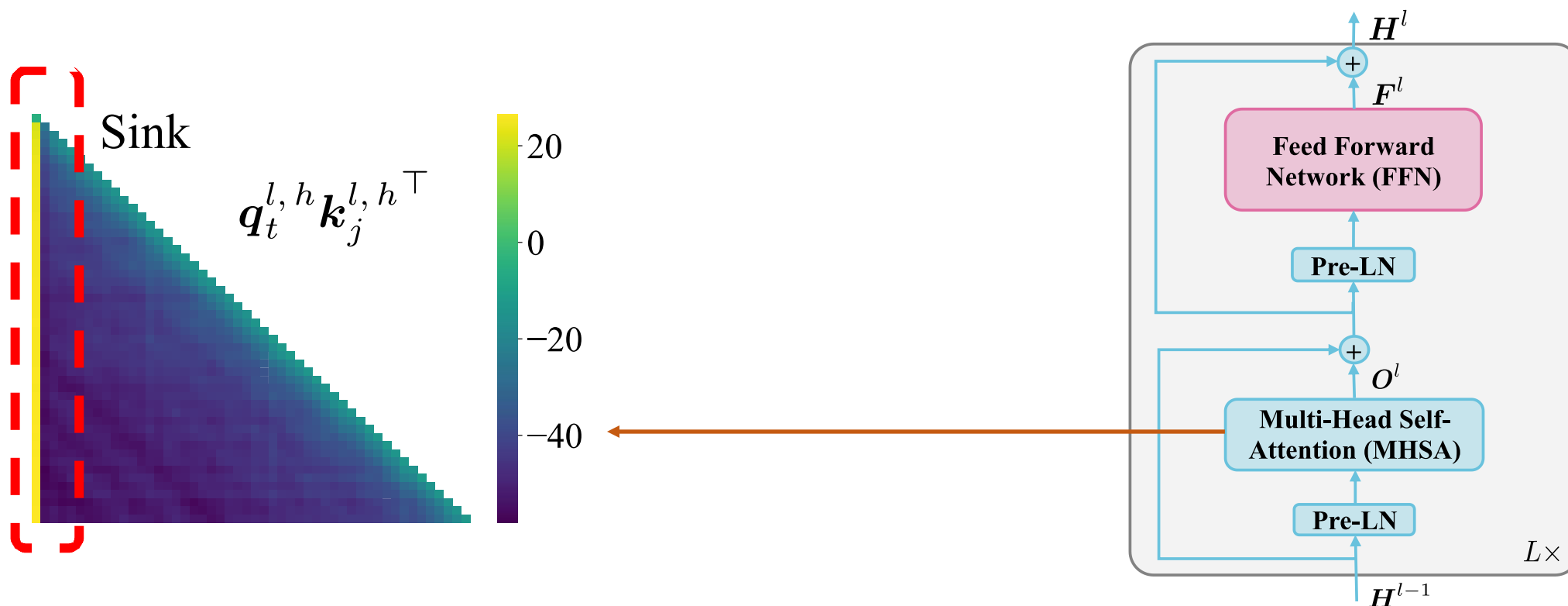


# When Attention Sink Emerges in Language Models: An Empirical View

**Xiangming Gu**, Tianyu Pang, Chao Du, Qian Liu,  
Fengzhuo Zhang, Cunxiao Du, Ye Wang, Min Lin

# What is attention sink?

- Attention sink refers to that Language Models (LMs) assign significant attention to the first token (Xiao et al. 2024)



