#### Mechanistically Demystifying Extreme-Token Phenomena in Large Language Models

Song Mei UC Berkeley

Peter Hall conference @ UC Davis Statisitcs in the age of AI

Joint work with Tianyu Guo, Druv Pai, Yu Bai, Jiantao Jiao, Mike Jordan

#### Extreme values in transformer-based LLMs



[Dettmers et al., 2022] [Xiao et al., 2023] [Sun et al., 2024] [Guo et al., 2024]

0

-1

-2

-3

-4

-5

#### Extreme values in transformer-based LLMs





[Dettmers et al., 2022] [Xiao et al., 2023] [Sun et al., 2024] [Guo et al., 2024]

#### Extreme values in transformer-based LLMs



[Dettmers et al., 2022] [Xiao et al., 2023] [Sun et al., 2024] [Guo et al., 2024]

- A universal phenomena: happens in almost all open-source LLMs including GPT-2, Llama-2, Llama-3, Pythia, Mixtrial, etc.
- Troublemaking: inference, quantization, interpretability...

#### **Transformer-based LLMs**

- A transformer is a sequence-to-sequence neural network  $\mathsf{TF}_{\theta} : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$ .
- Input sequence:  $H = [h_1, h_2, ..., h_N] \in \mathbb{R}^{D \times N}$ ; each  $h_i \in \mathbb{R}^D$  is called a token.

#### **Transformer-based LLMs**

- A transformer is a sequence-to-sequence neural network  $\mathsf{TF}_{\theta} : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$ .
- Input sequence:  $H = [h_1, h_2, ..., h_N] \in \mathbb{R}^{D \times N}$ ; each  $h_i \in \mathbb{R}^D$  is called a token.



#### **Transformer-based LLMs**

- A transformer is a sequence-to-sequence neural network  $\mathsf{TF}_{\theta} : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$ .
- Input sequence:  $H = [h_1, h_2, ..., h_N] \in \mathbb{R}^{D \times N}$ ; each  $h_i \in \mathbb{R}^D$  is called a token.



• A transformer is an iterative composition of MLP layers and Attention layers  $TF_{\theta}(\cdot) = (Id + MLP_{W^{(L)}}) \circ (Id + ATTN_{A^{(L)}}) \circ \cdots \circ (Id + MLP_{W^{(1)}}) \circ (Id + ATTN_{A^{(1)}})$ 

- A transformer is an iterative composition of MLP layers and Attention layers  $TF_{\theta}(\cdot) = (Id + MLP_{W^{(L)}}) \circ (Id + ATTN_{A^{(L)}}) \circ \cdots \circ (Id + MLP_{W^{(1)}}) \circ (Id + ATTN_{A^{(1)}})$
- MLP layer:  $MLP_W : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$

 $H' = W_2 \cdot \sigma(W_1 H)$ 

- A transformer is an iterative composition of MLP layers and Attention layers  $TF_{\theta}(\cdot) = (Id + MLP_{W^{(L)}}) \circ (Id + ATTN_{A^{(L)}}) \circ \cdots \circ (Id + MLP_{W^{(1)}}) \circ (Id + ATTN_{A^{(1)}})$
- MLP layer:  $MLP_W : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$  $H' = W_2 \cdot \sigma(W_1 H)$
- Attention layer:

$$\mathsf{ATTN}_{A}: \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$$

$$H' = \sum_{m=1}^{M} O_m V_m H \cdot \text{Softmax}(\text{mask}(H^{\top} K_m^{\top} Q_m H))$$

- A transformer is an iterative composition of MLP layers and Attention layers  $TF_{\theta}(\cdot) = (Id + MLP_{W^{(L)}}) \circ (Id + ATTN_{A^{(L)}}) \circ \cdots \circ (Id + MLP_{W^{(1)}}) \circ (Id + ATTN_{A^{(1)}})$
- MLP layer:  $MLP_W : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$  $H' = W_2 \cdot \sigma(W_1 H)$
- Attention layer:

$$ATTN_{A} : \mathbb{R}^{D \times N} \to \mathbb{R}^{D \times N}$$
$$H' = \sum_{m=1}^{M} O_{m} V_{m} H \cdot \text{Softmax}(\text{mask}(H^{\top} K_{m}^{\top} Q_{m} H))$$















#### Extreme-token phenomena



#### Extreme-token phenomena



## Explain why the extreme-token phenomena appear in transformer-based LLMs (Static and dynamic)

Tianyu Guo, Druv Pai, Yu Bai, Jiantao Jiao, Michael I. Jordan, and Song Mei. Active-Dormant Attention Heads: Mechanistically Demystifying Extreme-Token Phenomena in LLMs. arXiv preprint, arXiv: 2410.13835.



Different attention pattern on Wiki/GitHub



Different attention pattern on Wiki/GitHub



When zero-out the head, the performance of the transformer on GitHub data drops a lot.



Different attention pattern on Wiki/GitHub



When zero-out the head, the performance of the transformer on GitHub data drops a lot.

