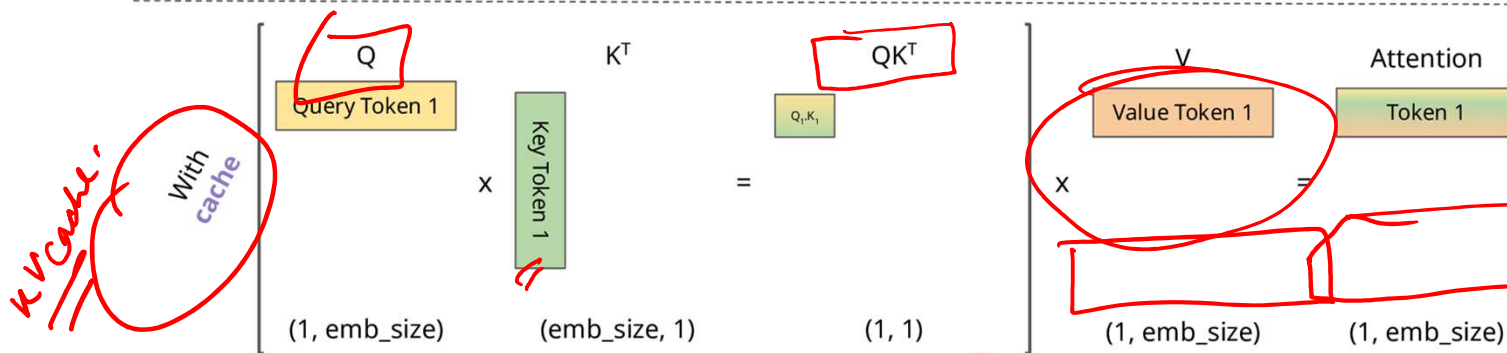


KV Cache Storage =  $O(N \times d)$

$(2 \times N \times d) \times \text{Size of a float}$

vs.

$(3 \times N \times d + N^2) \times \text{Size of a float}$



GOAL -

Source: <https://medium.com/@joaolages/kv-caching-explained-276520203249>

□ Values that will be masked    ■ Values that will be taken from cache



Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

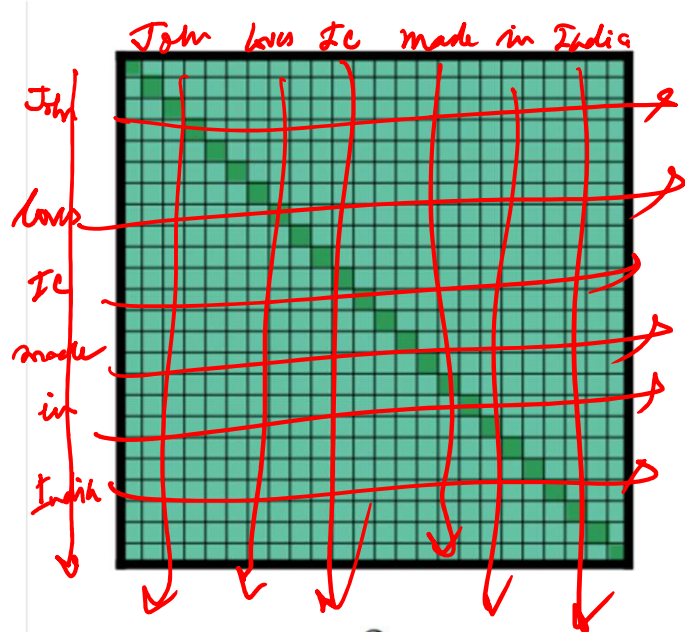
Year: 2020, Arxiv



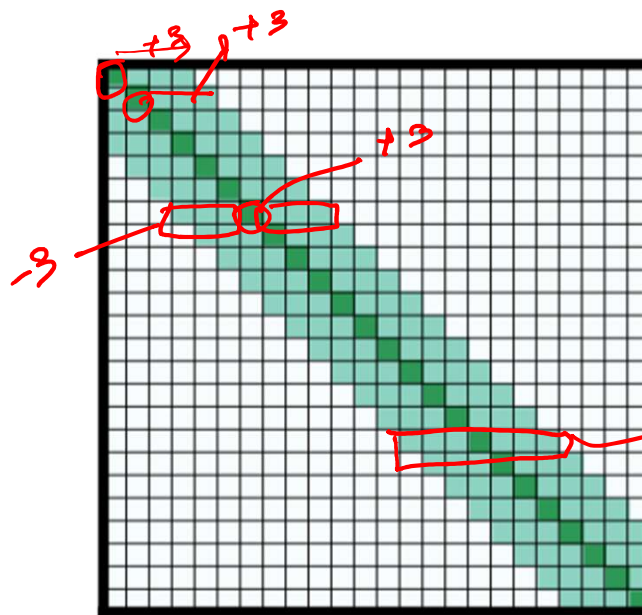
# Sliding Window Attention

Any LM (not just causal)

Window size = W (4, 5, 6)



(a) Full  $n^2$  attention



(b) Sliding window attention

Context window -

W = 7

~~apply to causal LM.~~

get information from the future.  
→ attention.



Year: 2020, Arxiv



# Sliding Window Attention

$$\text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V$$

Diagram illustrating the Sliding Window Attention mechanism. The input matrix  $Q$  (size  $(6, 4096)$ ) is multiplied by the input matrix  $K^T$  (size  $(4096, 6)$ ). The result is then passed through a softmax operation, scaled by  $\sqrt{4096}$ , to produce the attention weights.