Xiao et al. Efficient Streaming Language Models with Attention Sinks. ICLR 2024

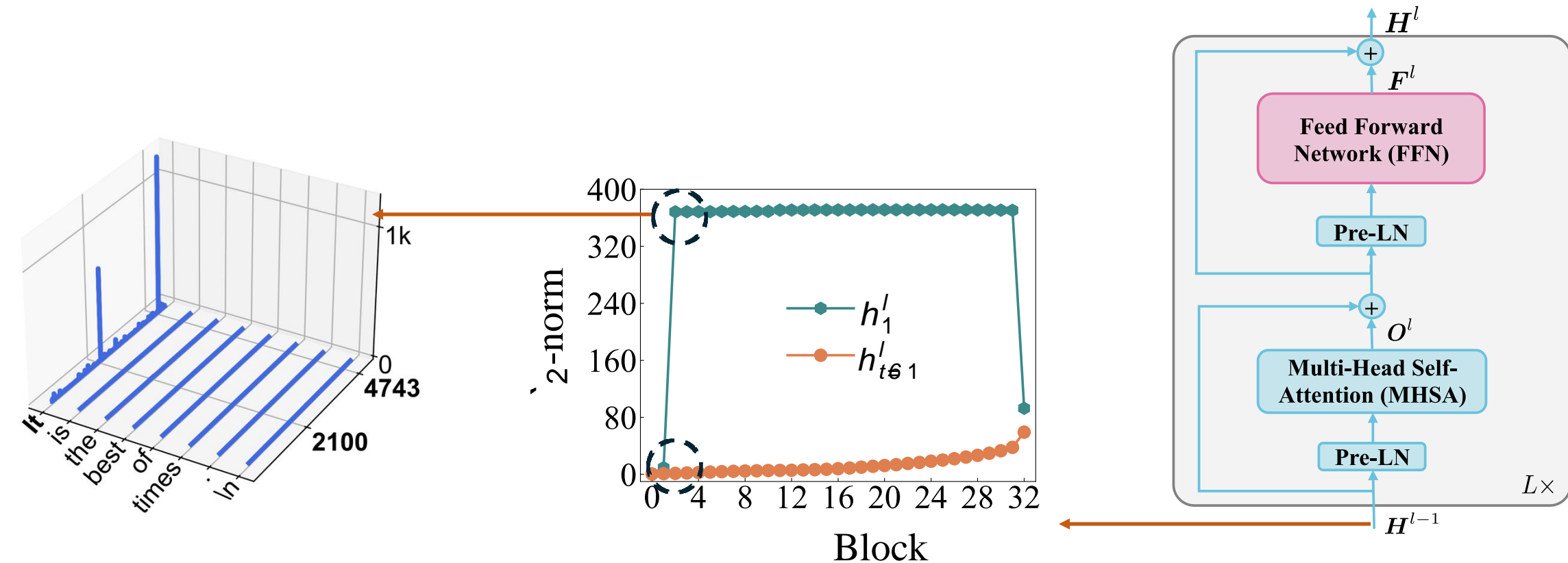
# What can we do with attention sink?

- Long context understanding / generation
- KV cache compression
- Model quantization
- etc

Attention sink is important!

# Mechanism of attention sink

- Massive Activations in hidden states of sink token: its L2-norm is significantly larger than that of other tokens (Cancedda 2024; Sun et al. 2024)



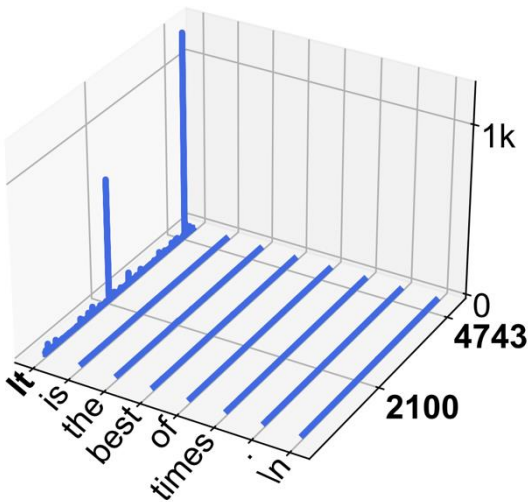
Cancedda, Nicola. Spectral filters, dark signals, and attention sinks. ACL 2024  
Sun et al. Massive activations in large language models. COLM 2024

