



# Efficient Streaming Language Models with Attention Sinks

**Guangxuan Xiao<sup>1</sup>, Yuandong Tian<sup>2</sup>, Beidi Chen<sup>3</sup>, Song Han<sup>1</sup>, Mike Lewis<sup>2</sup>**

Massachusetts Institute of Technology<sup>1</sup>

Meta AI<sup>2</sup>

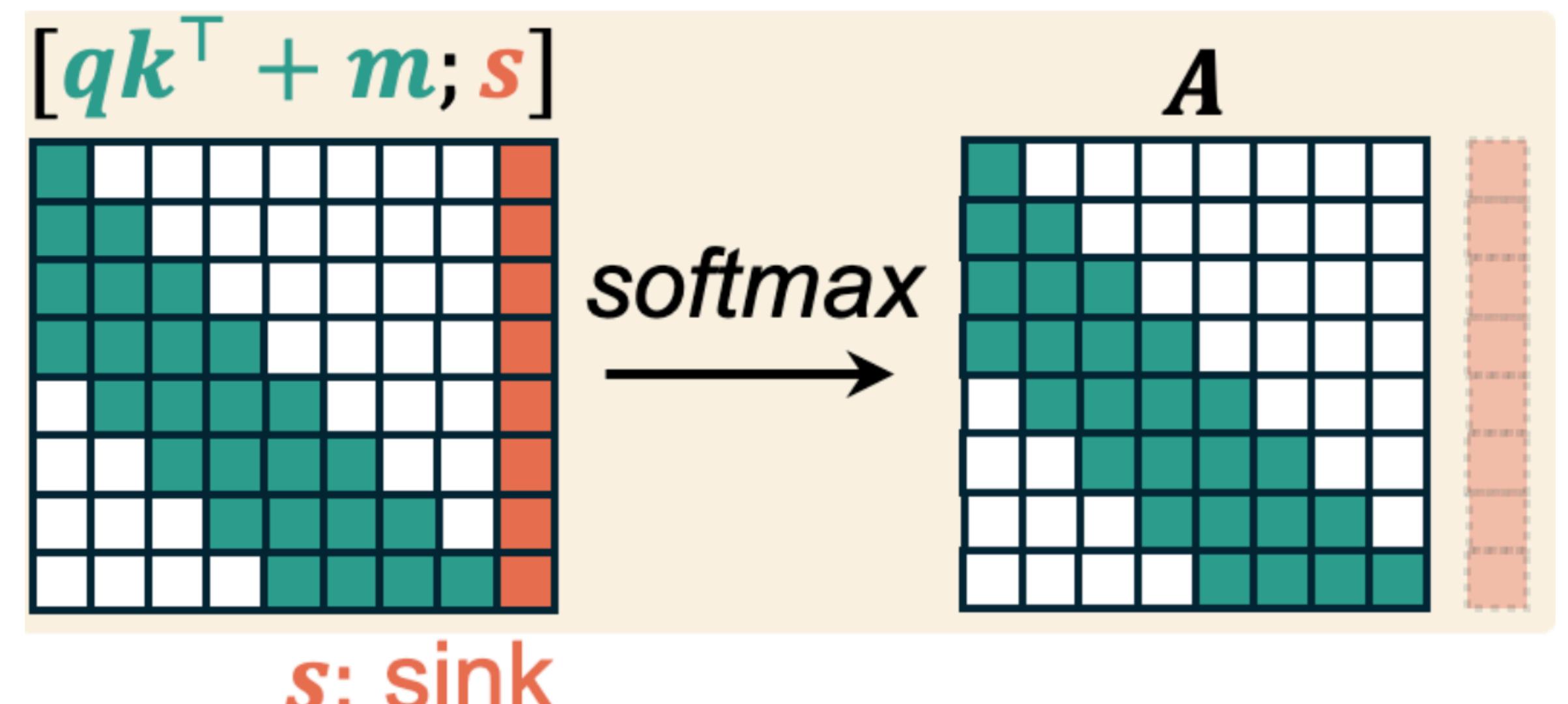
Carnegie Mellon University<sup>3</sup>

# Attention Sinks in GPT-OSS

- GPT-OSS added one learnable attention sink per attention head—enables the model to "**pay no attention to any tokens**" when needed
- Dedicated attention sinks can stabilize model convergence and make the model easier to quantize.

```
def sdpa(Q, K, V, S, sm_scale, sliding_window=0):
    # sliding_window == 0 means no sliding window
    n_tokens, n_heads, q_mult, d_head = Q.shape
    assert K.shape == (n_tokens, n_heads, d_head)
    assert V.shape == (n_tokens, n_heads, d_head)
    K = K[:, :, None, :].expand(-1, -1, q_mult, -1)
    V = V[:, :, None, :].expand(-1, -1, q_mult, -1)
    S = S.reshape(n_heads, q_mult, 1, 1).expand(-1, -1, n_tokens, -1)
    mask = torch.triu(Q.new_full((n_tokens, n_tokens), -float("inf")), diagonal=1)
    if sliding_window > 0:
        mask += torch.tril(
            mask.new_full((n_tokens, n_tokens), -float("inf")), diagonal=-sliding_window
        )
    QK = torch.einsum("qhmd,khmd->hmqk", Q, K)
    QK *= sm_scale
    QK += mask[None, None, :, :]
    QK = torch.cat([QK, S], dim=-1)
    W = torch.softmax(QK, dim=-1)
    W = W[..., :-1]
    attn = torch.einsum("hmqk,khmd->qhmd", W, V)
    return attn.reshape(n_tokens, -1)
```

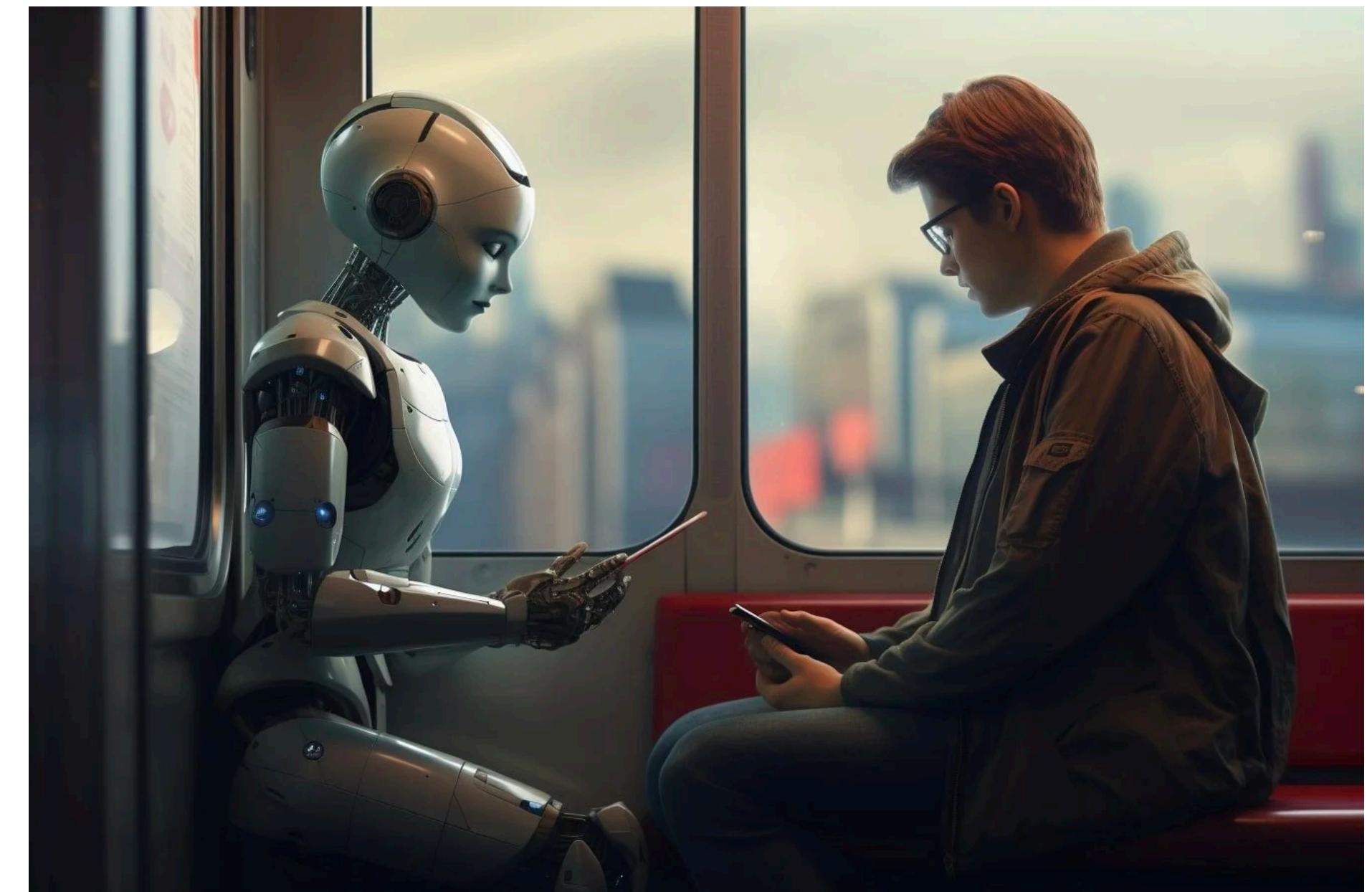
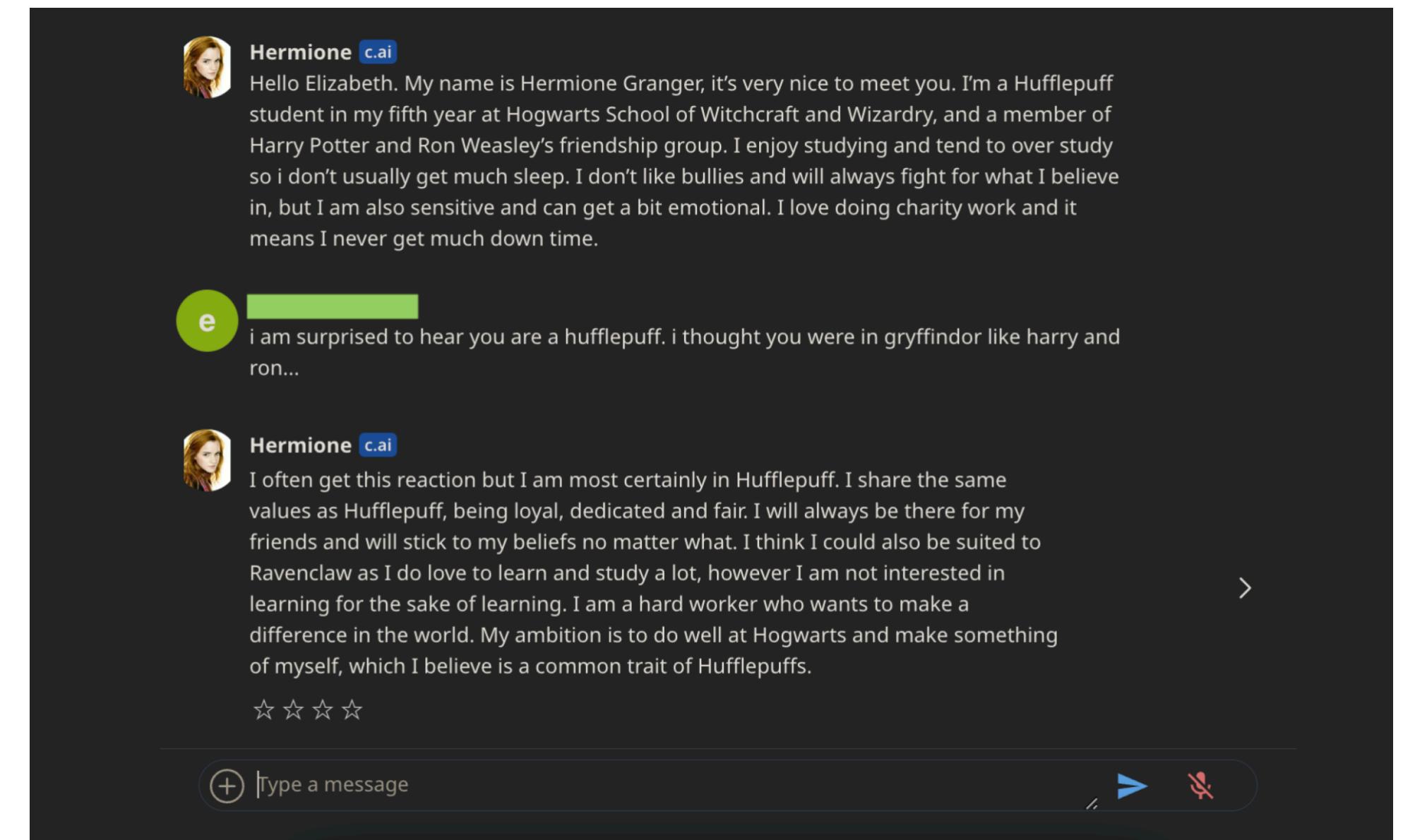
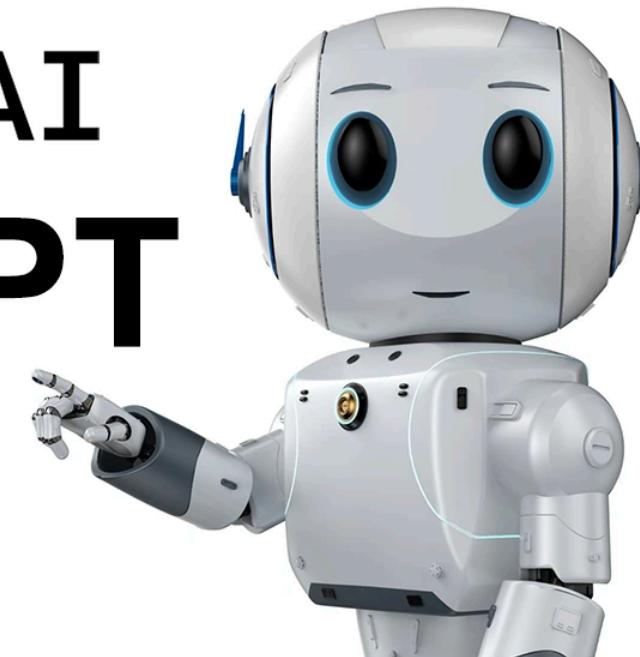
**Attention:** Following GPT-3, attention blocks alternate between banded window and fully dense patterns [10][11], where the bandwidth is 128 tokens. Each layer has 64 query heads of dimension 64, and uses Grouped Query Attention (GQA [12][13]) with 8 key-value heads. We apply rotary position embeddings [14] and extend the context length of dense layers to 131,072 tokens using YaRN [15]. Each attention head has a learned bias in the denominator of the softmax, similar to off-by-one attention and **attention sinks** [16][17], which enables the attention mechanism to pay no attention to any tokens.