	THE	CAT	IS	ON	A	CHAIR
THE	1.0	0	0	0	0	0
CAT	0.461	0.538	0	0	0	0
IS	0.3219	0.317	0.361	0	0	0
ON	0	0.316	0.341	0.343	0	0
A	0	0	0.326	0.323	0.351	0
CHAIR	0	0	0	0.313	0.331	0.356

$w = 3$



Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

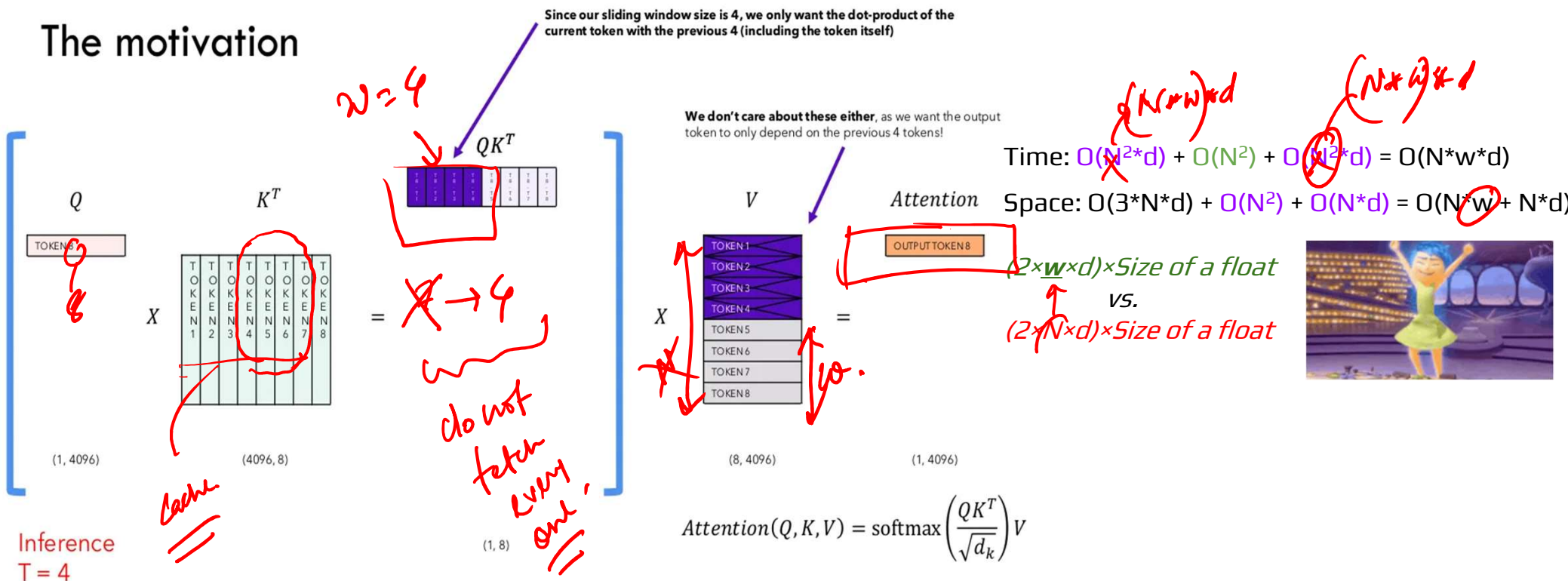


Year: 2020, Arxiv



# What happens to the KV Cache?

## The motivation



Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

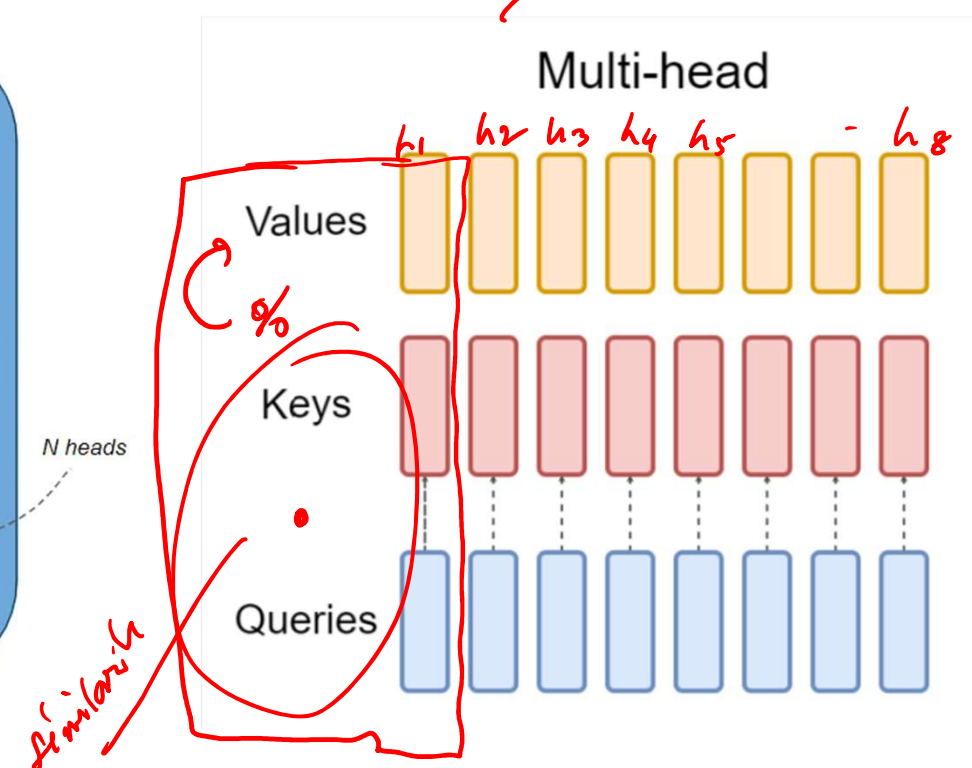




Google Brain Team

Make machines intelligent. Improve people's lives.