# Mechanism of attention sink

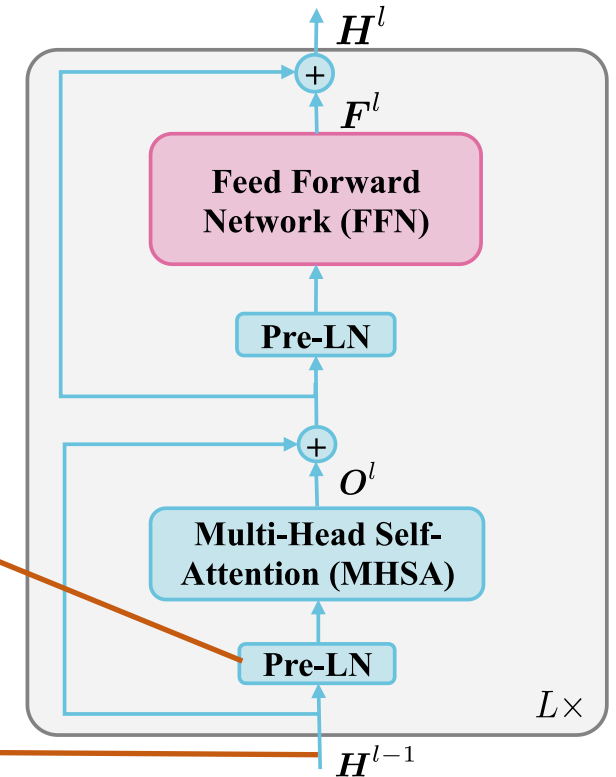
- What is the connection between attention sink and massive activations?

Layer norm retains values for spike dimensions for key of sink token

Key of the sink token is distributed in a different manifold (low-rank)



$$\mathbf{k}_t^{l,h} = \text{LN}(\mathbf{h}_t^{l-1}) \mathbf{W}_K^{l,h} \mathbf{R}_{\Theta, -t}$$
$$\text{LN}(\mathbf{h}) = \frac{\mathbf{h}}{\sqrt{\frac{1}{d} \sum_{i=1}^d h_i^2}} \odot \mathbf{g}$$



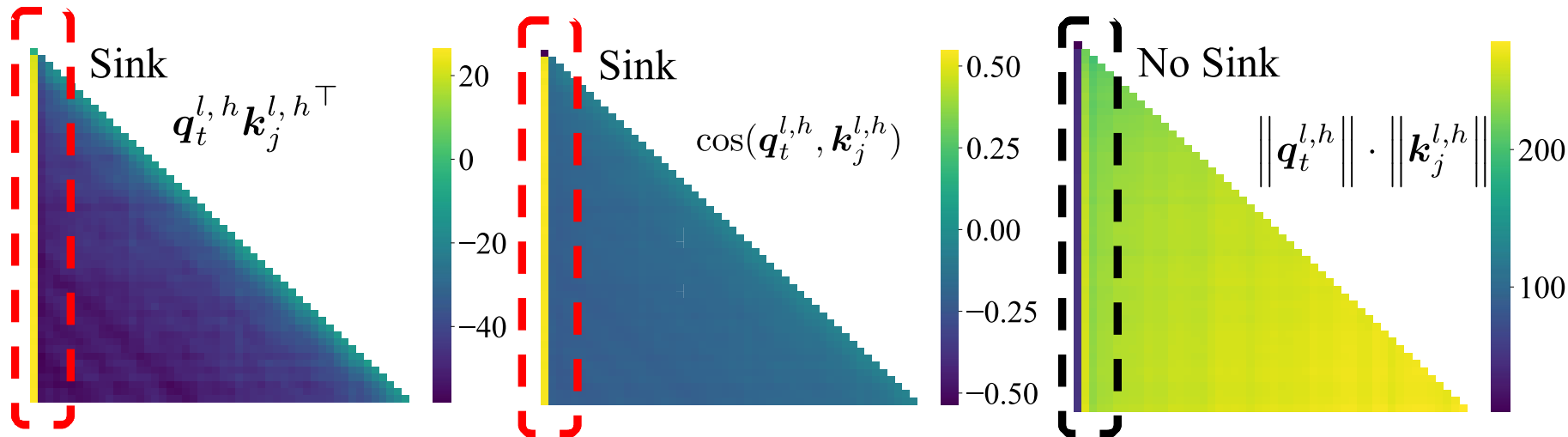
# Mechanism of attention sink

- What's the result of “first key is distributed in a different manifold”

