

---

# Active-Dormant Attention Heads: Mechanistically Demystifying Extreme-Token Phenomena in LLMs

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We investigate the mechanisms behind three puzzling phenomena observed in  
2 transformer-based large language models (LLMs): *attention sinks*, *value-state*  
3 *drains*, and *residual-state peaks*, collectively referred to the *extreme-token phenom-*  
4 *ena*. First, we demonstrate that these phenomena also arise in simpler architec-  
5 tures—transformers with one to three layers—trained on a toy model, the Bigram-  
6 Backcopy (BB) task. In this setting, we identify an *active-dormant mechanism* that  
7 causes attention heads to become attention sinks for certain domain-specific inputs  
8 while remaining non-sinks for others. We further develop a precise theoretical  
9 characterization of the training dynamics that lead to these phenomena, revealing  
10 that they are driven by a *mutual reinforcement mechanism*. By small interventions,  
11 we demonstrate ways to avoid extreme-token phenomena during pre-training. Next,  
12 we extend our analysis to pre-trained LLMs, including Llama and OLMo, revealing  
13 that many attention heads are governed by a similar active-dormant mechanism as  
14 in the BB task. We further show that the same mutual reinforcement mechanism  
15 drives the emergence of extreme-token phenomena during LLM pre-training. Our  
16 results study the mechanisms behind extreme-token phenomena in both synthetic  
17 and real settings and offer potential mitigation strategies.

## 18 1 Introduction

19 Recent analyses of transformer-based open-source large language models (LLMs), such as GPT-2  
20 [33], Llama-2 [41], Llama-3 [12], Mixtral [25], Pythia [4], and OLMo [18], have revealed several  
21 intriguing phenomena:

- 22 • **Attention sinks** [45]: In many attention heads, the initial token consistently attracts a large  
23 proportion of attention weights. In certain LLMs, other special tokens, such as the delimiter  
24 token, also draw significant attention. We refer to these as *sink tokens*.
- 25 • **Value state drains** [20]: The value states of sink tokens are consistently much smaller than  
26 those of other tokens.
- 27 • **Residual state peaks** [37]: The intermediate representations of sink tokens, excluding those  
28 from the first and last layers, exhibit a significantly larger norm than other tokens.

29 These phenomena often appear simultaneously, and we collectively refer to them as the **extreme-**  
30 **token phenomena**. Figure 1 illustrates these phenomena using a fixed prompt: “<bos> Summer  
31 is warm. Winter is cold.” in Llama-3.1-8B-Base, where the first token, <bos>, the beginning-of-  
32 sentence token, serves as the sink token. We note that the first token does not have to be <bos> to  
33 function as a sink token, as in GPT-2, where other tokens, being the initial token, can also serve this

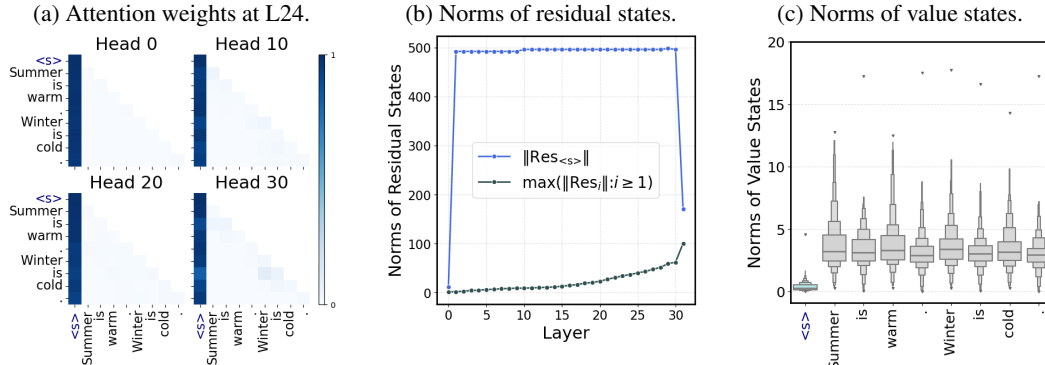


Figure 1: **Extreme-token phenomena in Llama 3.1-8B-Base.** We evaluate the sentence “<bos> Summer is warm. Winter is cold.” on the Llama 3.1-8B-Base model. *Left (a):* The value of the attention weights across multiple heads at Layer 24. We demonstrate that there are *attention sinks*: the key state associated with the <bos> token attracts the most attention from query states in these (and most) heads. *Middle (b):* The norm of the (residual stream) hidden states, measured at the output of each layer. We observe a *residual state peak* phenomenon: the <bos> token’s residual states have significantly larger norms than those of other tokens from layers 1 to 30. *Right (c):* The distribution of the norms of value states corresponding to each token at all layers and all heads. We observe the *value state drain* phenomenon: across many attention heads, the value state of the <bos> token is much smaller than those of other tokens on average.

34 role. Furthermore, in models like Llama-2, a delimiter token can also act as the sink token. Despite  
 35 the consistency of these observations, no prior work has provided a satisfying explanation for the  
 36 mechanisms behind these phenomena. As a tentative explanation, Xiao et al. [45] suggested that  
 37 models tend to dump unnecessary attention values to specific tokens.