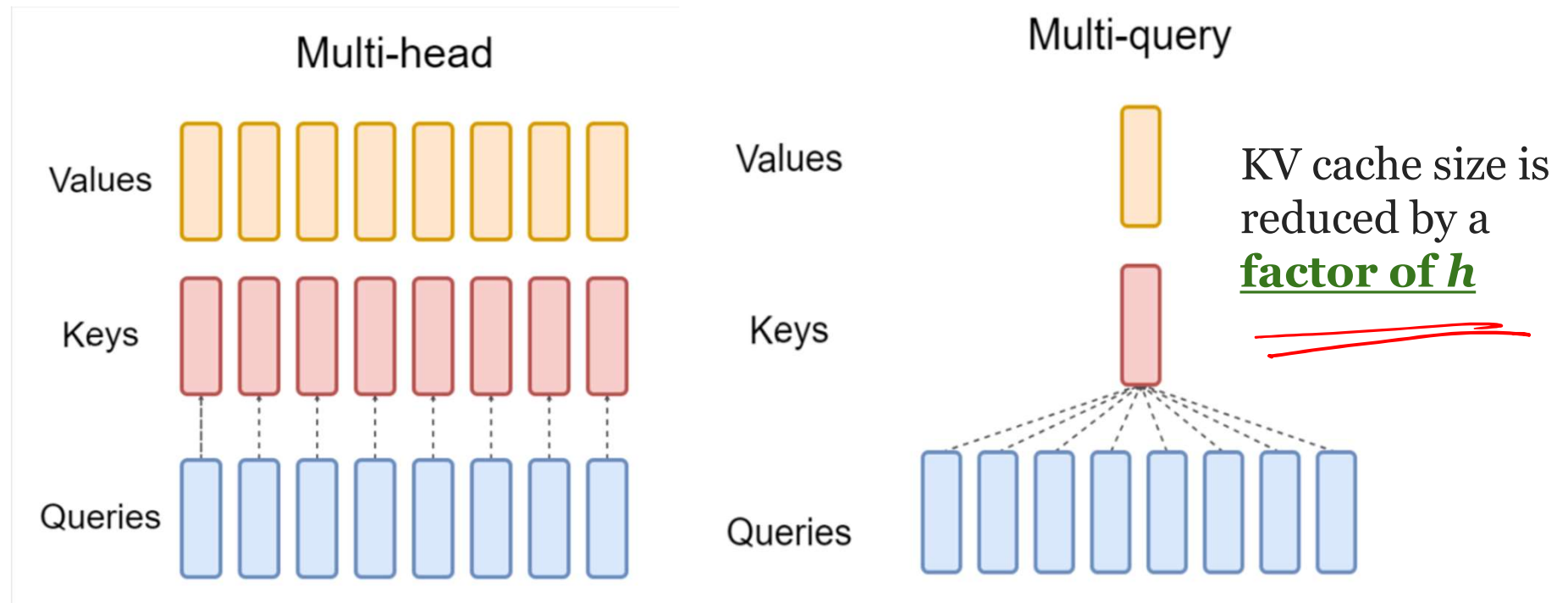
→ # heads:  $h$



Year: 2019, arxiv



# Multi-Query Attention (MQA)



Sourish Dasgupta

LLMs: Advanced Attention  
Mechanisms

# Do we lose out on something?

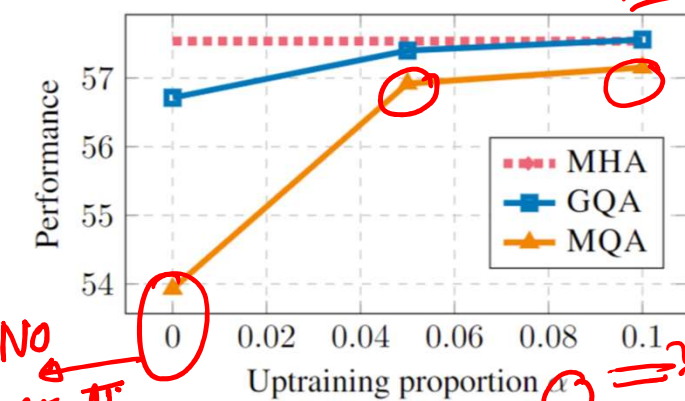
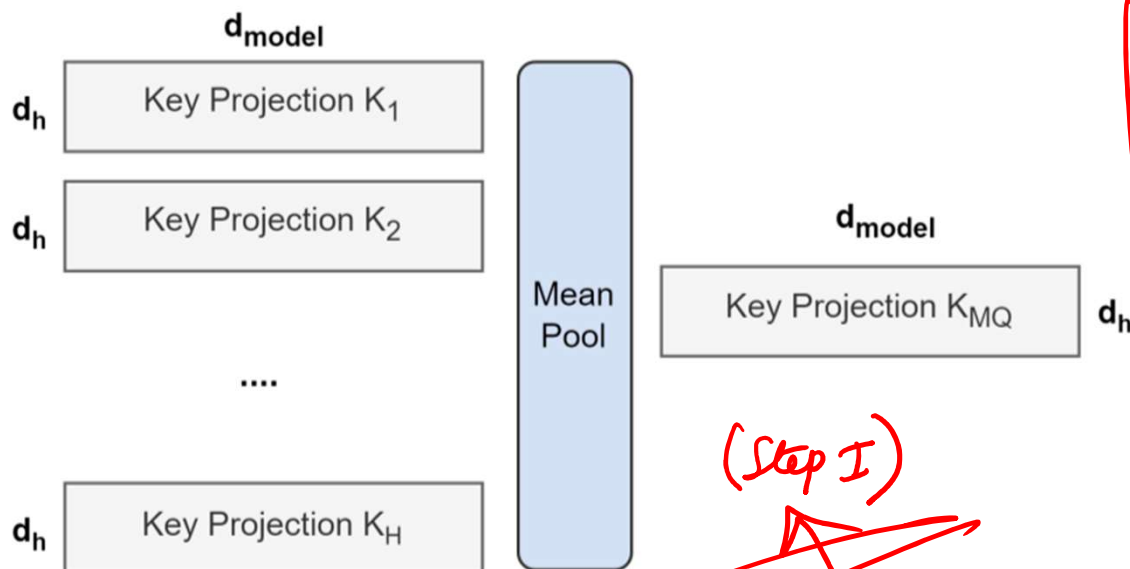
- Decline in performance quality
- Training instability



Year: 2023; ICLR



# Uptraining: Converting MHA to MQA



getting trained }  $k$  epochs -  
"further"  
[Step 2] - pre-train  
 $k' \rightarrow k$   
No step II.  
STOP the train @  $k'$  epoch.  
 $\alpha = \frac{k-k'}{k}$