Attention sink

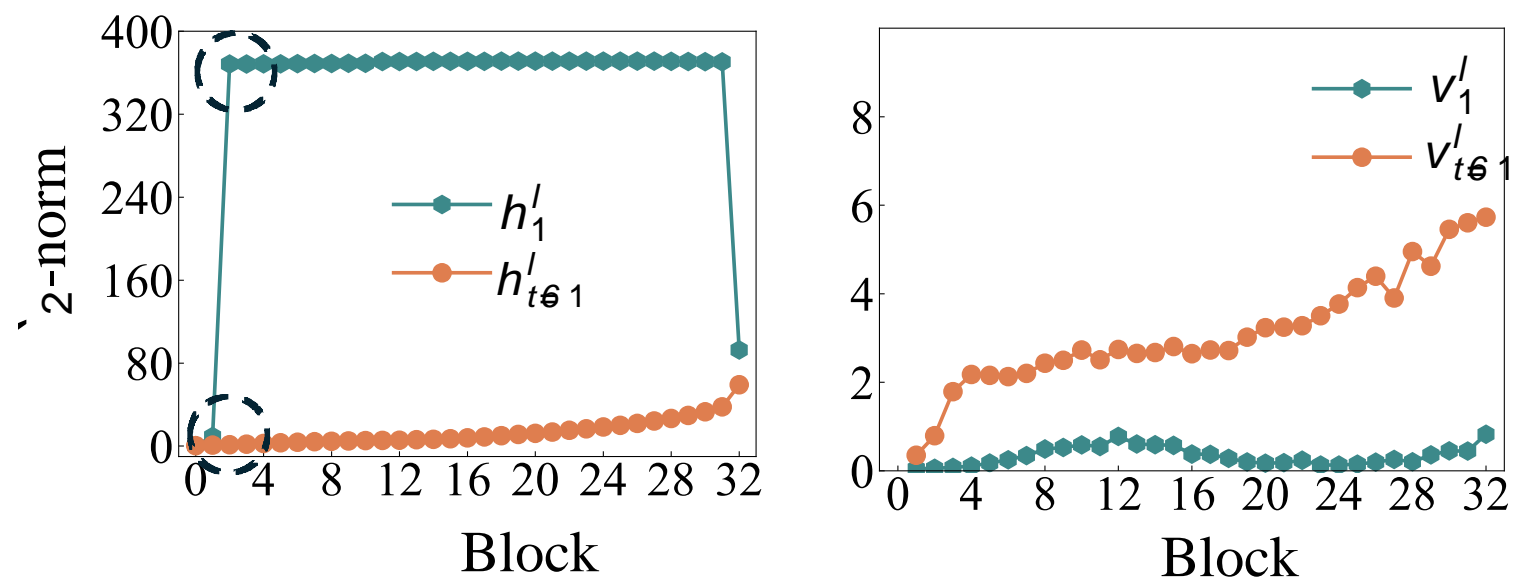
QK angle

$$\mathbf{q}_t^{l,h} \mathbf{k}_1^{l,h\top} \gg \mathbf{q}_t^{l,h} \mathbf{k}_{j \neq 1}^{l,h\top}$$
$$\cos(\mathbf{q}_t^{l,h}, \mathbf{k}_1^{l,h}) \gg \cos(\mathbf{q}_t^{l,h}, \mathbf{k}_{j \neq 1}^{l,h})$$