$$\text{attention\_probs} = \text{softmax}([\text{sink}, a_1, a_2, \dots, a_t])$$

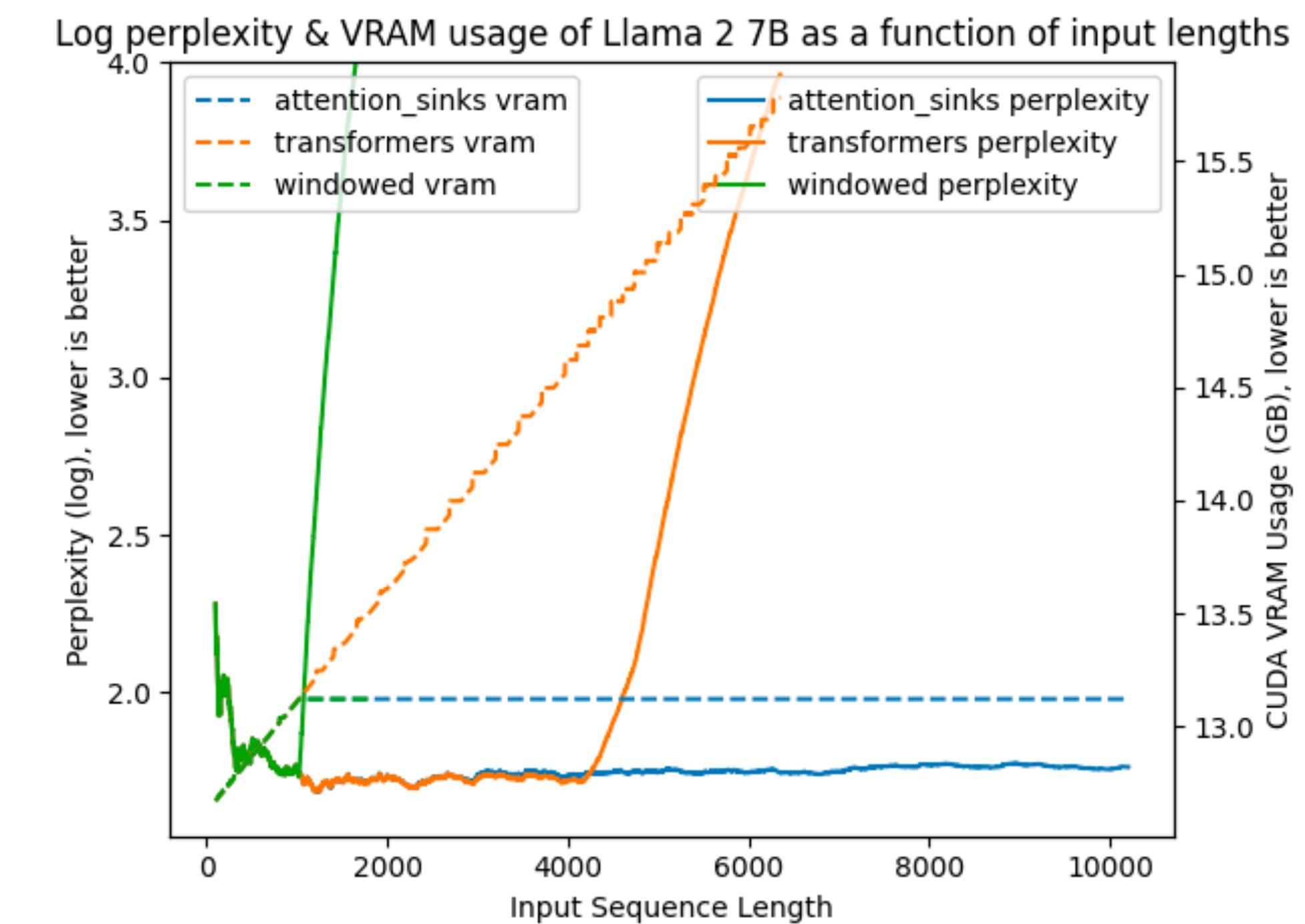
# Motivation: Use cases

 OpenAI  
ChatGPT



# Challenges of Deploying LLMs in Streaming Applications

- Urgent need for LLMs in streaming applications such as multi-round dialogues, where long interactions are needed.
- Challenges:
  - Extensive memory consumption during the decoding stage.
  - Inability of popular LLMs to generalize to longer text sequences.



[https://github.com/tomaarsen/attention\\_sinks](https://github.com/tomaarsen/attention_sinks)

# Challenges of Deploying LLMs in Streaming Applications

w/o StreamingLLM

```
(streaming) guangxuan@l29:~/workspace/streaming-llm$ CUDA_VISIBLE_DEVICES=0 python examples/run_streaming_llama.py  
Loading model from lmsys/vicuna-13b-v1.3 ...  
Loading checkpoint shards: 67%|██████████| 2/3 [00:09<00:04, 4.94s/it]
```

w/ StreamingLLM

```
(streaming) guangxuan@l29:~/workspace/streaming-llm$ CUDA_VISIBLE_DEVICES=1 python examples/run_streaming_llama.py --enable_streaming  
Loading model from lmsys/vicuna-13b-v1.3 ...  
Loading checkpoint shards: 67%|██████████| 2/3 [00:09<00:04, 4.89s/it]
```

# Challenges of Deploying LLMs in Streaming Applications

# w/o StreamingLLM

# Model Performance Breaks

ASSISTANT: 0000000-t-t-t-t'

USER: Write a C++ program to print Fibonacci number using recursion.

USER: Now we define a sequence of numbers in which each number is the sum of the three preceding ones. The first three numbers are 0, -1, -1. Write a program to find the nth number.