# What can still go wrong?

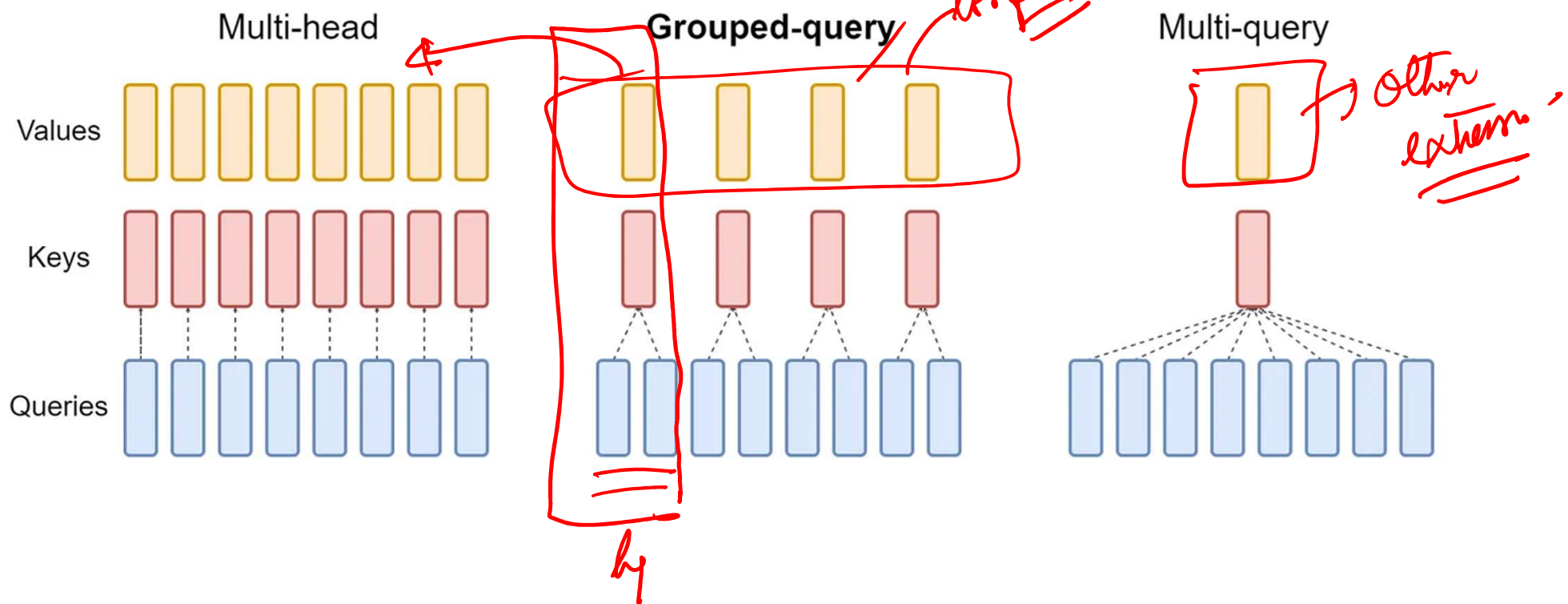
- Decline in performance quality
- ~~Training instability~~



Year: 2023; EMNLP



# Grouped Query Attention



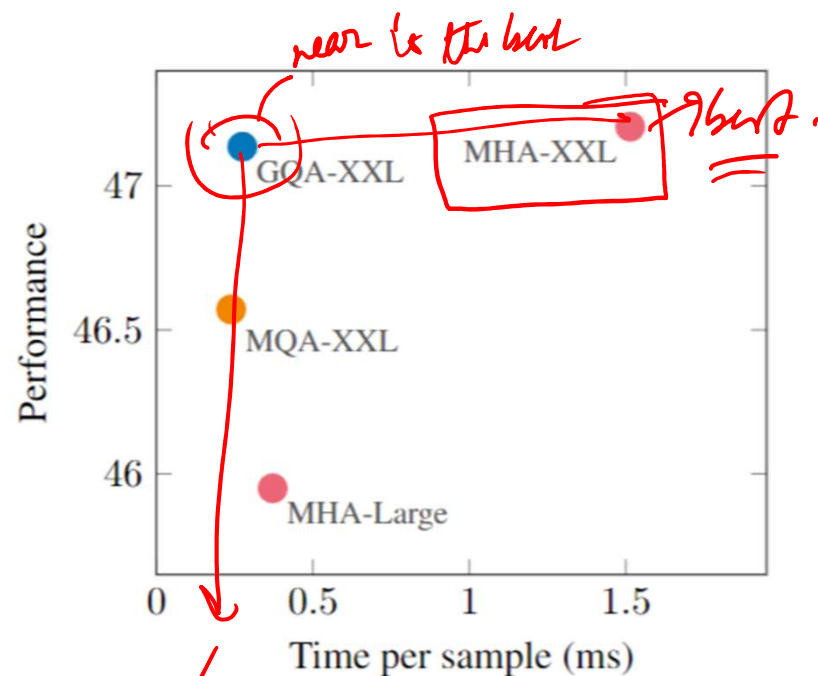
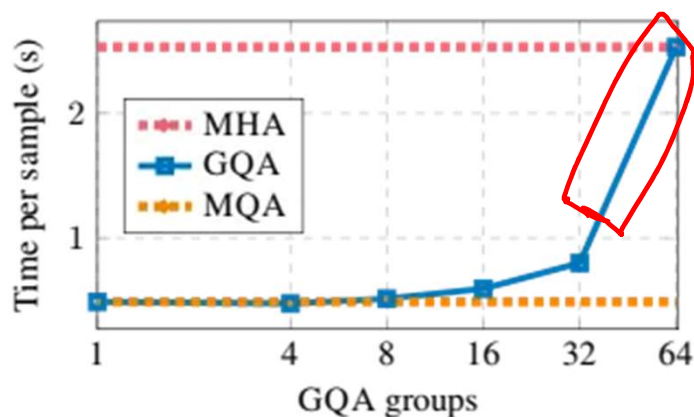
Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

# What did we gain?

*inference Time*

Model	$T_{\text{infer}}$	Average	CNN	arXiv	PubMed	MediaSum	MultiNews	WMT	TriviaQA
	s		$R_1$	$R_1$	$R_1$	$R_1$	$R_1$	BLEU	F1
MHA-Large	0.37	46.0	42.9	44.6	46.2	35.5	46.6	27.7	78.2
MHA-XXL	1.51	47.2	43.8	45.6	47.5	36.4	46.9	28.4	81.9
MQA-XXL	0.24	46.6	43.0	45.0	46.9	36.1	46.5	28.5	81.3
GQA-8-XXL	0.28	47.1	43.5	45.4	47.7	36.3	47.2	28.4	81.6