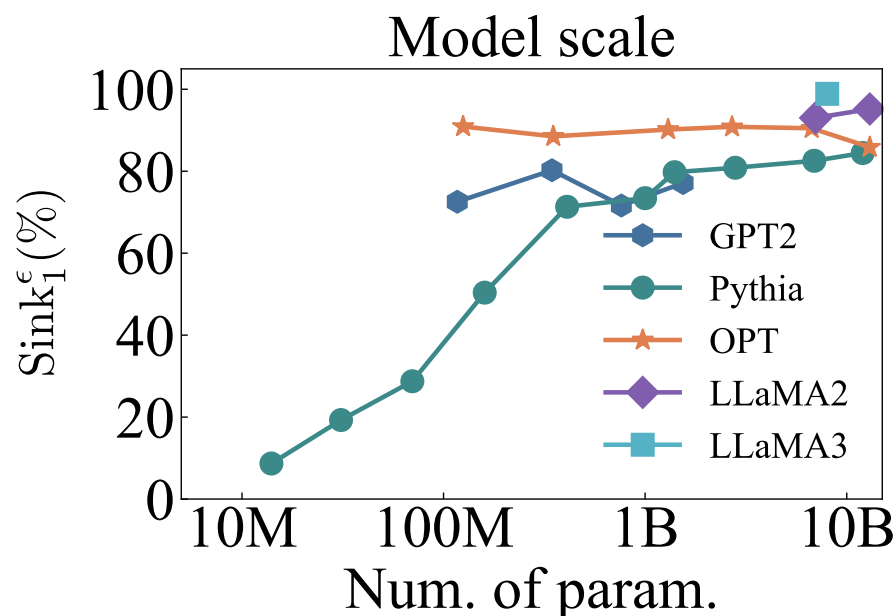
# Mechanism of attention sink

- Another property of sink token: small values



# Attention sink widespread appears in pre-trained LMs

- Attention sink appears widespread in various LMs, even in LMs with 14M params.
- Attention sink emerges in LM pre-training



- Sink metric refers to the number of attention sink heads in the whole LM

LLM	Sink <sub>1</sub> <sup>ϵ</sup> (%)	
	Base	Chat
Mistral-7B	97.49	88.34
LLaMA2-7B	92.47	92.88
LLaMA2-13B	91.69	90.94
LLaMA3-8B	99.02	98.85



# Attributing attention sink to LM pre-training

- LM pre-training objective  $\min_{\theta} \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}} [\mathcal{L}(p_{\theta}(\mathbf{X}))]$
- Experiments on LLaMA2-style models: Attributing model behavior to training recipes