ASSTSTANT: 0-a-a-a-eah000000000000

**USER:** Write a simple website in HTML. When a user clicks the button, it shows a random joke from a list of 4 jokes.

ASSISTANT: ■

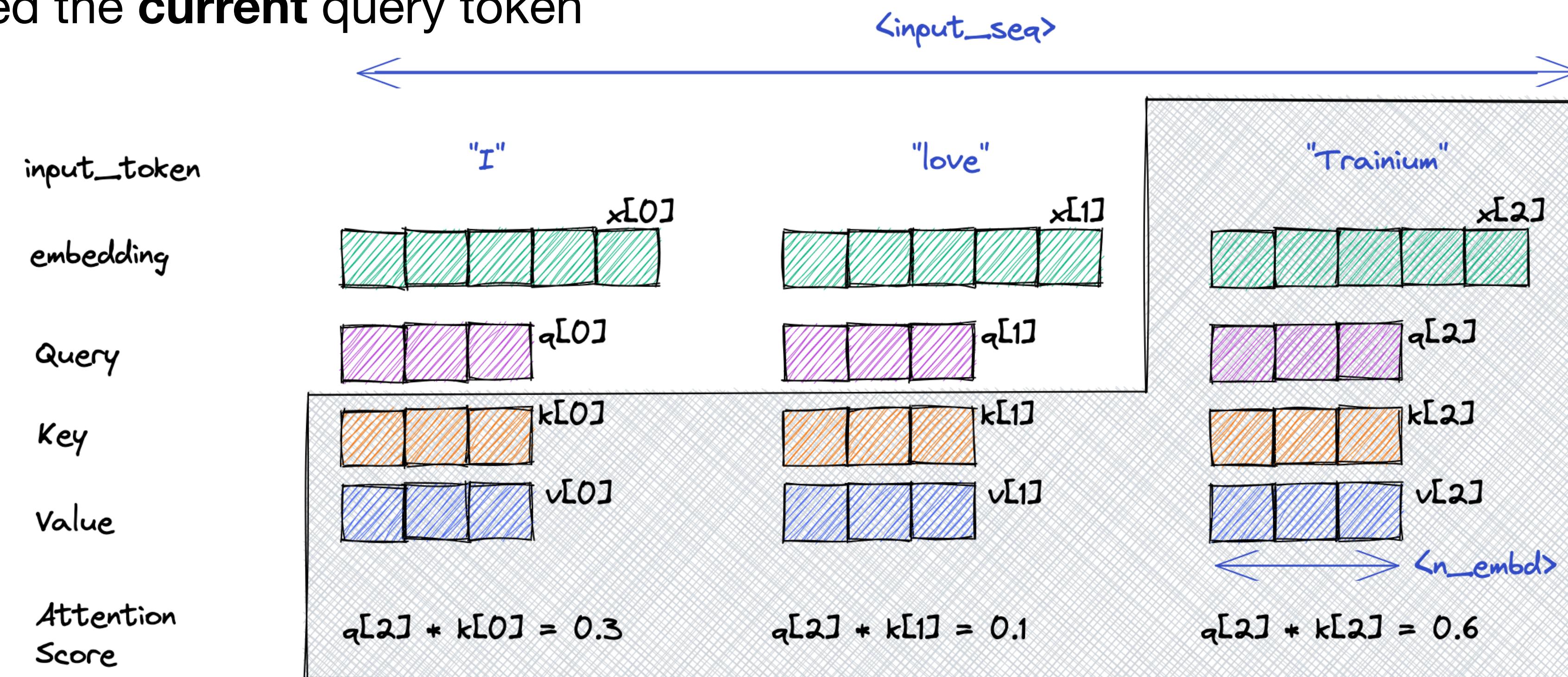
# w/o StreamingLLM

```
outputs = model(
    File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/torch/nn/modules/module.py", line 1501, in __call__impl
        return forward_call(*args, **kwargs)
    File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/transformers/models/llama/modeling_llama.py", line 820, in forward
    outputs = self.model(
        File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/torch/nn/modules/module.py", line 1501, in __call__impl
            return forward_call(*args, **kwargs)
        File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/transformers/models/llama/modeling_llama.py", line 708, in forward
            layer_outputs = decoder_layer(
                File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/torch/nn/modules/module.py", line 1501, in __call__impl
                    return forward_call(*args, **kwargs)
                File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/transformers/models/llama/modeling_llama.py", line 424, in forward
                    hidden_states, self_attn_weights, present_key_value = self.self_attn(
                        File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/torch/nn/modules/module.py", line 1501, in __call__impl
                            return forward_call(*args, **kwargs)
                        File "/home/guangxuan/miniconda3/envs/streaming/lib/python3.8/site-pac
kages/transformers/models/llama/modeling_llama.py", line 337, in forward
                            key_states = torch.cat([past_key_value[0], key_states], dim=2)
                        torch.cuda.OutOfMemoryError: CUDA out of memory. Tried to allocate 90.00
                            MiB (GPU 0; 47.54 GiB total capacity; 44.53 GiB already allocated; 81.0
                            6 MiB free; 46.47 GiB reserved in total by PyTorch) If reserved memory i
                            s >> allocated memory try setting max_split_size_mb to avoid fragmentati
                            on. See documentation for Memory Management and PYTORCH_CUDA_ALLOC_CONF
                            (streaming) guangxuan@l29:~/workspace/streaming-llm$
```

# The Problem of Long Context: Large KV Cache

The KV cache could be large with long context

- During Transformer decoding (GPT-style), we need to store the **Keys** and **Values** of **all previous** tokens so that we can perform the attention computation, namely the **KV cache**
  - Only need the **current** query token



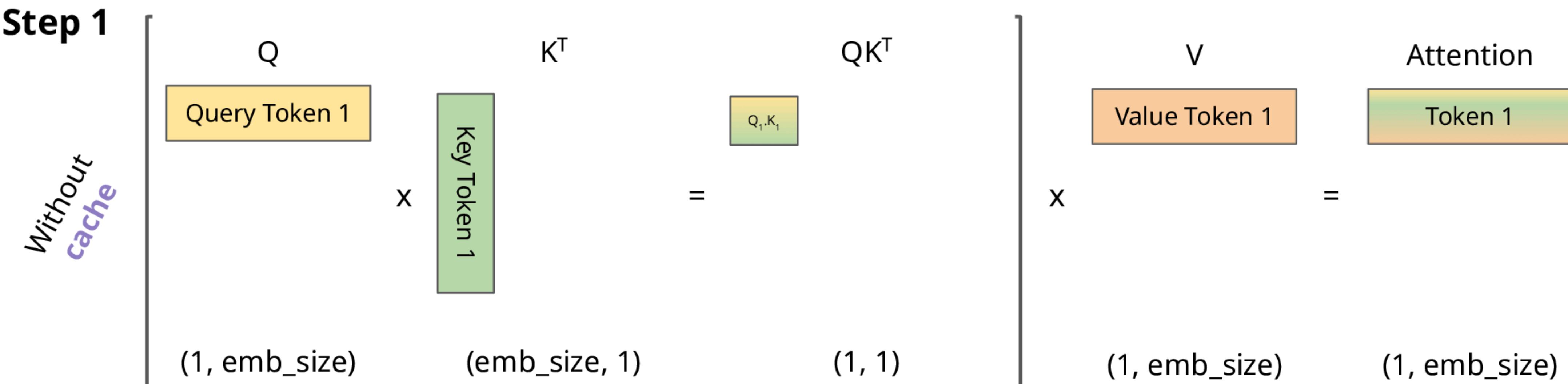
$$a_{ij} = \frac{\exp(q_i^\top k_j / \sqrt{d})}{\sum_{t=1}^i \exp(q_i^\top k_t / \sqrt{d})}, \quad o_i = \sum_{j=1}^i a_{ij} v_j$$

Image credit: <https://awsdocs-neuron.readthedocs-hosted.com/en/latest/general/appnotes/transformers-neuronx/generative-llm-inference-with-neuron.html>

# The Problem of Long Context: Large KV Cache

The KV cache could be large with long context

- During Transformer decoding (GPT-style), we need to store the **Keys** and **Values** of **all previous** tokens so that we can perform the attention computation, namely the **KV cache**
  - Only need the **current** query token



# The Problem of Long Context: Large KV Cache

**The KV cache could be large with long context**

- During Transformer decoding (GPT-style), we need to store the **Keys** and **Values** of **all previous** tokens so that we can perform the attention computation, namely the **KV cache**
  - Only need the **current** query token

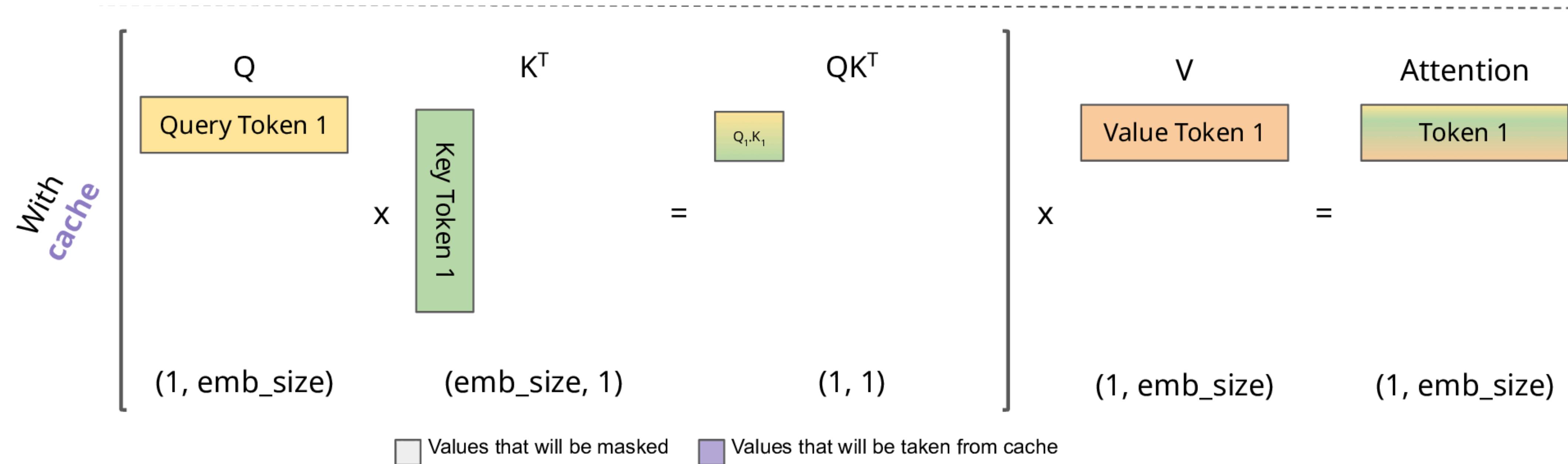


Image credit: <https://medium.com/@joaolages/kv-caching-explained-276520203249>

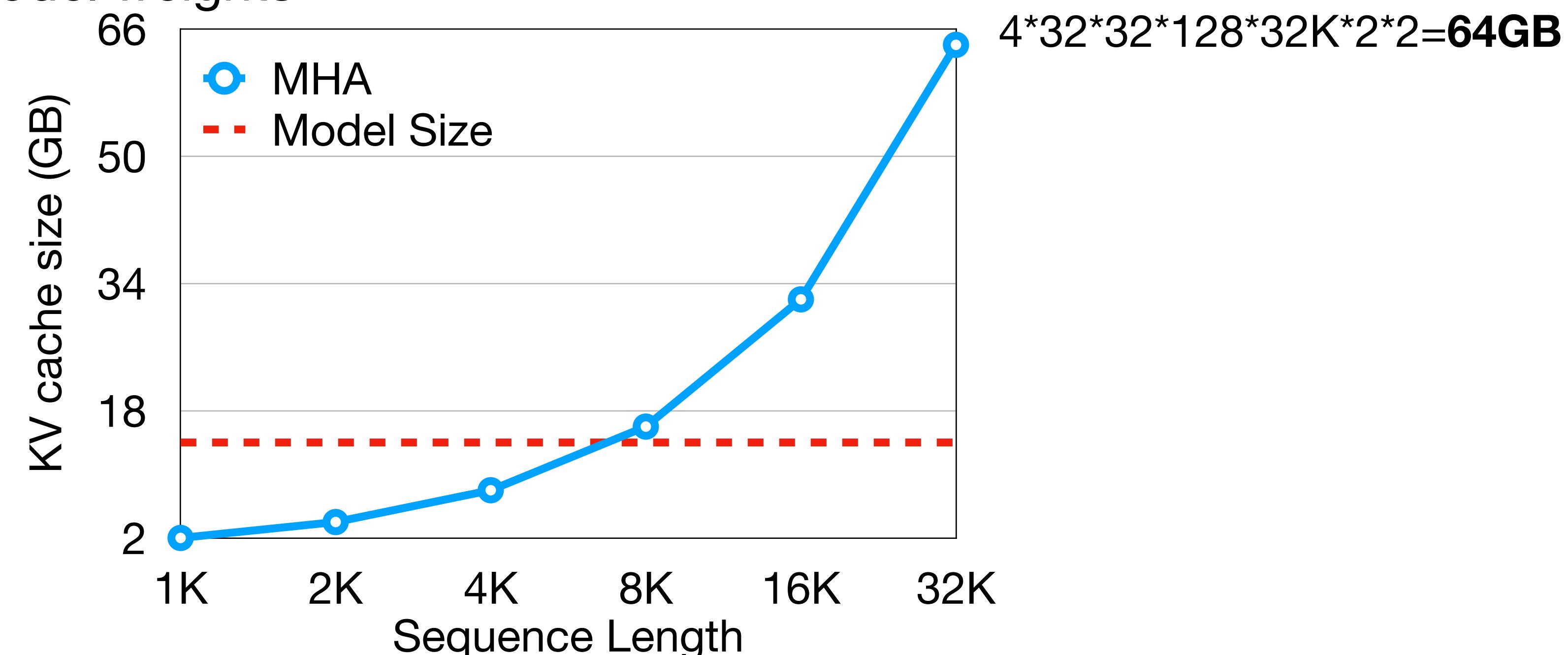
# The Problem of Long Context: Large KV Cache