*lot less time*





## So are we all set? Key

- GQA/MQA Aim: *To reduce the need for storing a large amount of KV cache*
- LLM server can handle more requests, larger batch sizes and increased throughput
- *Cannot significantly reduce the computational load* ~~/// BUT~~
- *Quality degradation remains*



# Large Language Models

## Advanced Attention Mechanisms - II

ELL881 · AIL821

Sourish Dasgupta

Assistant Professor, DA-IICT, Gandhinagar

<https://www.daiict.ac.in/faculty/sourish-dasgupta>



Semester 1, 2024-2025

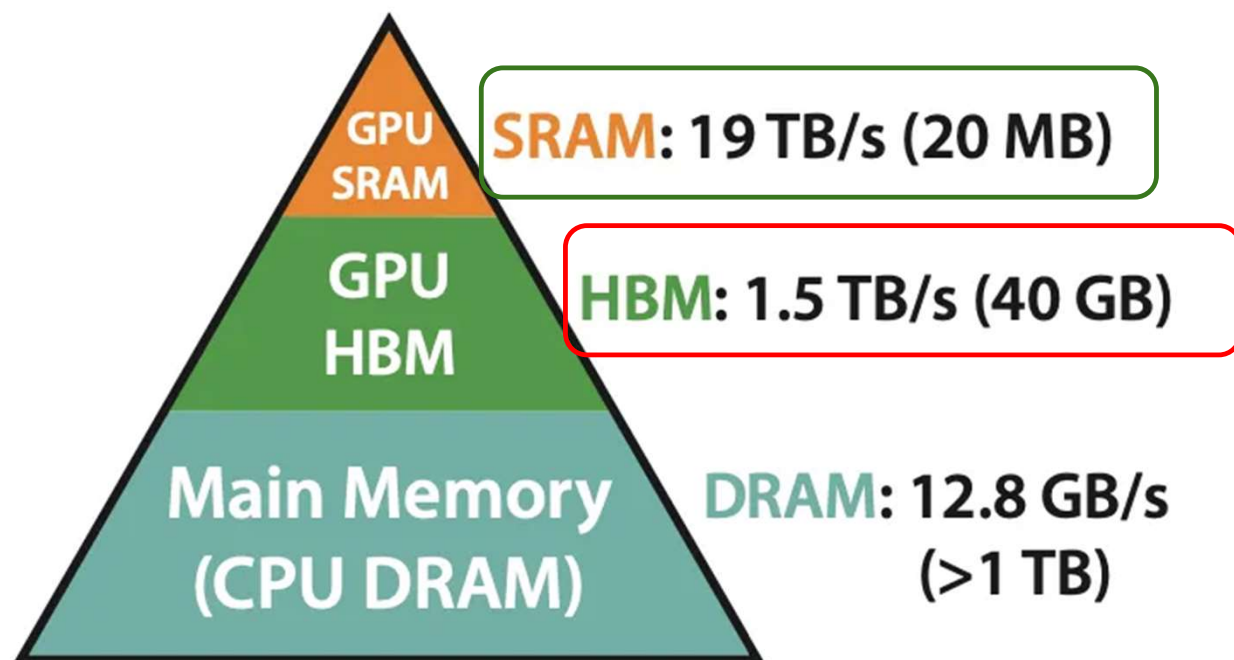
# Can we optimize without performance degradation?



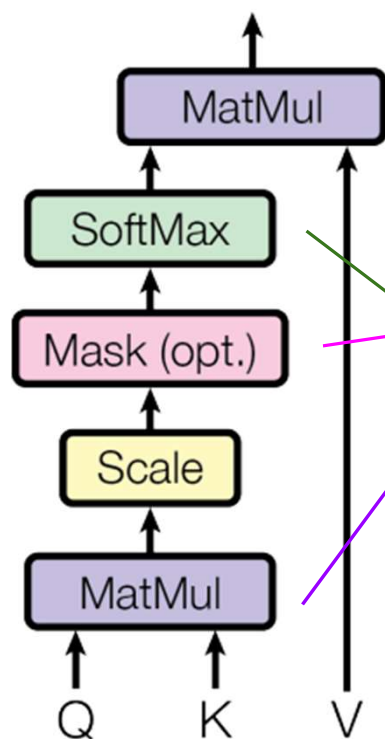
Sourish Dasgupta

LLMs: Advanced Attention  
Mechanisms

## A bit more about the GPU



# What was happening so far:



## 1. Matmul\_op (Q,K)

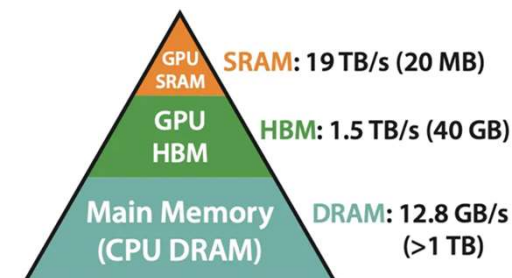
- Read Q,K to SRAM (read-op)
- Compute matmul  $A=Q \times K$  (compute-op)
- Write A to HBM (write-op)

## 2. Mask\_op

- Read A to SRAM (read-op)
- Mask A into A' (compute-op)
- Write A' to HBM (write-op)

## 3. Softmax\_op

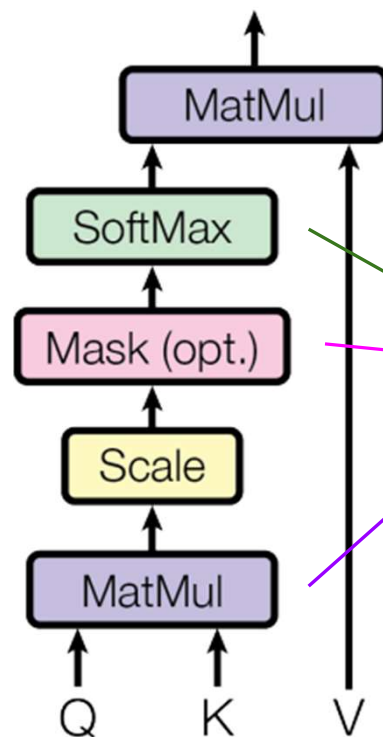
- Read A' to SRAM (read-op)
- Softmax A' into A'' (compute-op)
- Write A'' to HBM (write-op)