Optimization

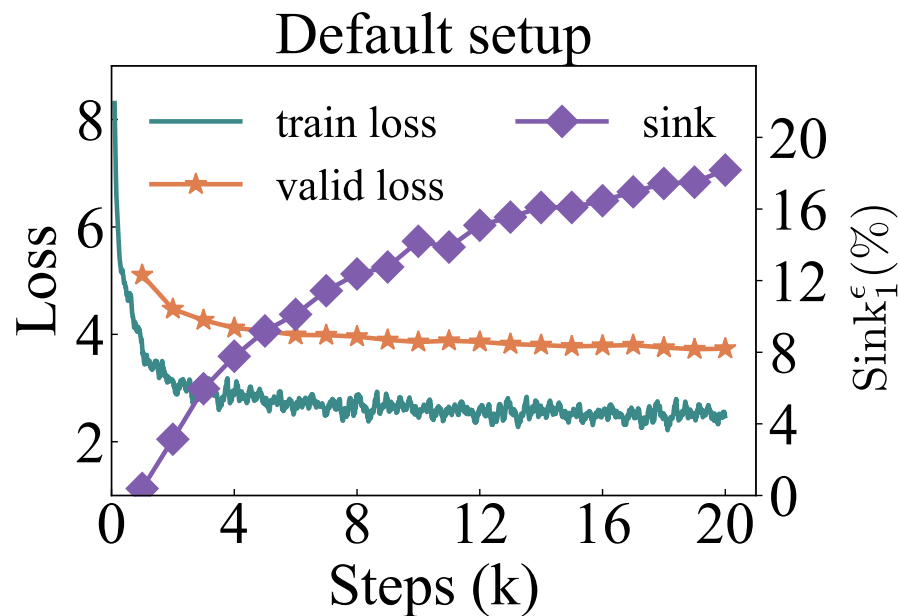
Data distribution

Loss function

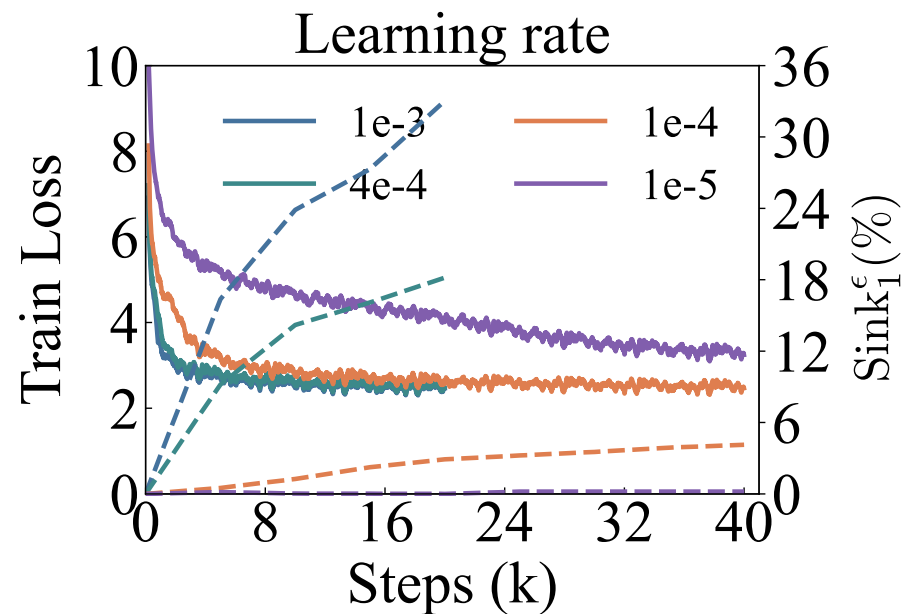
Model architecture

# Effects of optimization on attention sink

- Training steps



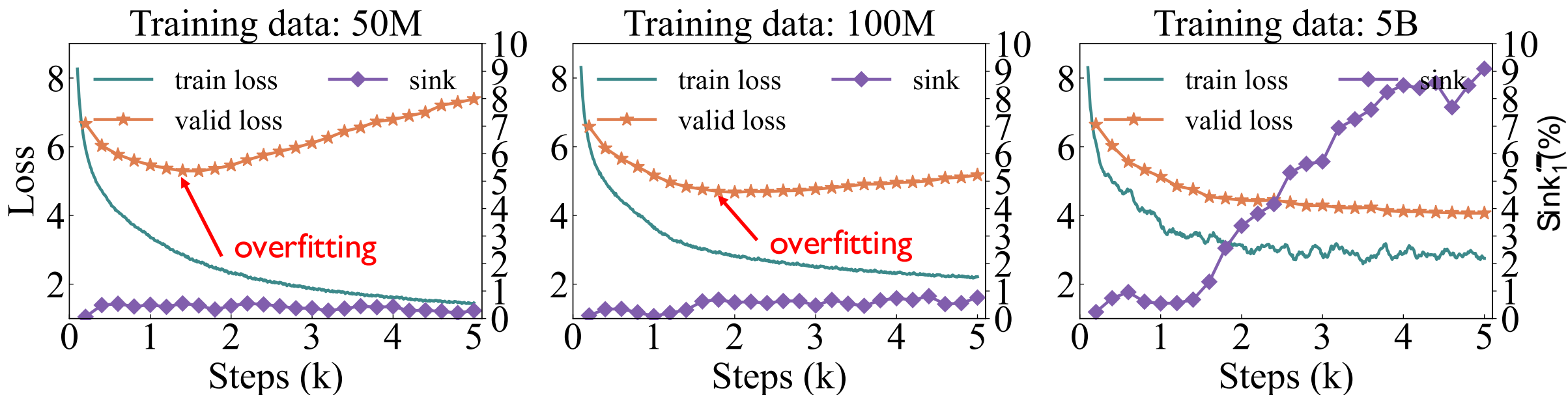
- Learning rate
- Small learning rates mitigate attention sink



# Effects of data distribution on attention sink

- Unique training data amount

Attention sink emerges after LMs are trained on **sufficient unique training data**, not really related to **overfitting**



# Effects of loss function on attention sink

- Auto-regressive loss

$$\mathcal{L} = \sum_{t=2}^C \log p_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t})$$

- Weight decay

$$\mathcal{L} = \sum_{t=2}^C \log p_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t}) + \gamma \|\theta\|_2^2$$