## The KV cache could be large with long context

- We can calculate the memory required to store the KV cache
- Take Llama-2-7B as an example

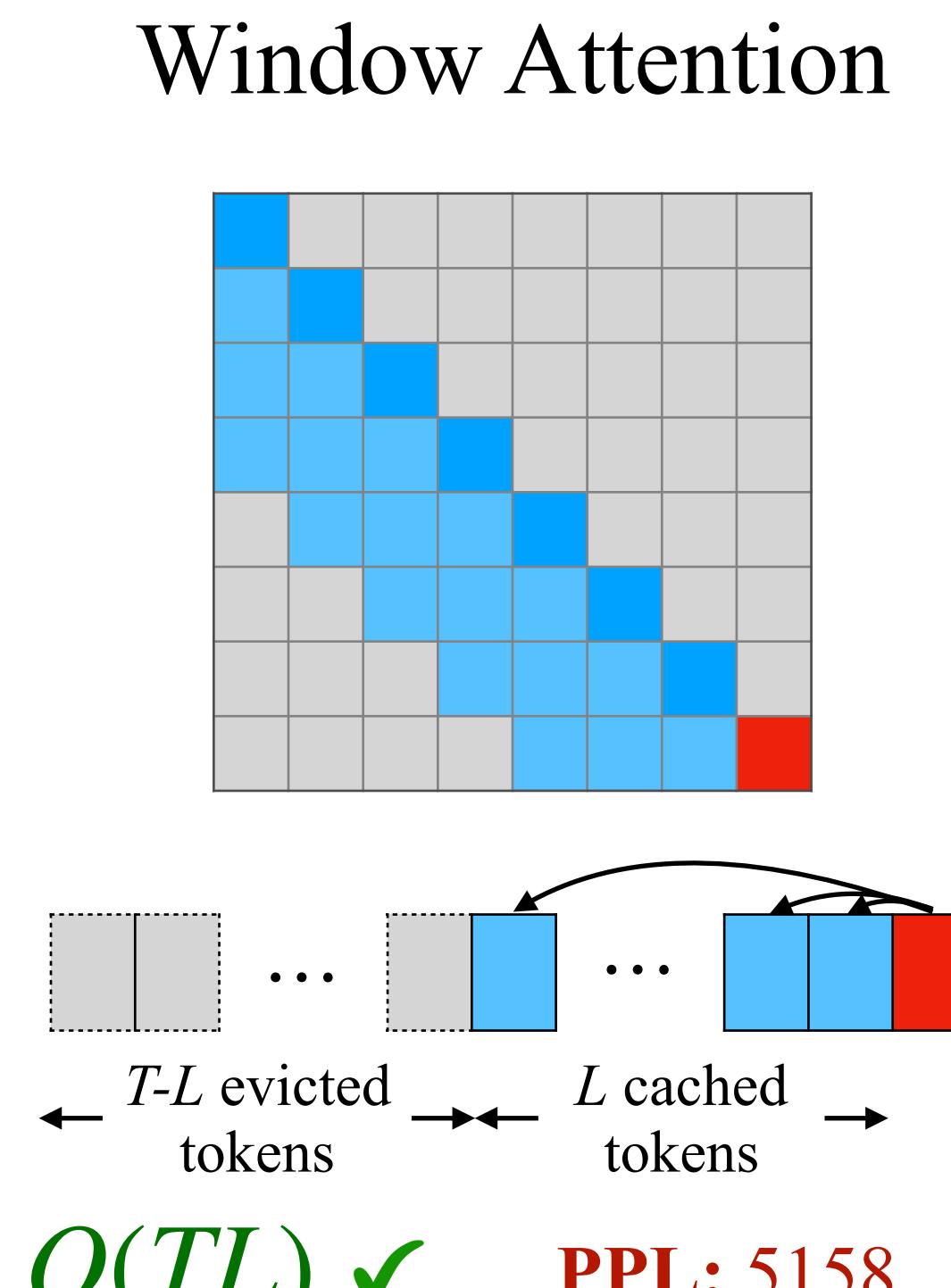
$$\underbrace{BS}_{batchsize} * \underbrace{32}_{layers} * \underbrace{32}_{kv-heads} * \underbrace{128}_{n_{emd}} * \underbrace{N}_{length} * \underbrace{2}_{K\&V} * \underbrace{2\text{bytes}}_{FP16} = 0.5\text{MB} \times BS \times N$$

- Now we calculate the KV cache size under  $BS = 4$  and different sequence lengths.
  - Quickly larger than model weights

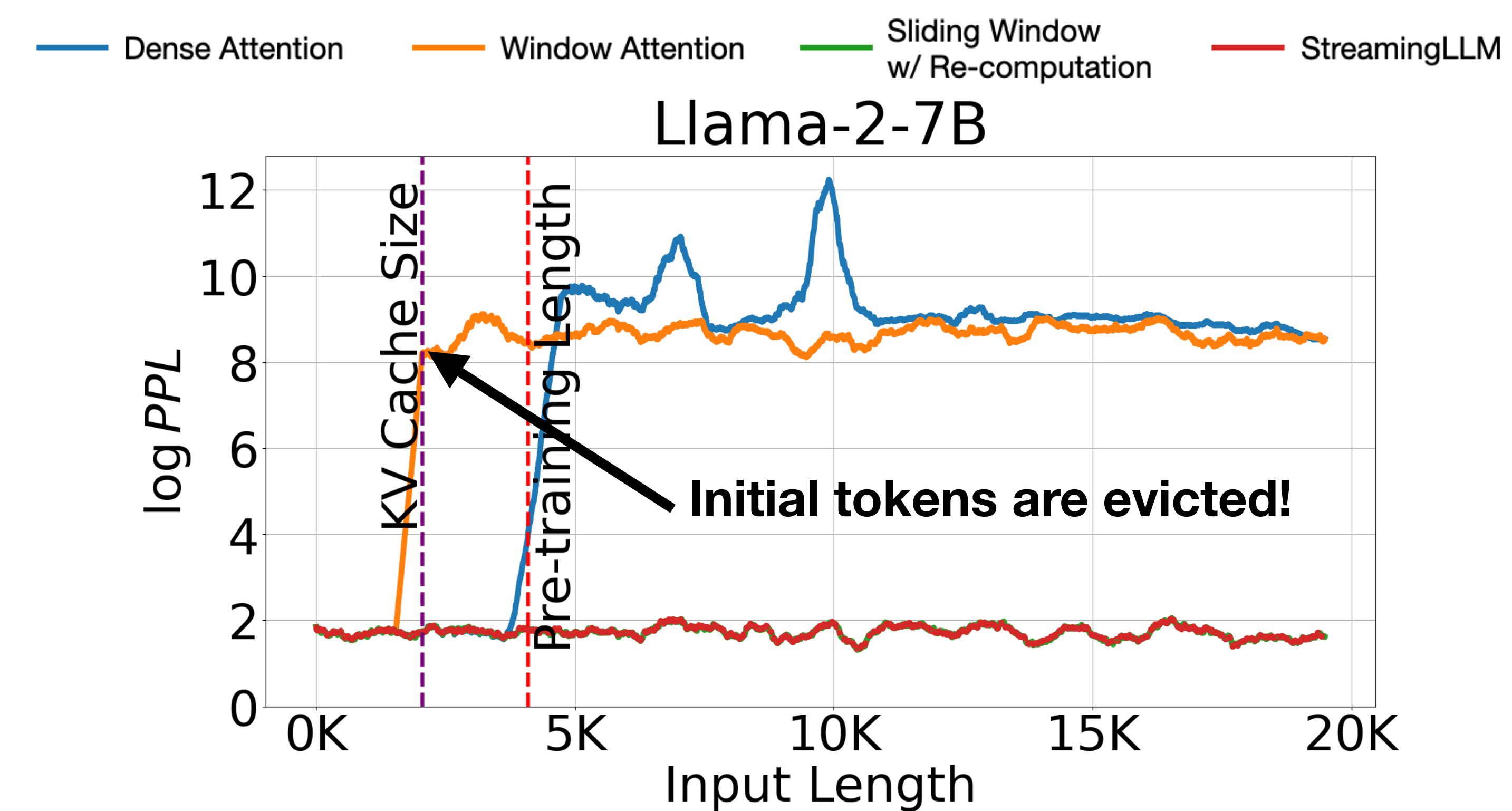


# The Limits of Window Attention

- A natural approach – window attention: caching only the most recent Key-Value states.
- Drawback: model collapses when the text length surpasses the cache size, when the initial token is evicted.