Year: 2022; NeurIPS



# The magic: Fused Kernel (GPU Operations)!



## Flash Attention

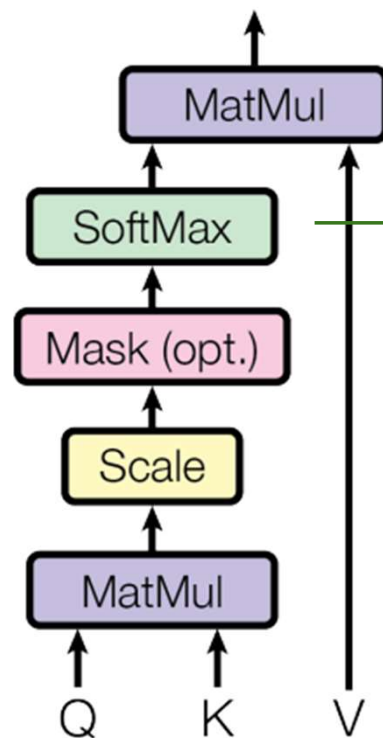
1. Read **Q, K** to **SRAM**
2. Compute  $A = Q \times K$
3. Mask  $A$  into  $A'$
4. Softmax  $A'$  into  $A''$
5. Write  $A''$  to **HBM**