L2 regularization



Larger weight decay encourages attention sink

$\gamma$	0.0	0.001	0.01	0.1	0.5	1.0	2.0	5.0
Sink <sub>1</sub> <sup>ε</sup> (%)	15.20	15.39	15.23	18.18	41.08	37.71	6.13	0.01
valid loss	3.72	3.72	3.72	3.73	3.80	3.90	4.23	5.24

# Effects of model architecture on attention sink

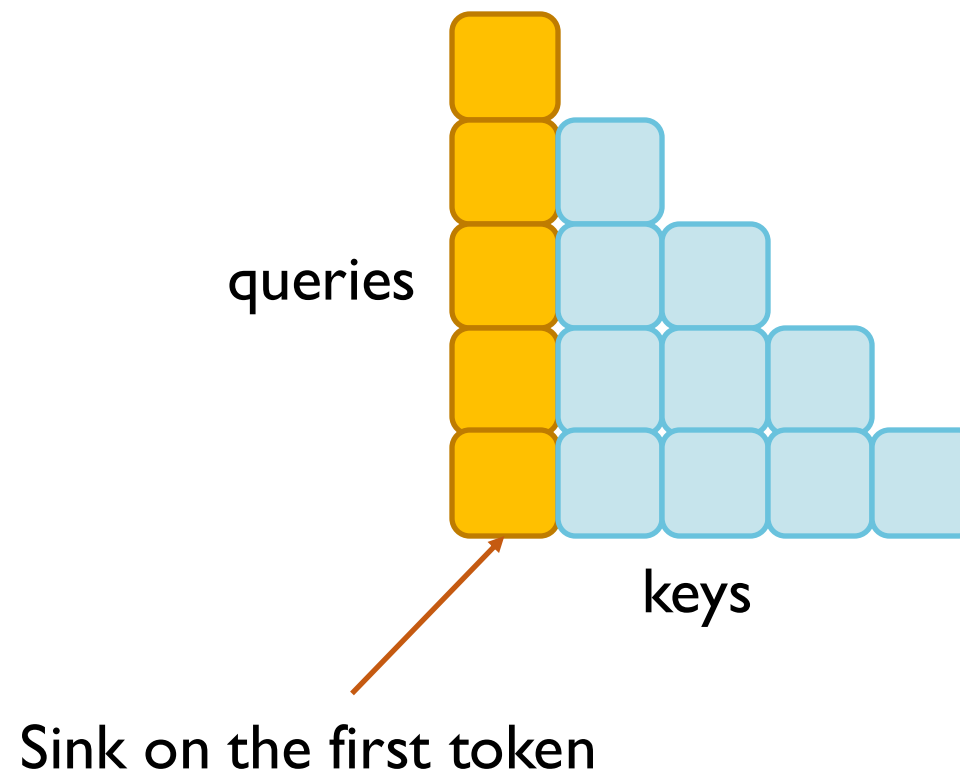
The following designs do not affect the emergence of attention sink

- Positional embeddings: including no positional embedding
- Pre-norm and post-norm design
- Feed forward networks (FFNs) with different activation functions
- Number of attention heads, how to combine multiple heads

# Effects of model architecture on attention sink

Standard softmax attention in  $h$ -th head  $l$ -th block

$$\text{Softmax} \left( \frac{1}{\sqrt{d_h}} \underset{\substack{\uparrow \\ \text{queries}}}{Q^{l,h}} \underset{\substack{\uparrow \\ \text{keys}}}{K^{l,h}{}^\top} + \underset{\substack{\uparrow \\ \text{casual mask}}}{M} \right) \underset{\substack{\uparrow \\ \text{values}}}{V^{l,h}}$$



# Effects of model architecture on attention sink

Softmax attention with learnable K biases

$$\text{Softmax} \left( \frac{1}{\sqrt{d_h}} Q^{l,h} \begin{bmatrix} \mathbf{k}^{*l,h\top} & \mathbf{K}^{l,h\top} \end{bmatrix} + \mathbf{M} \right) \begin{bmatrix} \mathbf{0} \\ \mathbf{V}^{l,h} \end{bmatrix}$$

Learnable K biases

V biases are all zeros

queries

keys

Sink on the learnable K bias

# Effects of model architecture on attention sink

- LM with K biases has no massive activations!
- Large attention score  $\neq$  important in semantic
- Sink token saves extra attention, adjusts the dependence among other tokens

Why need such a mechanism?