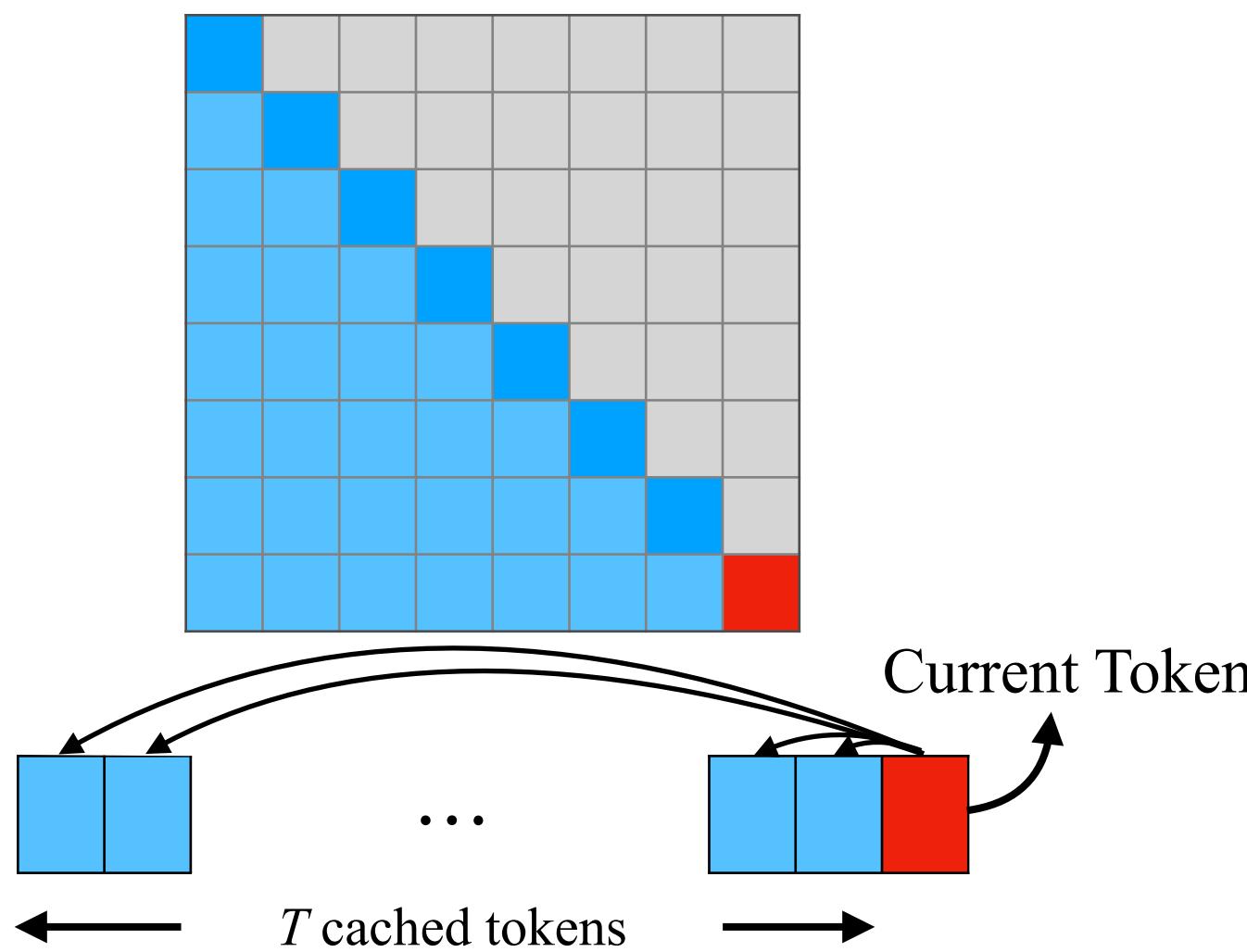


Breaks when initial tokens  
are evicted.



# Difficulties of Other Methods

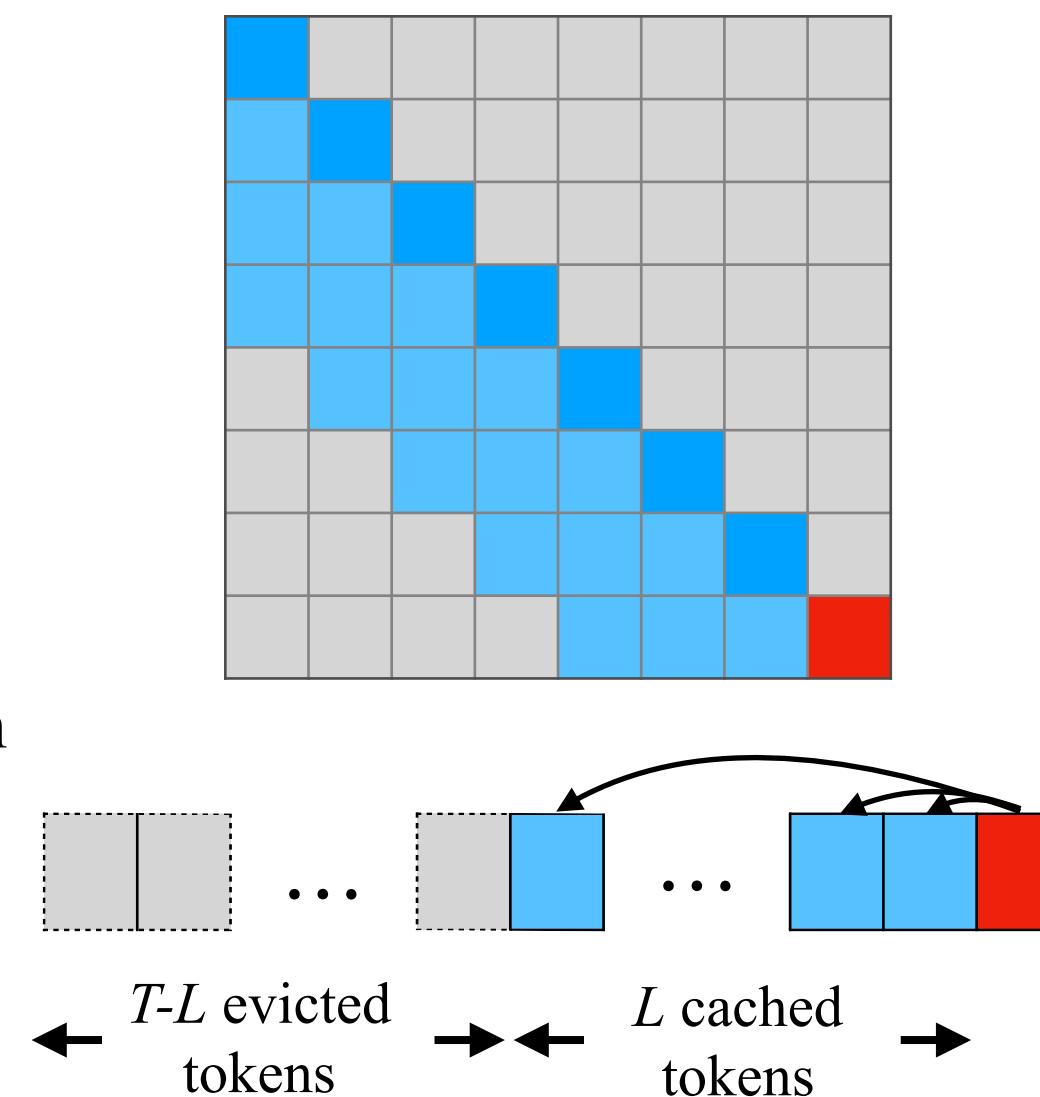
Dense Attention



$O(T^2) \times$  PPL: 5641  $\times$

Has poor efficiency and performance on long text.

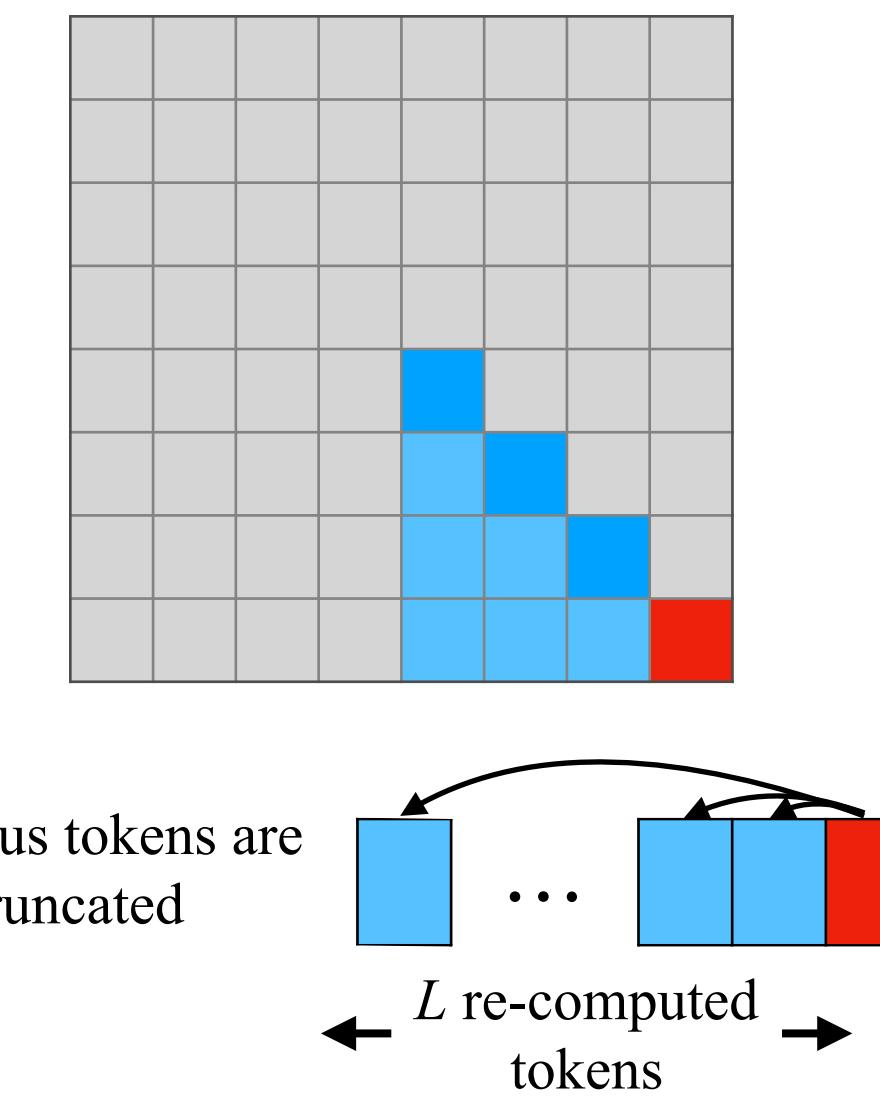
Window Attention



$O(TL) \checkmark$  PPL: 5158  $\times$

Breaks when initial tokens are evicted.

Sliding Window w/ Re-computation



$O(TL^2) \times$  PPL: 5.43  $\checkmark$

Has to re-compute cache for each incoming token.

# The “Attention Sink” Phenomenon

- **Observation:** initial tokens have large attention scores, even if they're not semantically significant.
- **Attention Sink:** Tokens that disproportionately attract attention irrespective of their relevance.

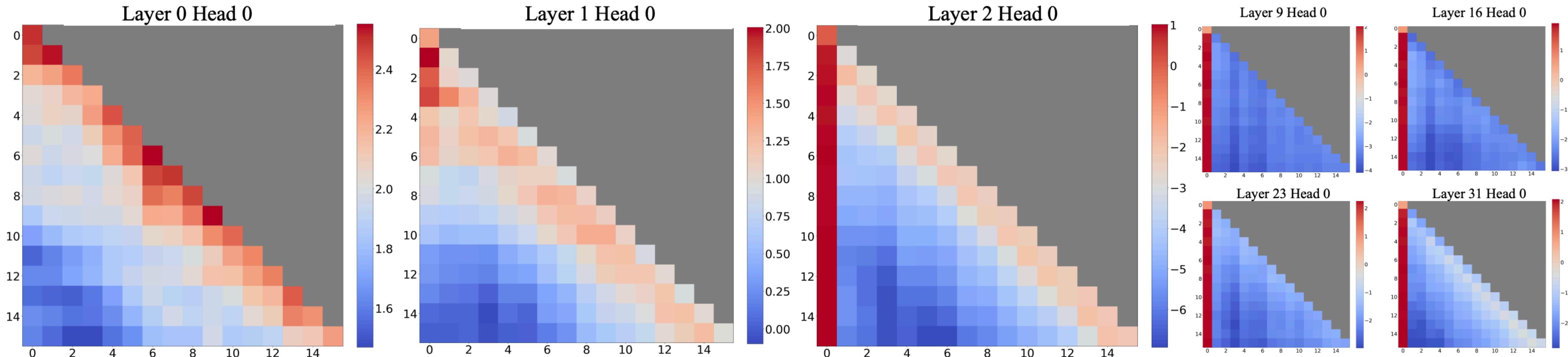


Figure 2: Visualization of the *average* attention logits in Llama-2-7B over 256 sentences, each with a length of 16. Observations include: (1) The attention maps in the first two layers (layers 0 and 1) exhibit the "local" pattern, with recent tokens receiving more attention. (2) Beyond the bottom two layers, the model heavily attends to the initial token across all layers and heads.