Year: 2022; NeurIPS



# The magic does not end here! More optimization



$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$

Gets  
computed  
for every row  
- *Problem!*

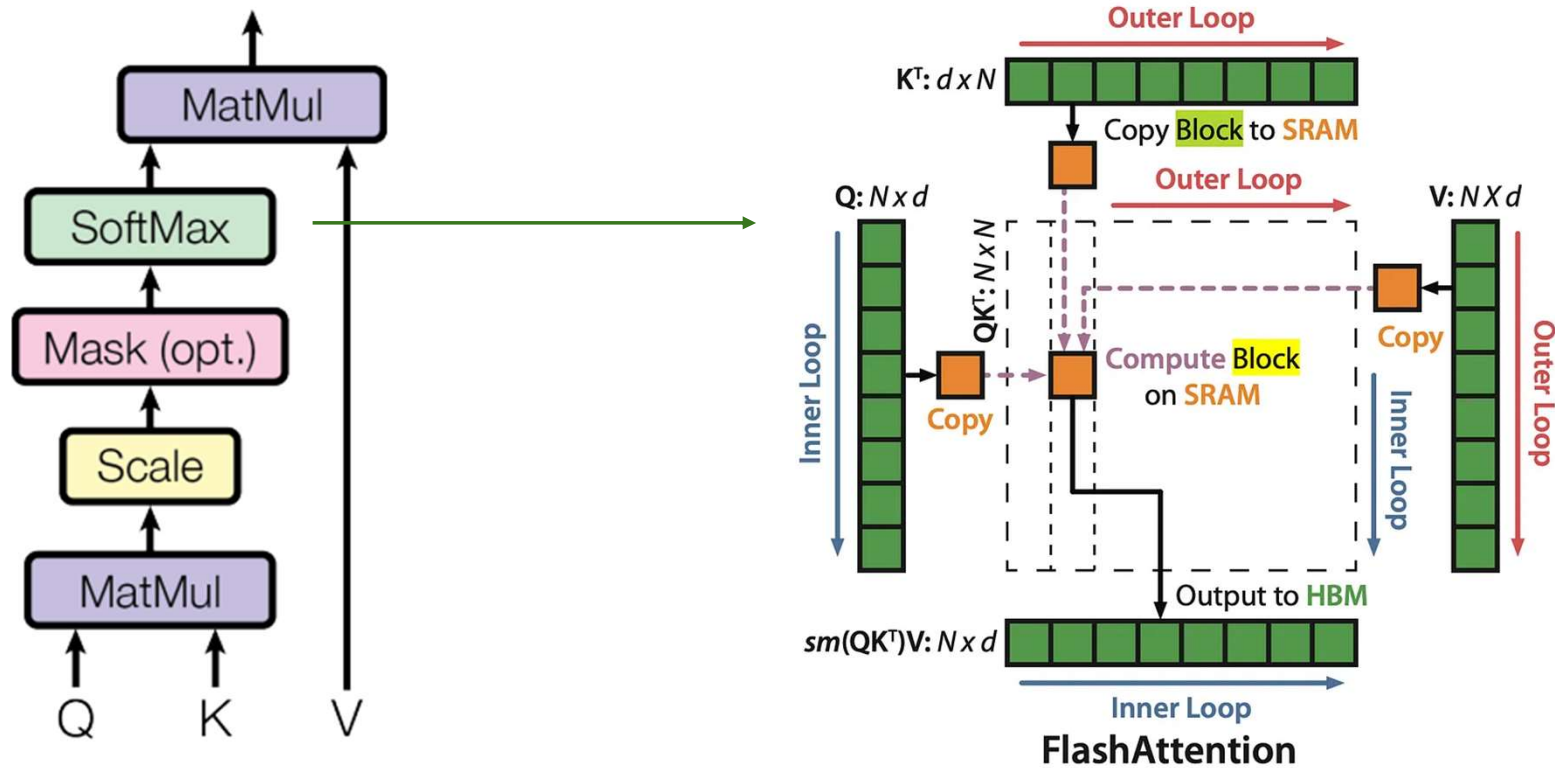




Year: 2022; NeurIPS



# The magic does not end here! Tiling



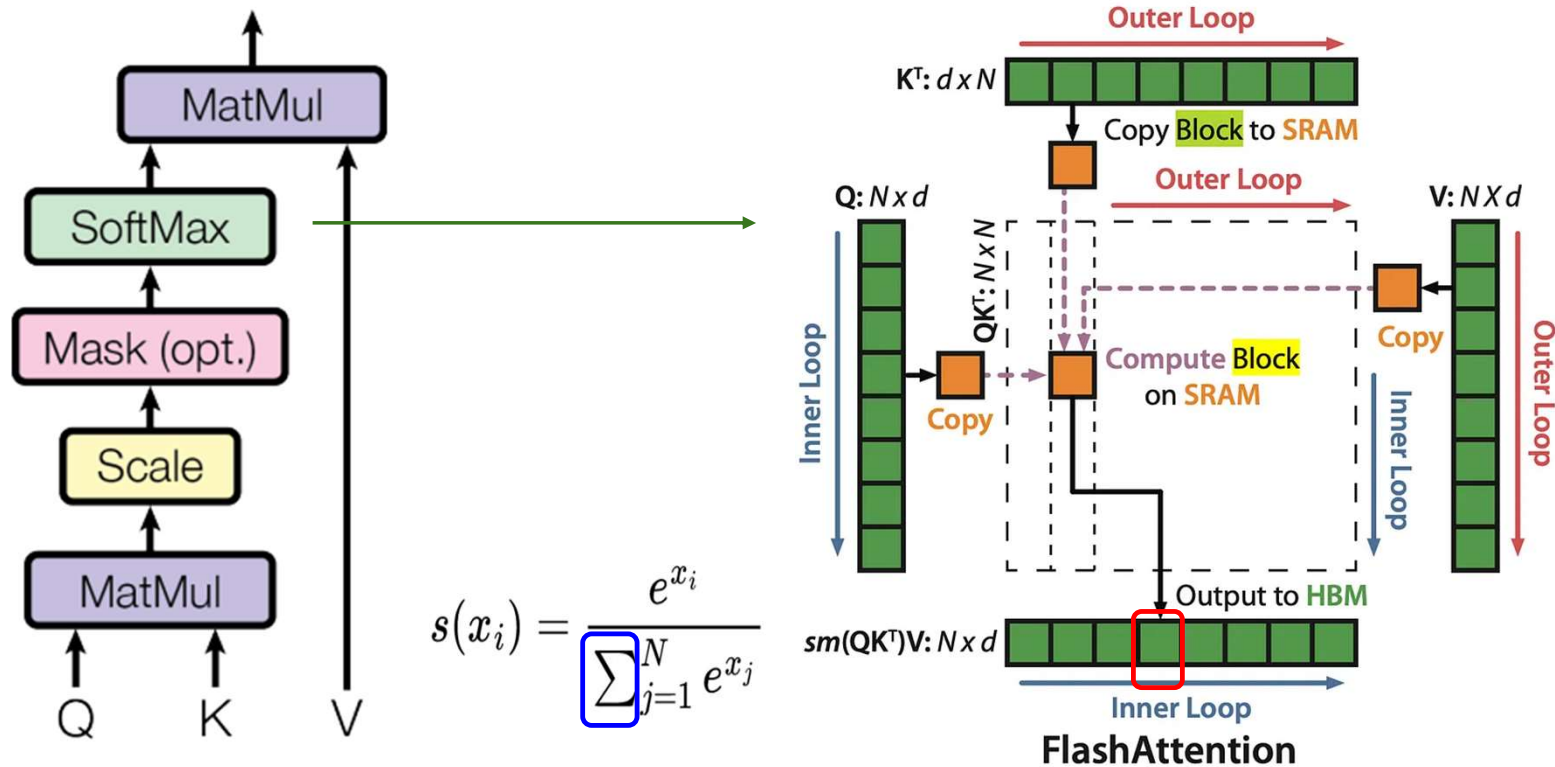
Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

Year: 2022; NeurIPS



# Does the story end here? What's the problem?



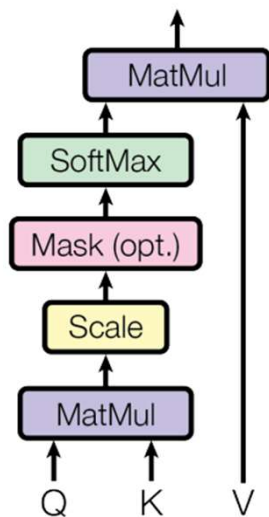
Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