Is it because attention score added up to one?



# Effects of model architecture on attention sink

Attention output

$$\mathbf{v}_i^\dagger = \sum_{j=1}^i \frac{\text{sim}(\varphi(\mathbf{q}_i), \varphi(\mathbf{k}_j))}{\sum_{j'=1}^i \text{sim}(\varphi(\mathbf{q}_i), \varphi(\mathbf{k}_{j'}))} \mathbf{v}_j$$

$$\text{sim}(\varphi(\mathbf{q}_i), \varphi(\mathbf{k}_j)) = \exp\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_h}}\right)$$

softmax

$$\mathbf{Z}_i = \sum_{j'=1}^i \text{sim}(\varphi(\mathbf{q}_i), \varphi(\mathbf{k}_{j'}))$$

normalization term

Perhaps normalization matters, as it forces the attention scores sum to one?

# Effects of model architecture on attention sink

- Relax tokens' inner dependence by removing normalization

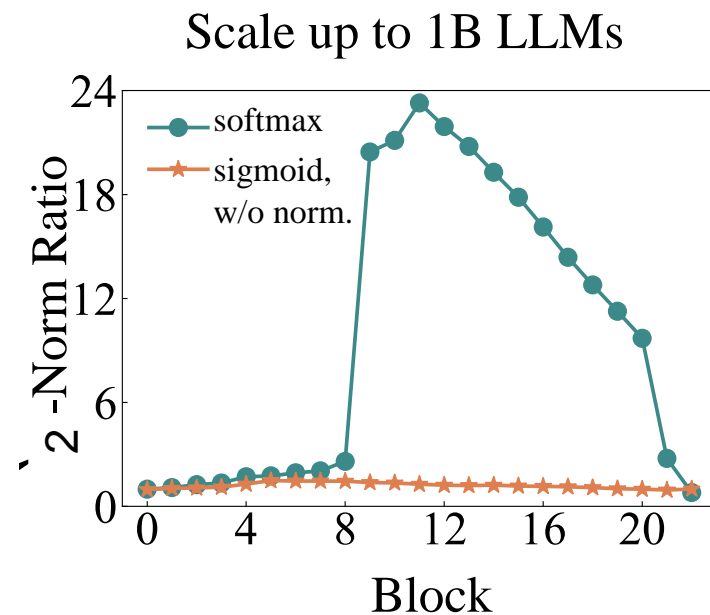
Sigmoid attention:

$$\mathbf{v}_i^\dagger = \sum_{j=1}^i \text{sigmoid}\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_h}}\right) \mathbf{v}_j$$

ELU plus one attention:

$$\mathbf{v}_i^\dagger = \sum_{j=1}^i \left( \text{elu}\left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d_h}}\right) + 1 \right) \mathbf{v}_j$$

No normalization -> No attention sink, no massive activations!  
Added back normalization -> Attention sink, massive activations!



# Effects of model architecture on attention sink

- Relax tokens' inner dependence by allowing negative attention scores

Linear attention, with a mlp kernel

$$\mathbf{v}_i^\dagger = \sum_{j=1}^i \frac{\text{mlp}(\mathbf{q}_i)\text{mlp}(\mathbf{k}_j)^\top}{\sqrt{d_h}} \mathbf{v}_j \rightarrow \text{No attention sink, no massive activations}$$

Add a normalization

$$\mathbf{Z}_i = \max \left( \left| \sum_{j'=1}^i \frac{\text{mlp}(\mathbf{q}_i)\text{mlp}(\mathbf{k}_{j'})^\top}{\sqrt{d_h}} \right|, 1 \right) \rightarrow \text{No attention sink, no massive activations}$$

# Takeaway

- Attention sink is a widespread phenomena across models and input
- Attention sink emerges during the LM pre-training
- Attention sink acts as key biases, storing extra attention and non-informative
- Softmax plays an important role in the emergence of attention sink

Please check our paper to see more interesting results!