$$\text{SoftMax}(x)_i = \frac{e^{x_i}}{e^{x_1} + \sum_{j=2}^N e^{x_j}}, \quad x_1 \gg x_j, j \in 2, \dots, N$$

What window attention evicts

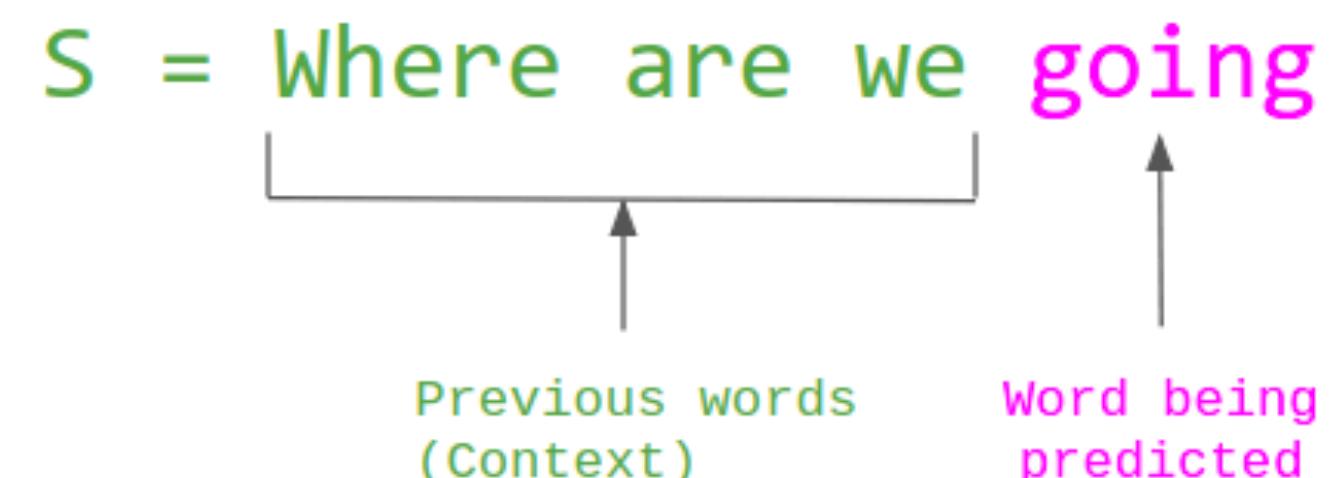
# Understanding Why Attention Sinks Exist

## The Rationale Behind Attention Sinks

- SoftMax operation's role in creating attention sinks – attention scores have to sum up to one for all contextual tokens.

$$\text{SoftMax}(x)_i = \frac{e^{x_i}}{e^{x_1} + \sum_{j=2}^N e^{x_j}}, \quad x_1 \gg x_j, j \in 2, \dots, N$$

- Initial tokens' advantage in becoming sinks due to their visibility to subsequent tokens, rooted in autoregressive language modeling.



- Does the importance of the initial tokens arise from their **position** or their **semantics**?
  - We found adding initial four "\n"s can also recover perplexity.
  - Therefore, it is **position!**

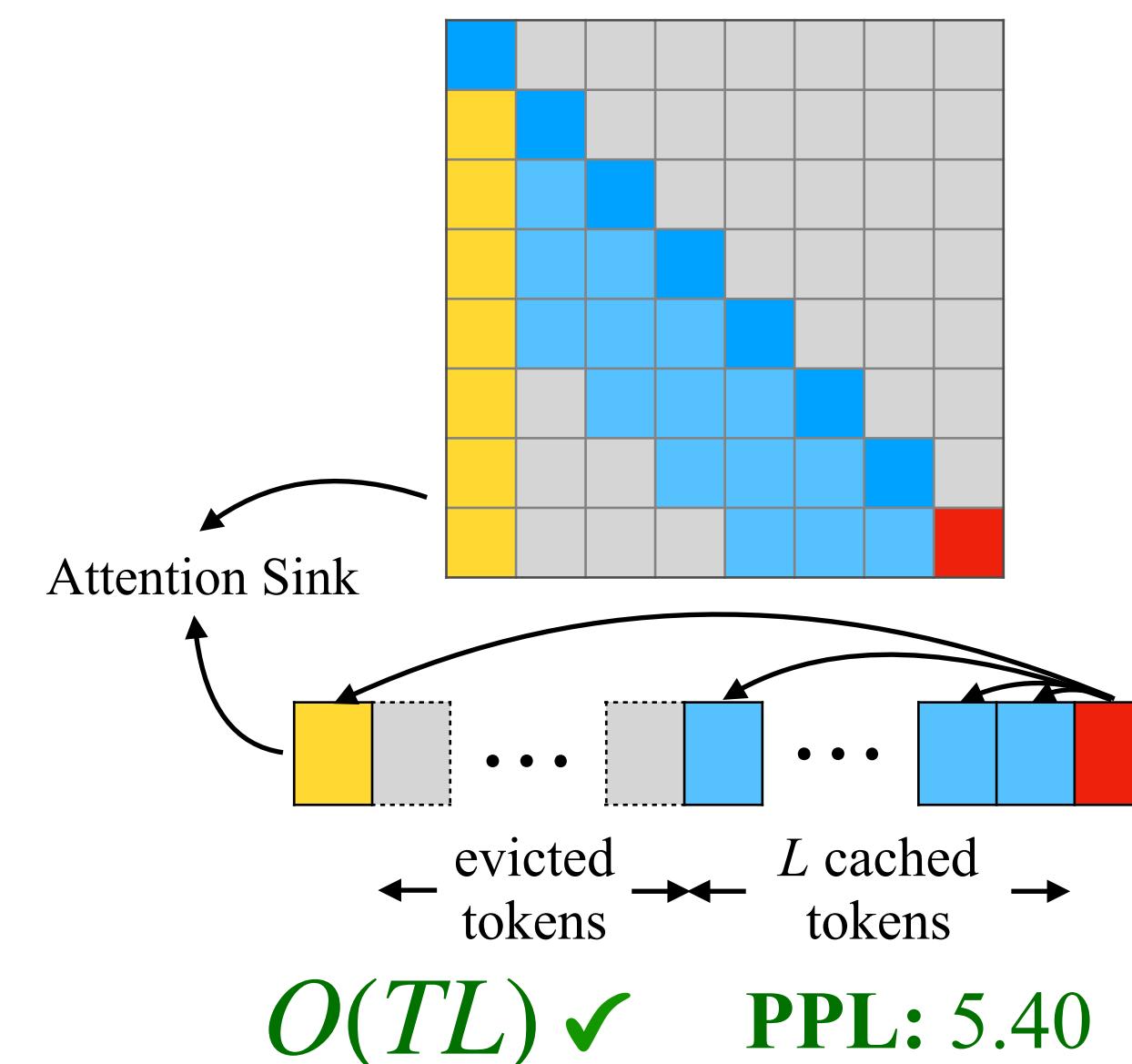
Llama-2-13B	PPL ( $\downarrow$ )
0 + 1024 (Window)	5158.07
4 + 1020	5.40
4"\n"+1020	5.60

$$P(S) = P(\text{Where}) \times P(\text{are} \mid \text{Where}) \times P(\text{we} \mid \text{Where are}) \times P(\text{going} \mid \text{Where are we})$$

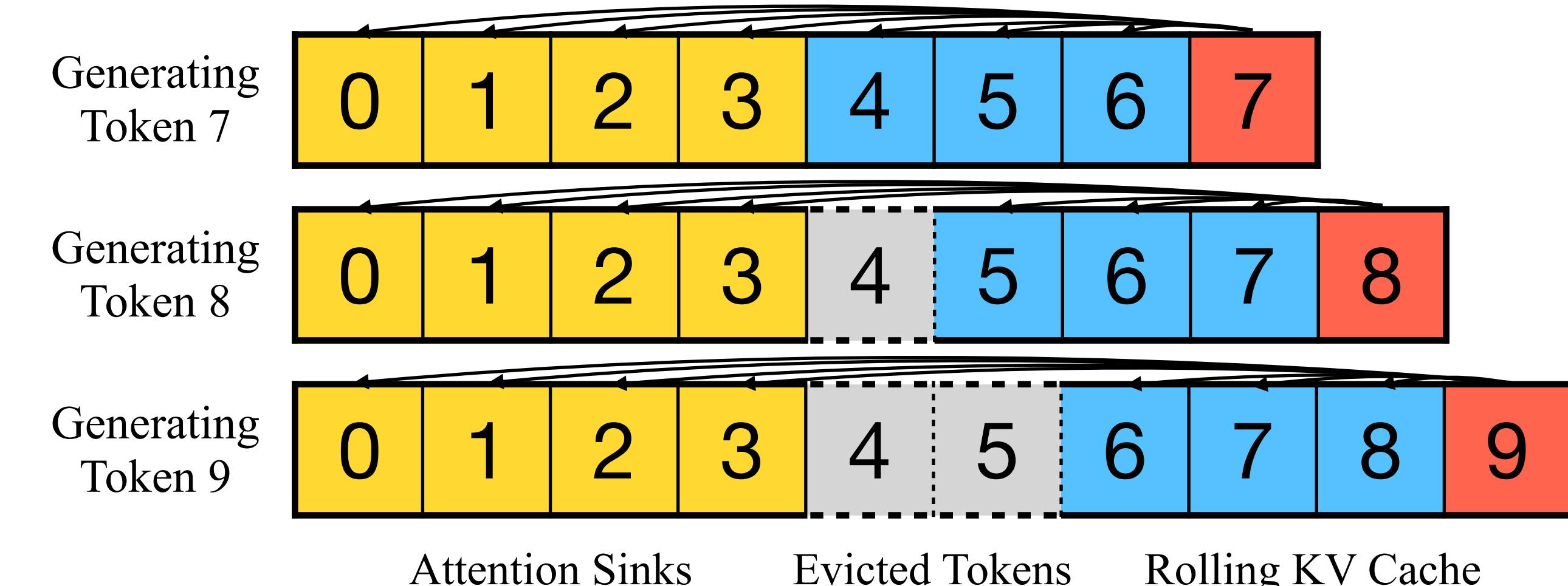
# StreamingLLM: Using Attention Sinks for Infinite Streams

- **Objective:** Enable LLMs trained with a finite attention window to handle infinite text lengths without additional training.
- **Key Idea:** **preserve the KV of attention sink tokens**, along with the sliding window's KV to stabilize the model's behavior.