Year: 2022; NeurIPS



# The softmax denominator problem



$$Q = [1] \quad K = [1, 2, 3] \quad V = [2, 4, 8]$$

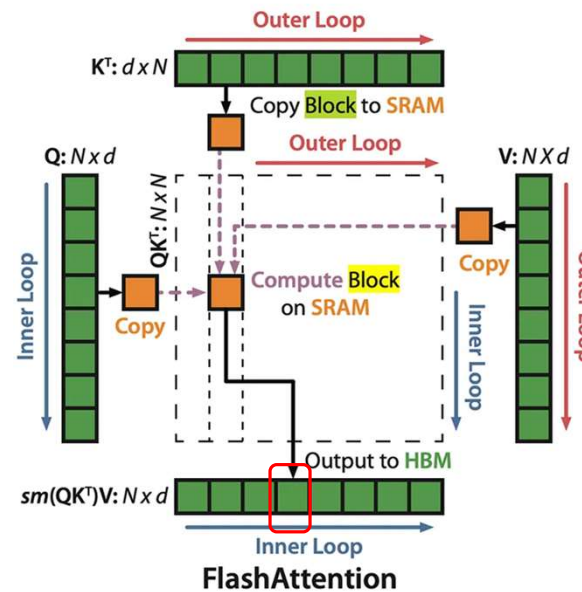
$$A = QK^T = [1, 2, 3]$$

$$V = [2, 4, 8]$$

$$O = \text{softmax}(A)V$$

$$O = \frac{N}{D} = \frac{2e^1 + 4e^2 + 8e^3}{e^1 + e^2 + e^3}$$

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^N e^{x_j}}$$



At  $i = 0$

$$D_b = e^1$$

$$N_b = 2e^1$$

$$-O = \frac{1}{e^1} [0 + 2e^1]$$



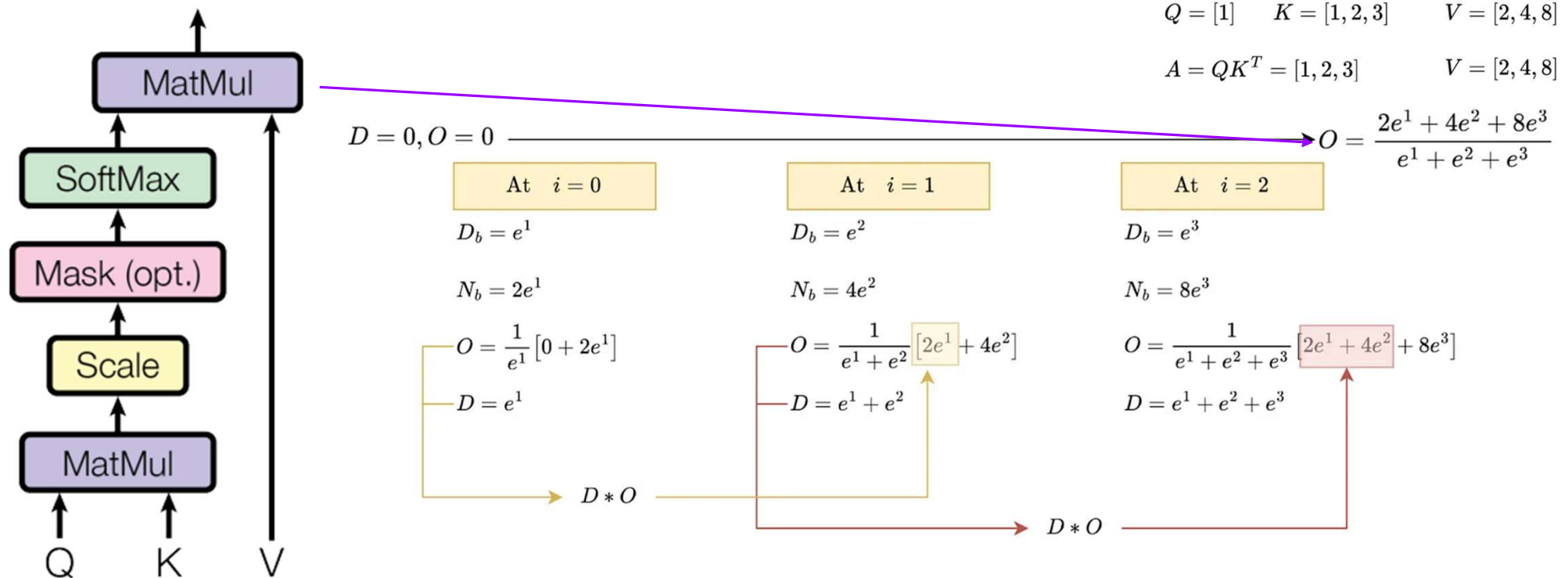
Sourish Dasgupta

LLMs: Advanced Attention Mechanisms

Year: 2022; NeurIPS



# Summary Statistics - the final touch!



**LCS**  
LABORATORY FOR  
COMPUTATIONAL SOCIAL SYSTEMS

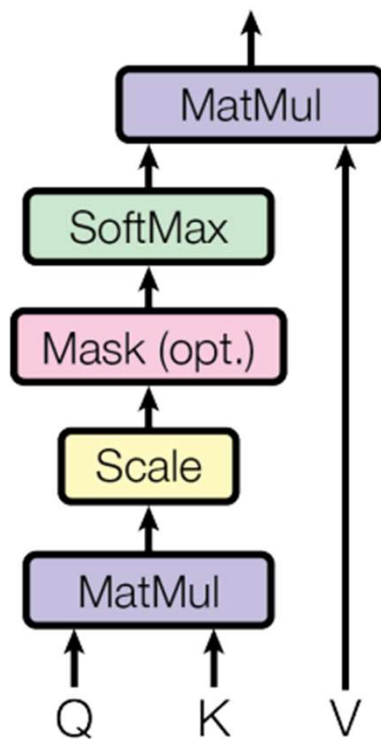
Sourish Dasgupta

LLMs: Advanced Attention  
Mechanisms

Year: 2022; NeurIPS



## Summary Statistics - the final touch!



$$Q = [1] \quad K = [1, 2, 3] \quad V = [2, 4, 8]$$

$$A = QK^T = [1, 2, 3] \quad V = [2, 4, 8]$$

$$O = \text{softmax}(A)V$$

$$O = \frac{N}{D} = \frac{2e^1 + 4e^2 + 8e^3}{e^1 + e^2 + e^3}$$

$$D = 0, O = 0$$

# Treat each element as a block,  
# so we have three blocks  
for i in range(3):

$$D_b = \exp(Q[i] \times K[i])$$

$$N_b = V[i] * \exp(Q[i] \times K[i])$$

$$O = \frac{1}{D + D_b} [D * O + N_b]$$

$$D = D + D_b$$



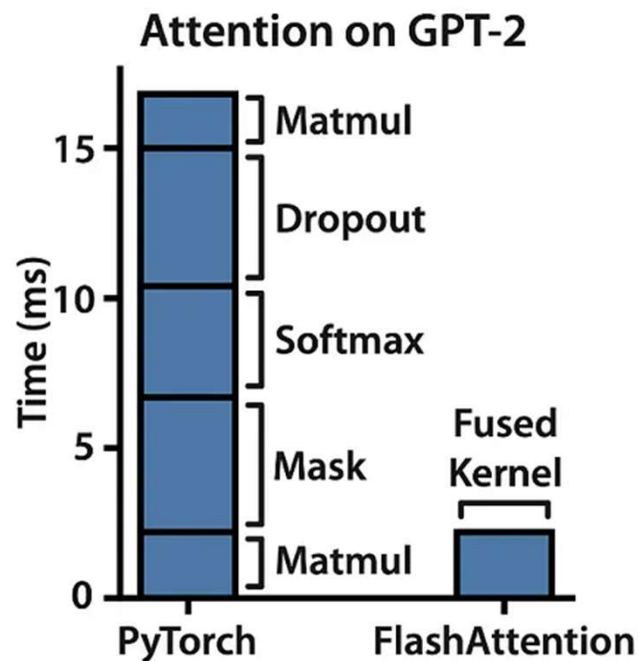
Sourish Dasgupta

LLMs: Advanced Attention  
Mechanisms

Year: 2022; NeurIPS



## How well did they do?



BERT Implementation	Training time (minutes)
Nvidia MLPerf 1.1 [58]	20.0 ± 1.5
FLASHATTENTION (ours)	<b>17.4 ± 1.4</b>





# Key Takeaways


- Avoid unnecessary HBM writes
- Maximize SRAM computation

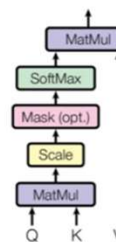




# Want more? Follow:

 Search 

[Browse State-of-the-Art](#) [Datasets](#) [Methods](#) [More](#) 



<https://paperswithcode.com/methods/category/attention-mechanisms>



Sourish Dasgupta

LLMs: Advanced Attention Mechanisms