38 This work aims to demystify the extreme-token phenomena in LLMs. We show that the extreme-token  
 39 phenomena are manifestations of the *active-dormant mechanism* of attention heads. We support  
 40 this claim through studies on simplified transformer architectures and tasks, a dynamical theory of  
 41 simplified models, and experiments on pre-trained LLMs. Our contributions are as follows:

- 42 1. In Section 2, we train one-to-three-layer transformers on the *Bigram-Backcopy* (BB) task, which  
 43 also exhibits extreme-token phenomena similar to those observed in LLMs. We show that  
 44 attention sinks and value-state drains are driven by the active-dormant mechanism mechanism.  
 45 Both theoretically and empirically, we demonstrate that the mutual reinforcement dynamics  
 46 underpin the extreme-token phenomena: attention sinks and value-state drains reinforce each  
 47 other, leading to a stable phase where all query tokens produce identical attention logits for the  
 48 keys of extreme tokens. Empirical evidence further shows that residual state peaks result from  
 49 the interaction between this mutual reinforcement mechanism and Adam.
- 50 2. In Section 3, we demonstrate the *active-dormant mechanism* mechanism in LLMs by identifying  
 51 an interpretable active-dormant head (Layer 16, Head 25 in Llama 2-7B-Base [41]), confirmed  
 52 through causal intervention analyses. We also discover circuits in LLMs related to extreme  
 53 tokens that partially align with models trained on the BB task. Examining the dynamics of  
 54 OLMo-7B-0424 [18], we observe the same mutual reinforcement mechanism and stable phase,  
 55 consistent with predictions from the BB task.
- 56 3. Through causal interventions, we isolate the extreme-token phenomena to architecture and  
 57 optimization strategy. Specifically, we show that replacing SoftMax with ReLU activations  
 58 in attention heads can eliminate extreme-token phenomena in the BB task, and switching  
 59 from Adam to SGD removes the residual-state peak phenomenon in the BB task. Our work  
 60 demonstrates potential classes of modifications to mitigate extreme-token phenomena in LLMs.

## 61 1.1 Notation

62 We denote the SoftMax attention layer with a causal mask as `attn`, the MLP layer as `m1p`, and the  
 63 transformer block as `TF`. The query, key, value states, and residuals of a token  $v$  are represented as  
 64  $\text{Qry}_v$ ,  $\text{Key}_v$ ,  $\text{Val}_v$ , and  $\text{Res}_v$ , respectively, with the specific layer and head indicated in context. We  
 65 use <bos> to refer to the "Beginning of Sequence" (bos) token. Throughout the paper, we employ

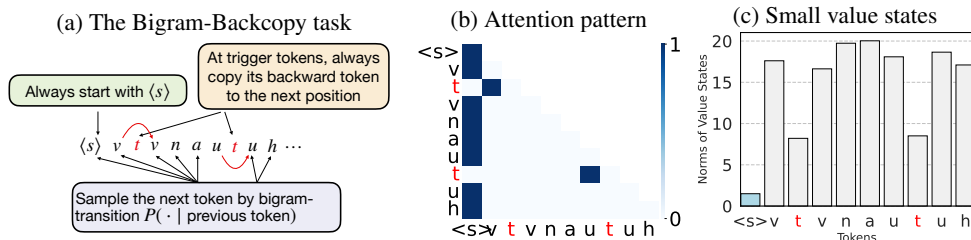


Figure 2: **Experiments on the Bigram-Backcopy task.** *Left (a):* We illustrate the data generation procedure for the Bigram-Backcopy task, where we fix 't', 'e', and the space character (' ') as trigger tokens. The BB task samples bigram transitions for non-trigger tokens and backcopies for trigger tokens. *Middle (b):* We present the attention weight heat map of a given prompt, with trigger tokens marked in red. Non-trigger tokens act as attention sinks. *Right (c):* We plot the value state norms for the prompt, where the <bos> token has a tiny norm.

66 zero-indexing (i.e., attention head and layer indices start from 0 rather than 1) for consistency between  
 67 code and writing.

## 68 2 The Bigram-Backcopy Task

69 The Bigram-Backcopy task consists of two sub-tasks: *Bigram-transition* and *Backcopy*. Each input  
 70 sequence begins with a <bos> token, followed by tokens sampled according to a pre-determined  
 71 bigram transition probability  $P$ . When some special trigger tokens are encountered, instead of  
 72 sampling, the preceding token is copied to the next position. Following Bietti et al. [5], we select the  
 73 transition  $P$  and the vocabulary  $\mathcal{V}$  with  $|\mathcal{V}| = V = 64$  based on the estimated character-level bigram  
 74 distribution from the tiny *Shakespeare* dataset. In all experiments, the set of trigger tokens  $\mathcal{T}$  is fixed  
 75 and consists of the  $|\mathcal{T}| = 3$  most frequent tokens in the unigram distribution. Thus, the non-trigger  
 76 token set,  $\mathcal{V} \setminus \mathcal{T}$ , comprises 61 tokens.