StreamingLLM (ours)

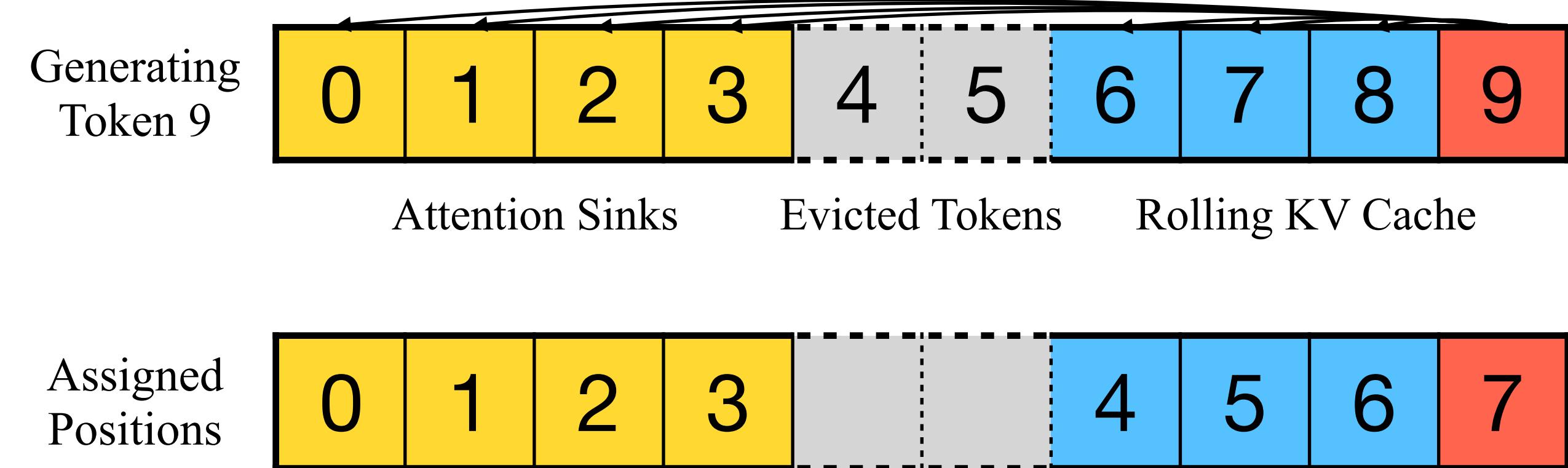


Can perform efficient and stable language modeling on long texts.



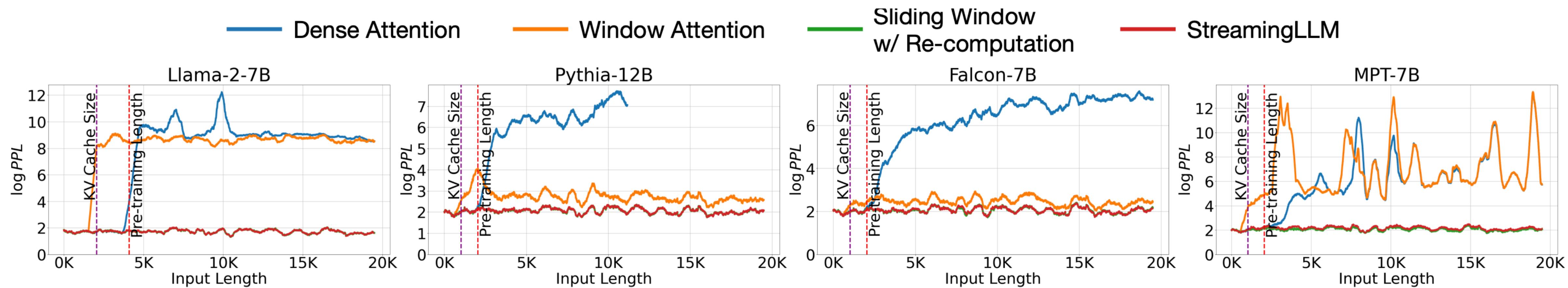
# Positional Encoding Assignment

- Use positions *in the cache* instead of those *in the original text*.



# Streaming Performance

- Comparison between dense attention, window attention, and sliding window w/ re-computation.

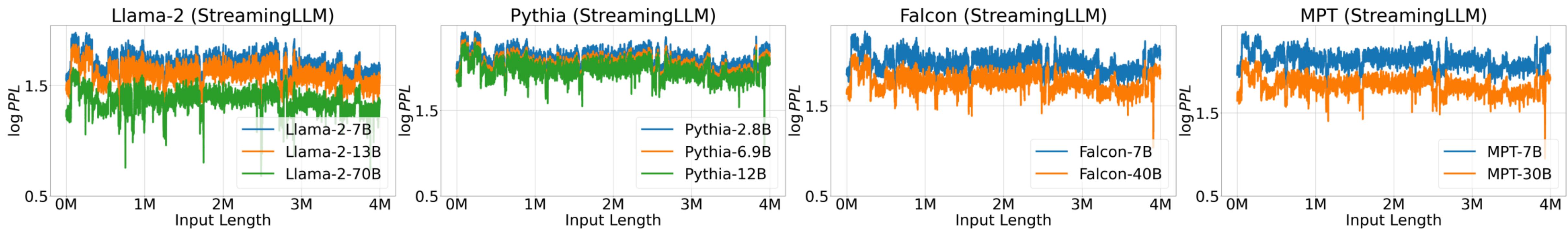


- Dense attention fails beyond pre-training attention window size.
- Window attention fails after input exceeds cache size (initial tokens evicted).
- StreamingLLM shows stable performance; perplexity close to sliding window with re-computation baseline.

# Streaming Performance

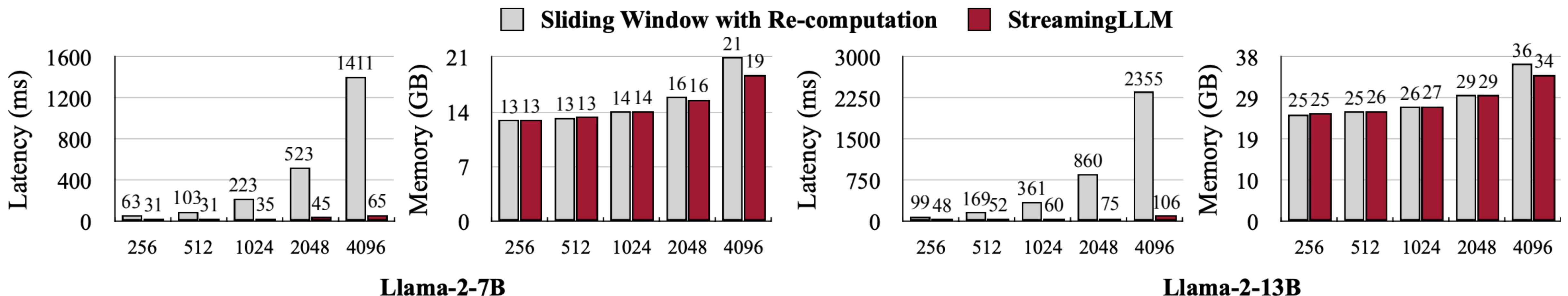
## Super Long Language Modeling

- With StreamingLLM, model families include Llama-2, MPT, Falcon, and Pythia can now effectively model up to 4 million tokens.



# Efficiency

- **Comparison baseline:** The sliding window with re-computation, a method that is computationally heavy due to quadratic attention computation within its window.
- StreamingLLM provides up to 22.2x speedup over the baseline, making LLMs for real-time streaming applications feasible.



# Ablation Study: #Attention Sinks

- The number of attention sinks that need to be introduced to recover perplexity.
  - 4 attention sinks are generally enough.

Cache Config	0+2048	1+2047	2+2046	4+2044	8+2040
Falcon-7B	17.90	12.12	12.12	12.12	12.12
MPT-7B	460.29	14.99	15.00	14.99	14.98
Pythia-12B	21.62	11.95	12.09	12.09	12.02

Cache Config	0+4096	1+4095	2+4094	4+4092	8+4088
Llama-2-7B	3359.95	11.88	10.51	9.59	9.54

# Pre-training with a Dedicated Attention Sink Token

- **Idea: Why 4 attention sinks?** Can we train a LLM that need only one single attention sink? **Yes!**
- **Method:** Introduce an extra learnable token at the start of all training samples to act as a dedicated attention sink.
- **Result:** This pre-trained model retains performance in streaming cases with just this single sink token, contrasting with vanilla models that require multiple initial tokens.