#### L16H25 in Llama-2-7B is active on GitHub data and dormant on Wikipedia data

Bigram-Backcopy model

 $\mathcal{V} = \{\mathbf{t}, \mathbf{e}\} \cup \{a, b, \cdots d, f, \cdots, s, u, \dots\}$ 



**Bigram-Backcopy model** 

 $\mathcal{V} = \{ \underline{t}, \underline{e} \} \cup \{ a, b, \cdots d, f, \cdots, s, u, \dots \}$ 



Bigram-Backcopy model

 $\mathcal{V} = \{ \underline{t}, \underline{e} \} \cup \{ a, b, \cdots d, f, \cdots, s, u, \dots \}$ 



Bigram-Backcopy model

 $\mathcal{V} = \{ \underline{t}, \underline{e} \} \cup \{ a, b, \cdots d, f, \cdots, s, u, \dots \}$ 



Dataset = 
$$\{H^{(i)}\} \sim_{iid} H$$

Bigram-Backcopy model  $\mathcal{V} = \{\mathbf{t}, \mathbf{e}\} \cup \{a, b, \cdots d, f, \cdots, s, u, \dots\}$ Backcopy: Copy the backward token at trigger tokens Begin-of-sequece  $\langle s \rangle v t v n a u t u h \cdots$ H =**Bigram**: Sample the next token by Markov transition  $P(\cdot | \text{ previous token})$ 

Dataset = 
$$\{H^{(i)}\} \sim_{iid} H$$

Single-head single-layer transformer

$$H = \text{tokenize}(\langle s \rangle v t v n a u t u h)$$

$$ATTN(H) = OVH \cdot \text{Softmax}(\text{mask}(H^{\mathsf{T}}K^{\mathsf{T}}QH))$$

$$MLP(H) = W_1 \cdot \text{ReLU}(W_2H)$$

$$TF(H) \qquad H_{\text{shifted}}$$

$$= v t v n a u t u h$$

$$CrossEntropy$$

$$Loss$$

Training TF by running Adam algorithm on CrossEntropy(TF(H),  $H_{shifted}$ ) achieves optimal risk

#### The extreme-token phenomena in simplified model





**Attention sinks** 







Pre-trained transformer: Full model





Pre-trained transformer: No MLP





Pre-trained transformer: No Attn





Pre-trained transformer: Full model











Attention head is active on trigger letters and dormant on non-trigger letters





When zero-out the head, the performance of the transformer on GitHub data drops a lot.

L16H25 in Llama-2-7B is active on GitHub data and dormant on Wikipedia data

How does attention sinks and value-state drains form along pre-training dynamics?

#### The formation of attention sink along pre-training



### The formation of attention sink along pre-training



• Phase I: Bigram excess risk achieves near 0; Backcopy risk drops to a small value.

### The formation of attention sink along pre-training



- Phase I: Bigram excess risk achieves near 0; Backcopy risk drops to a small value.
- Phase II: Backcopy risk continues to decrease; attention sinks and value-state drains begin to form, and go extreme even if Backcopy risk has been very close to 0.

#### **The Mutual Reinforcement Mechanism**



#### **The Mutual Reinforcement Mechanism**



### **The Mutual Reinforcement Mechanism**



#### Theorem [Our result; Informal]:

If the model can accurately predict the next token without using the attention head, but adding any value state from previous tokens worsens the prediction, the attention head will become dormant and form an attention sink. Dynamically, this arises from a mutual reinforcement mechanism:

- The SoftMax function shifts attention weights towards tokens that exhibit valuestate drains, reinforcing these tokens as attention sinks.
- Attention sinks on these extreme tokens further suppress their value states, reinforcing their role as value-state drains.

## Examine mutual reinforcement in LLMs: the OLMo model

Models	arc challenge	arc easy	boolq	hella- swag	open bookqa	piqa	sciq	wino- grande	avg.
StableLM 1.6B	43.8	63.7	76.6	68.2	45.8	74.0	94.7	64.9	66.5
Pythia 1B	33.1	50.2	61.8	44.7	37.8	69.1	86.0	53.3	54.5
TinyLlama 1.1B	34.8	53.2	64.6	58.7	43.6	71.1	90.5	58.9	59.4
OLMo-1B	34.5	58.1	60.7	62.5	46.4	73.7	88.1	58.9	60.4
Falcon-7B	47.5	70.4	74.6	75.9	53.0	78.5	93.9	68.9	70.3
LLaMA 7B	44.5	67.9	75.4	76.2	51.2	77.2	93.9	70.5	69.6
Llama 2 7B	48.5	69.5	80.2	76.8	48.4	76.7	94.5	69.4	70.5
MPT-7B	46.5	70.5	74.2	77.6	48.6	77.3	93.7	69.9	69.8
Pythia 6.9B	44.1	61.9	61.1	63.8	45.0	75.1	91.1	62.0	63.0
<b>RPJ-INCITE-7B</b>	42.8	68.4	68.6	70.3	49.4	76.0	92.9	64.7	66.6
OLMo-7B	48.5	65.4	73.4	76.4	50.4	78.4	93.8	67.9	69.3

#### **Mutual Reinforcement Mechanism in OLMo**



#### **Mutual Reinforcement Mechanism in OLMo**

### The attention sinks start increase when value states norms start decreasing



#### **Mutual Reinforcement Mechanism in OLMo**

The attention sinks stop increase when value states norms stop decreasing



# Replace softmax to ReLU elimintates the extreme-token phenomena



In practice, people have trained small scale ReLU transformer (up to 1 Billion parameters), demonstrating similar performance as Softmax transformer.

### Summary



- Pretrained transformer-based LLMs exhibit extreme-token phenomena.
- As predicted from the Bigram-Backcopy model, these extreme-token phenomena are governed by the Active-Dormant mechanism and mutual reinforcement mechanism.
- Replacing Softmax to ReLU attention eliminate the extreme-token phenomena.

Tianyu Guo, Druv Pai, Yu Bai, Jiantao Jiao, Michael I. Jordan, and Song Mei. Active-Dormant Attention Heads: Mechanistically Demystifying Extreme-Token Phenomena in LLMs. arXiv preprint, arXiv: 2410.13835.