### 77 2.1 One-layer transformer shows attention sinks and value-state drains.

78 On the Bigram-Backcopy task, we pre-train a standard one-layer transformer with only one softmax  
 79 `attn` head and one `mpl` layer. Unless otherwise specified, the model is trained with Adam for 10,000  
 80 steps. We relegate the training details in Appendix C. Figure 2b shows that the trained transformer  
 81 also exhibits the attention sink phenomenon, where the <bos> token captures a significant proportion  
 82 of the attention weights. More importantly, the attention weights reveal interpretable patterns: all  
 83 non-trigger tokens exhibit attention sinks, while the attention for trigger tokens is concentrated on  
 84 their preceding positions. Furthermore, Figure 2c reveals a value state drain phenomenon similar  
 85 to LLMs, indicating that on non-trigger tokens, the `attn` head adds a minimal value to the residual  
 86 stream.

87 **The active-dormant mechanism of the attention head:** Inspired by the observed interpretable  
 88 attention weight patterns, we propose the *active-dormant mechanism*. For any given token, an  
 89 attention head is considered *active* if it contributes significantly to the residual state, and *dormant* if  
 90 its contribution is minimal. As illustrated in Figure 2b, trained on the BB task, the attention head is  
 91 active on trigger tokens and dormant on non-trigger tokens.

92 Figure 3a demonstrates that the `mpl` layer is responsible for the Bigram task whereas the `attn` head  
 93 takes care of the Backcopy task. When the `mpl` layer is zeroed out, the backcopy loss remains signifi-  
 94 cantly better than a random guess, but the bigram loss degrades to near-random levels. Conversely,  
 95 when the `attn` layer is zeroed out, the backcopy loss becomes worse than a random guess, while the  
 96 bigram loss remains unaffected. This suggests that on trigger tokens, the `attn` head is active and  
 97 handles the backcopy task, whereas on non-trigger tokens, the `attn` head is dormant, allowing the  
 98 `mpl` layer to handle the Bigram task. We summarize the active-dormant mechanism of the `attn` head  
 99 in Claim 1.

100 **Claim 1.** *In the BB task, the `attn` head demonstrates active-dormant mechanism, alternating*  
 101 *between two phases:*

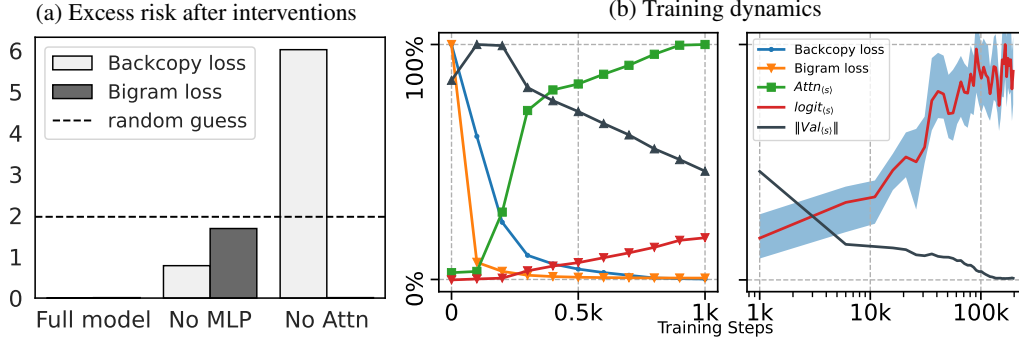


Figure 3: **Interventions and dynamics of one-layer transformer on the Bigram-Backcopy task.** *Left (a):* We display the excess risks for a one-layer model trained on the Bigram-Backcopy (BB) task under various interventions. *Right (b):* We plot the excess risks, attention weights, attention logits, and value state norms for the  $\langle \text{bos} \rangle$  token along the training dynamics. Each curve is rescaled to fall within a 0 to 1 range, though the trends remain consistent without rescaling. On the right side of (b), the horizontal axis is logarithmically scaled. The  $\text{logit}_{\langle \text{bos} \rangle}$  curve denotes the mean of attention logits from all given non-trigger query tokens  $v$  on the  $\langle \text{bos} \rangle$  token, normalized by the mean of attention logits on other tokens. The shaded area gives the 90% confidence interval on the distribution over all non-trigger tokens.

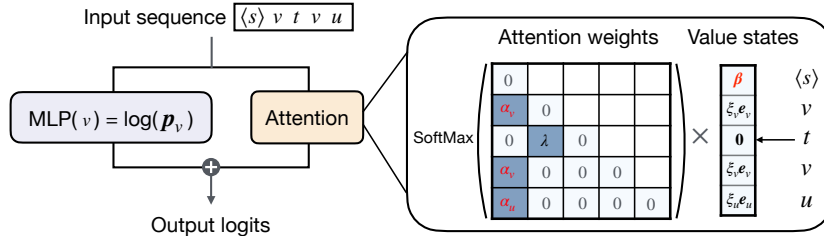


Figure 4: **The simplified transformer architecture with one MLP-layer and one attn head in parallel.** The predicted probability is the softmax of the output. Assume that the trainable variables are  $(\alpha, \beta) \in \mathbb{R}^V \times \mathbb{R}^V$ , which stands for the attention logits and value states of the  $\langle \text{bos} \rangle$  tokens.