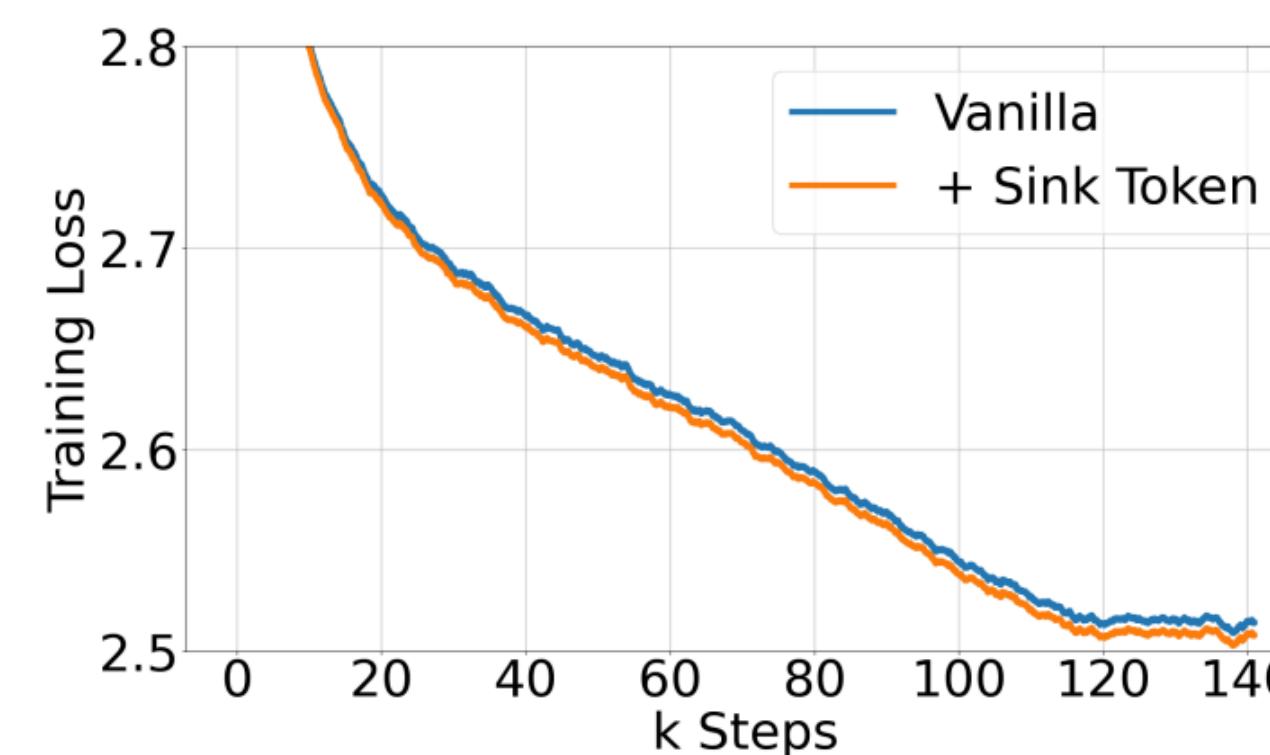


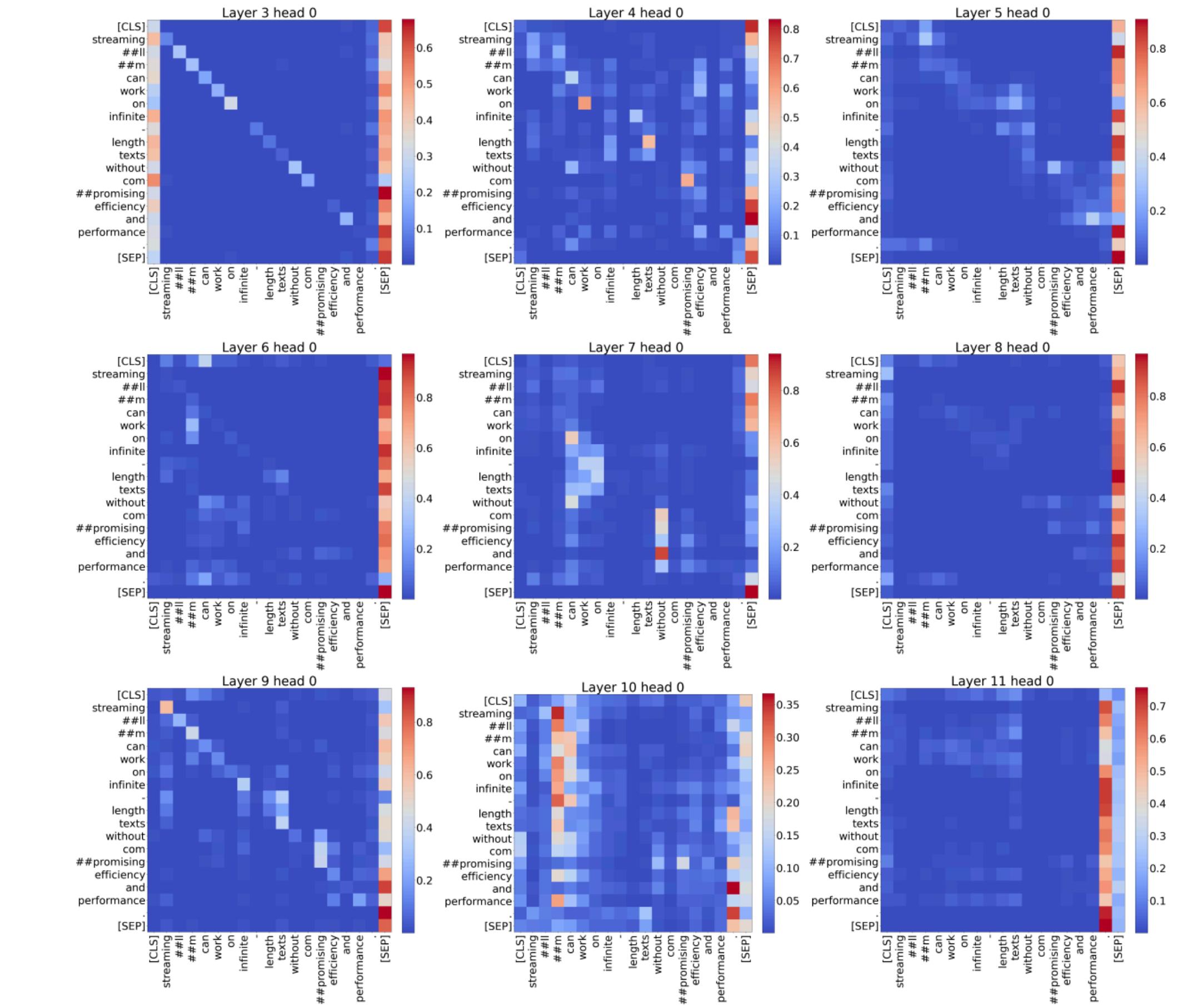
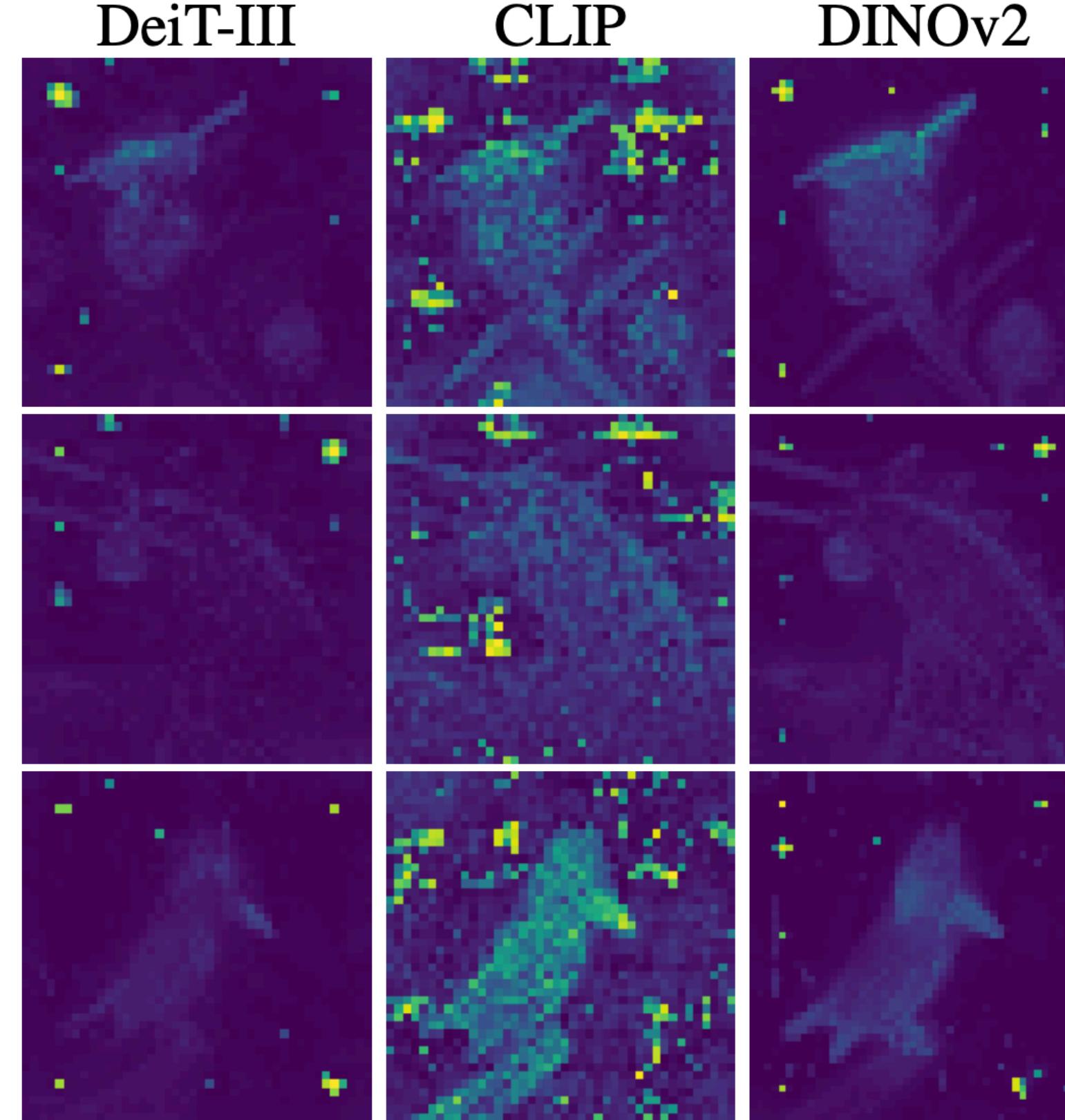
Figure 6: Pre-training loss curves of models w/ and w/o sink tokens. Two models have a similar convergence trend.

Methods	ARC-c	ARC-e	HS	LBD	OBQA	PIQA	WG
Vanilla	18.6	45.2	29.4	39.6	16.0	62.2	50.1
+Sink Token	<b>19.6</b>	<b>45.6</b>	<b>29.8</b>	<b>39.9</b>	<b>16.6</b>	<b>62.6</b>	<b>50.8</b>

Cache Config	0+1024	1+1023	2+1022	4+1020
Vanilla	27.87	<b>18.49</b>	<b>18.05</b>	<b>18.05</b>
Zero Sink	29214	<b>19.90</b>	18.27	18.01
Learnable Sink	1235	<b>18.01</b>	18.01	18.02

# Attention Sinks in Other Transformers

## Encoder Models: ViT and BERT



**ViT: attention sinks emerge in low semantic background pixels.**

**BERT: attention sinks are [SEP] tokens at the end of all sentences**

# Can StreamingLLM give us infinite context?

- Non-stop chatting ≠ Infinite context
- Tokens that are evicted from cache cannot be attended.

**Input Content**

Below is a record of lines I want you to remember.  
The REGISTER\_CONTENT in line 0 is <8806>  
[omitting 9 lines...]  
The REGISTER\_CONTENT in line 10 is <24879>  
[omitting 8 lines...]  
The REGISTER\_CONTENT in line 20 is <45603>  
**Query: The REGISTER\_CONTENT in line 0 is**  
The REGISTER\_CONTENT in line 21 is <29189>  
[omitting 8 lines...]  
The REGISTER\_CONTENT in line 30 is <1668>  
**Query: The REGISTER\_CONTENT in line 10 is**  
The REGISTER\_CONTENT in line 31 is <42569>  
[omitting 8 lines...]  
The REGISTER\_CONTENT in line 40 is <34579>  
**Query: The REGISTER\_CONTENT in line 20 is**  
[omitting remaining 5467 lines...]

**Desired Output**  
["<8806>", "<24879>", "<45603>", ...]

Llama-2-7B-32K-Instruct		Cache Config			
Line Distances	Token Distances	4+2044	4+4092	4+8188	4+16380
20	460	85.80	84.60	81.15	77.65
40	920	80.35	83.80	81.25	77.50
60	1380	79.15	82.80	81.50	78.50
80	1840	75.30	77.15	76.40	73.80
100	2300	0.00	61.60	50.10	40.50
150	3450	0.00	68.20	58.30	38.45
200	4600	0.00	0.00	62.75	46.90
400	9200	0.00	0.00	0.00	45.70
600	13800	0.00	0.00	0.00	28.50
800	18400	0.00	0.00	0.00	0.00
1000	23000	0.00	0.00	0.00	0.00

How to solve this issue?

# Thanks for Listening!

- We propose StreamingLLM, enabling the streaming deployment of LLMs.
- Paper: <https://arxiv.org/abs/2309.17453>
- Code: <https://github.com/mit-han-lab/streaming-llm>
- Demo: <https://youtu.be/UgDcZ3rvRPg>
- Blog: <https://hanlab.mit.edu/blog/streamingllm>

w/o StreamingLLM

```
(streaming) guangxuan@l29:~/workspace/streaming-llm$ CUDA_VISIBLE_DEVICES=0 python examples/run_streaming_llama.py
Loading model from lmsys/vicuna-13b-v1.3 ...
Loading checkpoint shards: 67% [██████] | 2/3 [00:09<00:04, 4.94s/it]
```

w/ StreamingLLM

```
(streaming) guangxuan@l29:~/workspace/streaming-llm$ CUDA_VISIBLE_DEVICES=1 python examples/run_streaming_llama.py --enable_streaming
Loading model from lmsys/vicuna-13b-v1.3 ...
Loading checkpoint shards: 67% [██████] | 2/3 [00:09<00:04, 4.89s/it]
```