- 102 • **Dormant phase:** On non-trigger tokens, the attn head puts dominant weights to the  $\langle \text{bos} \rangle$  token, adding minimal value to the residual stream, having little impact on the model's output.
- 103
- 104 • **Active phase:** On trigger tokens, the attn head puts dominant weights to the relevant context tokens, adding substantial value states to the residual stream, resulting in a significant impact on the model's output.
- 105
- 106

107 **The growth of attention logits on the  $\langle \text{bos} \rangle$  token and the decrease in the norm of its value state.** Figure 3b displays the training dynamics of excess risks, attention weights, attention logits, and value state norms for the  $\langle \text{bos} \rangle$  token. All values are rescaled to highlight the trends. The backcopy excess risk and the bigram excess risk both drop to zero within the first 1000 steps. As the backcopy risk decreases, the attention weights on the  $\langle \text{bos} \rangle$  token increase, suggesting a relationship between the formation of attention sinks and the functional development of the attention heads. For each token  $v_n$  at position  $n$  in the prompt, we compute  $\text{logit}_{\langle \text{bos} \rangle} = \text{mean}_n[\langle \text{Qry}_{v_n}, \text{Key}_{\langle \text{bos} \rangle} \rangle - \text{mean}_i[\langle \text{Qry}_{v_n}, \text{Key}_{v_i} \rangle]]$ , which serves as a progress measure for attention sinks. Even after the attention weights on the  $\langle \text{bos} \rangle$  token is nearly 1,  $\text{logit}_{\langle \text{bos} \rangle}$  continues to increase. Simultaneously, the norm of the value state of the  $\langle \text{bos} \rangle$  token continues to decrease to a small value.

## 117 2.2 Analysis of a minimally-sufficient transformer architecture

118 In this section, we analyze the training dynamics on the BB task by simplifying the architecture while preserving the attention sinks and value state drains phenomena. Let  $\mathcal{V}$  denote the set of all tokens except the  $\langle \text{bos} \rangle$  token, and  $\mathcal{T}$  denote the set of all trigger tokens. Given any  $v \in \mathcal{V}$ , we denote  $p_{vk} = P(k|v)$  to be the next token Markov transition probability, and  $\mathbf{p}_v = [p_{v1}, \dots, p_{vV}]$

122 be the row vector in the simplex. We assume that the tokens are embedded into  $V$ -dimensional  
123 space using one-hot encoding, and for notation simplicity, we abuse  $v$  to stand for its one-hot  
124 encoding vector  $e_v \in \mathbb{R}^V$  which is a row vector. The predicted probability of the  $n+1$  token is  
125 given by  $\text{SoftMax}(\text{TF}([\langle \text{bos} \rangle; v_{1:n-1}; v]_n))$ , where transformer architecture is given by  $\text{TF}(\cdot) =$   
126  $\text{attn}(\cdot) + \text{mlp}(\cdot)$ . Here  $\text{attn}(\cdot) = \text{SoftMax}(\text{mask}(\text{Qry}(\cdot)\text{Key}(\cdot)^\top))\text{Val}(\cdot)$  and  $(\text{Qry}, \text{Key}, \text{Val})$   
127 are linear maps from  $\mathbb{R}^V \rightarrow \mathbb{R}^V$ . Since the  $\text{mlp}$  layer could handle the Bigram task, we assume that  $\text{mlp}$   
128 outputs the Markov transition probabilities  $\mathbf{p}_v$  on non-trigger tokens  $v$  and zero on trigger tokens. For  
129 the  $\text{attn}$  head, we assume that the attention logits on the  $\langle \text{bos} \rangle$  key-token are  $(\alpha_{v_1}; \dots; \alpha_{v_n})$ , the  
130 attention logits on any trigger query-token are  $(0, \dots, \lambda, 0)$  where the second last coordinate is  $\lambda$ , and  
131 assume other logits are zero. Assume that the value state of  $\langle \text{bos} \rangle$  is  $\beta \in \mathbb{R}^V$ , and the value state of  
132 each non-trigger token  $v$  is a one-hot encoding vector  $e_v$  multiplied by  $\xi_v \geq 0$ . Figure 4 illustrates this  
133 simplified transformer architecture. These assumptions are summarized in the following equations.

$$\begin{aligned} \text{mlp}(v) &= \log \mathbf{p}_v \cdot \mathbf{1}\{v \notin \mathcal{T}\} \quad \text{for } v \in \mathcal{V}, \\ \langle \text{Qry}(v), \text{Key}(\langle \text{bos} \rangle) \rangle &= \alpha_v \cdot \mathbf{1}\{v \notin \mathcal{T}\} \quad \text{for } v \in \mathcal{V}, \\ \langle \text{Qry}(v), \text{Key}(v') \rangle &= \lambda \cdot \mathbf{1}\{v \in \mathcal{T}, v' \text{ is the former token of } v\} \quad \text{for } v, v' \in \mathcal{V}, \\ \text{Val}(v) &= \xi_v e_v \quad \text{with } \xi_v = 0 \text{ for } v \in \mathcal{T}, \text{ and } \xi_v \geq 0 \text{ for } v \in \mathcal{V} \setminus \mathcal{T}. \end{aligned} \tag{1}$$

134 Theorem 2 demonstrates the existence of a transformer structure that is equivalent to the simplified  
135 version. We relegate the proof in Section B.

136 **Theorem 2.** *For any parameters  $(\alpha \in \mathbb{R}^V, \beta \in \mathbb{R}^V, \xi \in \mathbb{R}^V, \lambda \in \mathbb{R})$ , there is a one-layer  
137 transformer  $(\text{mlp}, \text{Qry}, \text{Key}, \text{Val})$  such that Eq. (1) holds. The transformer gives ground truth  
138 transition of the BB model if  $\min_{v \in \mathcal{V}} \alpha_v \rightarrow \infty$ ,  $\min_{v \in \mathcal{V}} \xi_v \rightarrow \infty$ ,  $\lambda \rightarrow \infty$ , and  $\beta = 0$ .*

139 Throughout we adopt Eq. (1) as our assumption. We further define  $W_k = \sum_{i=1}^n \mathbf{1}\{v_i = k\}$ ,  
140  $\mathbf{W} = (W_1, \dots, W_V)$ , and  $W = \sum_{k \in \mathcal{V}} W_k = n$ . Then for a non-trigger token  $v$ , the output of  
141 attention layer with input sequence  $[\langle \text{bos} \rangle; v_{1:n-1}; v]$  gives (denoting  $\xi_k = 0$  for  $k \in \mathcal{T}$ )

$$\text{TF}([\langle \text{bos} \rangle; v_{1:n-1}; v]_n) = \log \mathbf{p}_v + \frac{e^{\alpha_v}}{e^{\alpha_v} + W} \beta + \sum_{k=1}^V \frac{W_k \xi_k}{e^{\alpha_v} + W} \cdot e_k.$$

142 Therefore, on the non-trigger token  $v$ , the cross-entropy loss between the true Markov transition  $\mathbf{p}_v$   
143 and predicted transition  $\text{SoftMax}(\text{TF}([v_{1:n-1}; v]_n))$  is given by

$$\text{loss}_v(\alpha_v, \beta) = \sum_{k=1}^V p_{vk} \left\{ \log \left[ \sum_{i=1}^V p_{vi} \exp \left( \frac{e^{\alpha_v} \beta_i + W_i \xi_i}{e^{\alpha_v} + W} \right) \right] - \frac{e^{\alpha_v} \beta_k + W_k \xi_k}{e^{\alpha_v} + W} - \log p_{vk} \right\}.$$

144 For simplicity, we neglect the loss on trigger tokens and assume that  $(\{W_i\}_{i \in [V]}, W)$  are fixed  
145 across different positions in the input sequences<sup>1</sup>, and consider the total loss to be the losses on each  
146 non-trigger token averaged with its proportion in the stable distribution  $\{\pi_v\}_{v \in \mathcal{V}}$ , given by

$$\text{loss}(\alpha, \beta) = \sum_{v \in \mathcal{V} \setminus \mathcal{T}} \pi_v \text{loss}_v(\alpha_v, \beta).$$

147 **Theorem 3.** *Consider the gradient flow of the loss function  $\text{loss}(\alpha, \beta)$ . Assume  $\xi_v \geq 0$  for any  $v$ ,  
148 and  $\{W_i \cdot \xi_i\}_{i \in \mathcal{V}}$  are not all equal.*

149 • (Attention logits grow logarithmically reinforced by small value states) Fix  $\beta = \beta \cdot \mathbf{1}$  for a  
150 constant  $\beta$ , and consider the gradient flow over  $\alpha$ . With any initial value  $\alpha(0)$ , there exists  $\mathbf{r}(t)$   
151 with norm uniformly bounded in time such that

$$\alpha(t) = \frac{1}{2} \log t \cdot \mathbf{1} + \mathbf{r}(t).$$

152 • (Value state shrinks to a small constant vector reinforced by large attention logits) Fix  $\alpha = \alpha \cdot \mathbf{1}$   
153 for a constant  $\alpha$ , and define  $\bar{\beta}(0) = V^{-1} [\sum_v \beta_v(0)]$ . Consider the gradient flow over  $\beta$ . As  
154  $t \rightarrow \infty$ , we have

$$\beta(t) \rightarrow \beta^* = \bar{\beta}(0) \cdot \mathbf{1} - e^{-\alpha} \cdot \mathbf{W} \circ \xi.$$

<sup>1</sup>We note that [34] makes similar simplification in analyzing induction heads.

155 • (Stable phase: identical attention logits) Consider the gradient flow over variables  $(\alpha, \beta)$ . Any  
 156 vector of the following form

$$\alpha = \alpha \cdot \mathbf{1}, \quad \beta = c \cdot \mathbf{1} - e^{-\alpha} \cdot \mathbf{W} \circ \xi, \quad \alpha, c \in \mathbb{R}$$

157 is a stationary point. These are all global minimizers of  $\text{loss}(\alpha, \beta)$ .

158 The proof of Theorem 3 is provided in Appendix B.2. We give three key remarks: (1) As  $\alpha_v \rightarrow \infty$ , a  
 159 Taylor expansion of the gradient  $\partial \text{loss} / \partial \alpha_v$  suggests that  $d\alpha_v / dt \propto \exp(-2\alpha_v)$ , which leads to the  
 160 logarithmic growth of  $\alpha_v$ . Similar logarithmic growth exists in the literature under different setups  
 161 [39, 22]. (2) For a fixed  $\alpha = \alpha \mathbf{1}$ , under additional assumptions on the initial value  $\beta(0)$ , we can  
 162 prove a linear convergence for  $\beta$ . (3) The stable phase described in Theorem 3 seems to imply that  
 163 the system could be stable without attention sinks, as it does not require  $\alpha$  to be large. However, in  
 164 practice, models trained on the BB task tend to converge to a stable phase where  $\alpha$  is relatively large.

165 **The Formation of Attention Sinks and Value State Drains.** When  $\beta = \mathbf{0}$ , the attention logits on  
 166 the `<bos>` token increase monotonically. This demonstrates that the presence of a small value state  
 167 of the `<bos>` token reinforces the formation of attention sinks. When  $\alpha = \alpha \cdot \mathbf{1}$ , with  $\alpha$  sufficiently  
 168 large,  $\beta(t) \rightarrow \bar{\beta}(0)\mathbf{1}$ . Given the random Gaussian initialization,  $\|\bar{\beta}(0)\mathbf{1}\|_2 \approx \|\beta(0)\|_2 / \sqrt{d}$ , where  
 169  $d$  is the hidden dimension. This demonstrates that the presence of attention sinks reinforces the  
 170 formation of value states drains.

171 **Experimental verification.** Revisiting Figure 3b, which shows the dynamics of a full transformer  
 172 model trained with Adam, we observe that both  $\text{logit}_{\langle \text{bos} \rangle}$  and  $\|\text{Val}_{\langle \text{bos} \rangle}\|_2$  exhibit growth rates  
 173 consistent with Theorem 3. The  $\text{logit}_{\langle \text{bos} \rangle}$  is equivalent to  $\alpha$  in this context, as all other attention  
 174 logits are assumed to be zero under the setup of Theorem 3. When plotted on a logarithmic scale, the  
 175  $\text{logit}_{\langle \text{bos} \rangle}$  curve grows approximately linearly between 1,000 and 10,000 steps, then accelerates before  
 176 stabilizing around 100,000 steps. Meanwhile, the norm of the value state decreases monotonically.  
 177 The simultaneous increase in attention weights and decrease in value-state norms suggest that these  
 178 phases occur together during the training process. To further validate Theorem 3, we construct a  
 179 simplified model that aligns with Equ. (1), and train the parameters  $(\alpha \in \mathbb{R}^V, \beta \in \mathbb{R}^V, \xi \in \mathbb{R}^V, \lambda \in$   
 180  $\mathbb{R})$  with Adam. The resulting training curves are similar to those of a one-layer transformer, also  
 181 exhibiting the mutual reinforcement mechanism.

182 Combining theoretical insights and experimental evidence, we summarize the formation of attention  
 183 sinks and value state drains as a mutual reinforcement mechanism.

184 **Claim 4** (Mutual reinforcement mechanism). *For any attention head given a specific prompt, if*  
 185 *the model can accurately predict the next token without the attention head, but adding any value*  
 186 *state from previous tokens worsens the prediction, the attention head becomes dormant, forming an*  
 187 *attention sink, leading to the mutual reinforcement of attention sinks and value state drains:*

- 188 1. *The SoftMax mechanism pushes the attention weights to the value state drains, reinforcing*  
 189 *attention sinks.*
- 190 2. *The attention sinks on the value state drains further pushes down the value state, reinforcing*  
 191 *value state drains.*

192 *The mutual reinforcement stabilizes at the phase when all tokens have identical large attention*  
 193 *logits on the value state drains. Finally, due to the causal mask, the training dynamics favor the*  
 194 *<bos> token to become an extreme token.*

195 We expect that the formation of extreme tokens in LLMs follows a similar mutual reinforcement  
 196 mechanism. Indeed, although Theorem 3 focuses on a specific BB task with a simplified architecture  
 197 and loss function, the same principles can be applied to more general scenarios. Specifically, for  
 198 an attention head `attn`, we assume that  $(\text{LLM} \setminus \text{attn})(v) = \log p_v$ , meaning that the LLM, even  
 199 if we zeroed out `attn`, can still output an accurate next token prediction. Furthermore, we assume  
 200  $\text{Val}(v) = \xi_v e_v$ , indicating that adding the value state from any previous tokens performs a specific  
 201 function. Under these assumptions, we expect the same theoretical results to apply to LLMs. In  
 202 Section 3, we will explore the formation of attention sinks and value state drains along the training  
 203 dynamics of LLMs, where we find empirical evidence that aligns with the theory.

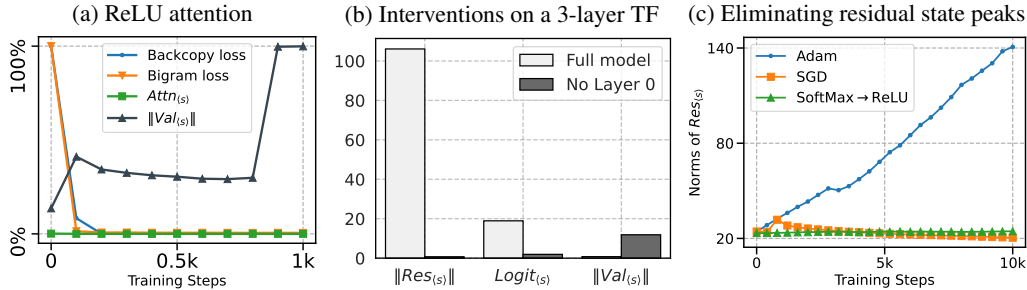


Figure 5: **Experiments on massive norms with multi-layer transformers trained on the Bigram-Backcopy task.** *Left (a):* We present the training dynamics of the ReLU attention for the first 1,000 steps. *Middle (b):* We plot the intervention results on the `attn+mlp+attn+mlp+mlp` structure. *Right (c):* We plot the evolution of massive norms in a three-layer transformer trained with Adam, SGD, and using a ReLU attention structure. Notably, only the three-layer model with softmax attention trained using Adam results in the emergence of residual state peaks.

204 **Replacing SoftMax by ReLU attention removes extreme-token phenomena.** As an implication  
 205 of our theory, we predict that training with ReLU attention instead of SoftMax attention will eliminate  
 206 the extreme-token phenomena. Without the SoftMax, the dynamics no longer push the attention  
 207 weights on the `<bos>` token, which remains zero along the training dynamics. Without attention sink,  
 208 the dynamics no longer push down the value state norm, and the mutual reinforcement mechanism  
 209 breaks. Figure 5a illustrates the training experiment on the BB task replacing SoftMax with ReLU,  
 210 showing that both the Bigram and Backcopy risk match the Bayes risk after 200 training steps, but the  
 211 attention logits of `<bos>` do not grow, and the value state does not shrink, confirming the prediction.

### 212 2.3 The emergence of residual state peaks

213 **The residual state peaks require a three-layer structure.** No residual state peaks appear in a  
 214 one-layer transformer trained on the BB task. We train various models on the BB task and track the  
 215 `<bos>` token’s residual state norms after layer 0. We relegate the experimental results to Appendix  
 216 C. We find that a three-layer transformer is enough to produce residual state peaks. If we allow to  
 217 skip some `mlp` or `attn` layers, the “`attn+mlp+attn+mlp+mlp`” combination becomes the simplest  
 218 model that produces residual state peaks (Figure 10). Circuit analysis also reveals that LLMs typically  
 219 add a large vector in the first layer and cancel it in the last layer. We propose that the add-then-cancel  
 220 mechanism is essential for residual state peaks and requires at least three layers.

221 **Residual state peak reinforces attention sinks and value state drains in trained models.** Figure  
 222 5b presents the intervention results on the “`attn+mlp+attn+mlp+mlp`” model. We recenter the  
 223  $\|Res_{<bos>}\|_2$  by subtracting the average norm of other tokens from the `<bos>` token norm. The  
 224  $logit_{<bos>}$  and  $\|Val_{<bos>}\|$  are computed in layer 1 following the same ways as in Figure 3b. When  
 225 layer 0 is zeroed out, the residual norm returns to normal, attention logits decrease, and the value  
 226 state norm rises. It verifies that the residual state peak contributes to the attention sink and value state  
 227 drain phenomenon in the trained transformer.

228 **Replacing Adam by SGD removes the linear growth of residual state norm.** Figure 5c shows  
 229 the `<bos>`’s residual state norms at the output of layer 0 of three-layer transformers with different  
 230 configurations. Adam leads to a linear increase in residual norms. In contrast, with SGD, attention  
 231 sinks persist, but residual state peaks vanish. The ReLU attention, which lacks the active-dormant  
 232 mechanism, shows no residual state peaks.

## 233 3 Extending Predictions of the BB Model to LLMs

234 In this section, we examine extreme-token phenomena in open-source pre-trained LLMs. In Sec-  
 235 tion 3.1, we analyze the static behavior of these phenomena in Llama 2-7B-Base [41], confirming  
 236 that certain attention heads in LLMs exhibit both active and dormant phases. Notably, we identify  
 237 a specific head that is active on GitHub samples but dormant on Wikipedia samples, illustrating

238 the *active-dormant mechanism*. In Section 3.2, we explore the dynamic behavior of extreme-token  
239 phenomena during the pre-training process of OLMo-7B [18]. We show that the attention logits,  
240 value state norms, and residual state norms of the sink token(s) in OLMo mirror their behavior in the  
241 simpler BB model. Specifically, the simultaneous formation of attention sinks and value state drains  
242 gives evidence for the *mutual reinforcement mechanism*.

### 243 3.1 Active-Dormant Mechanism in LLMs

244 Our study of the BB model leads to the following prediction about the extreme-token phenomena,  
245 which we hypothesize also applies to LLMs:

246 *Attention heads are controlled by an active-dormant mechanism. Attention sinks and value state*  
247 *drains indicate that an attention head is in dormant phase.*

248 This hypothesis suggests that in LLMs, attention heads become sinks or not depending on the  
249 context: the value vector can be totally non-informative towards picking likely next tokens for token  
250 distributions (e.g., tasks) in a particular context but not in others. This is a concrete instantiation  
251 vis-a-vis large-scale LLMs of the active-dormant dichotomy in Section 2, where this phenomenon  
252 was shown to occur in the context of small next-token predictors and the BB task.

253 Accordingly, we strive to find instances of heads in pretrained LLMs which satisfy this principle, i.e.,  
254 which are dormant on some domains and active on others. In Figure 6, we show a particular attention  
255 head – Layer 16 Head 25 of Llama 2-7B-Base [41] — which has an extremely clear active-dormant  
256 distinction across two distinct contexts (e.g., tokens from RedPajama [8] drawn from the GitHub  
257 subset versus the Wikipedia subset). While there are many such attention heads which are context-  
258 dependent — we provide some in Appendix D — we demonstrate this one because the conditions  
259 under which it is active are simple and interpretable, while others have more involved or complex  
260 criteria to become active. We observe that this attention head is *dormant* (i.e., an attention sink) on  
261 samples from Wikipedia, which more closely resemble prose, and *active* (i.e., not an attention sink)  
262 on samples from Github, which more closely resemble code. We also observe that this attention  
263 head, in general, contributes significantly to the performance of the model on code sequences, but has  
264 negligible impact on the performance of the model on prose sequences (Figure 6b). This is a further  
265 justification, from a practical perspective, of why this head is sometimes dormant and sometimes  
266 active — in some contexts we can ablate it from the model entirely with no effect, but in other  
267 contexts ablating the head leads to huge performance drops. We include more detail in Appendix E,  
268 where we extract a circuit for extreme-token phenomena in order to analyze the dormant-active  
269 mechanism and its interaction with the semantics of the input tokens.

### 270 3.2 Training Dynamics of Extreme-Token Phenomena in LLMs

271 Our study of the BB model leads to the following prediction about the dynamical behavior of the  
272 extreme-token phenomena, which we hypothesize also applies to LLMs:

273 *The attention heads go through a attention-increasing and value-state-shrinking phase. They then go*  
274 *into a stable phase, with identical attention logits on the <bos> token. Meanwhile, the residual state*  
275 *norm of the <bos> token linearly increases during pre-training.*

276 We confirm these predictions below. To observe the training dynamics of a large-scale LLM, we use  
277 the setup of OLMo-7B-0424 [18] (henceforth just referred to as OLMo), who have open-sourced  
278 weights at several steps during their training run. For our analysis, we inspect OLMo at a variety of  
279 training steps: every 500 steps throughout the first 10,000 steps, then 25,000 steps, then 50,000 steps,  
280 then every 50,000 steps until 449,000 steps (which is roughly the end of their training). Again, we  
281 use the input “Summer is warm. Winter is cold.”<sup>2</sup> Notice that in this prompt, token 3, namely “.”,  
282 is not very semantically meaningful; it becomes a sink token along with token 0 (c.f. Section 3.1,  
283 Appendix E, Appendix F.2).

284 In Figure 7, we confirm that attention heads go through an attention-increasing and value-state-  
285 shrinking phase, and that the residual state norm of the <bos> token increases linearly during

---

<sup>2</sup>Note that OLMo does not have a <bos> token, but attention sinks still form in the majority of heads. In particular, the first token behaves similarly to an attention sink. We discuss this in Appendix F.2.



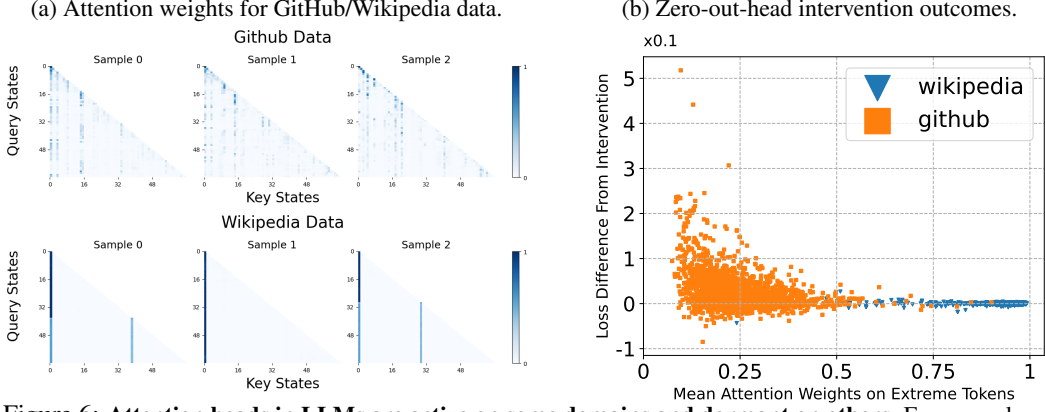


Figure 6: **Attention heads in LLMs are active on some domains and dormant on others.** For example, on Llama 2-7B-Base, we identify that Layer 16 Head 25 is active when the context contains many tokens related to programming, and dormant in other contexts such as prose. We use RedPajama-1T [8] Wikipedia and Github subsets for our data in this figure, truncating all samples to 64 tokens for demonstration purposes. *Left:* Sample weights from four randomly selected samples from each domain. *Right:* Result of an intervention study, i.e., change in cross-entropy of the input sequence when the attention head’s output (concretely, the value states for this head) is manually set to zero, across sequences in both domains. We observe that the model’s performance, measured by cross-entropy, strongly depends on the output of the attention head on coding data.

286 pre-training. We show that, at Layer 24 of OLMo, the average attention on extreme tokens (token  
 287 0 and token 3) increases rapidly at the beginning of training and converges to a constant, while the  
 288 value state norms of extreme tokens decrease rapidly. Also, the residual states of extreme tokens  
 289 also increase linearly, while the rest quickly converge. In Figure 8 we show that attention heads  
 290 converge to a stable phase, and that all logits corresponding to the first token’s value states (i.e., all  
 291 tokens’ value of  $\text{logit}_0$ , except possibly the value of  $\text{logit}_0$  corresponding to token 0 itself) have  
 292 similar distributions. These confirm our dynamics insights from the BB model (c.f. Figure 3).

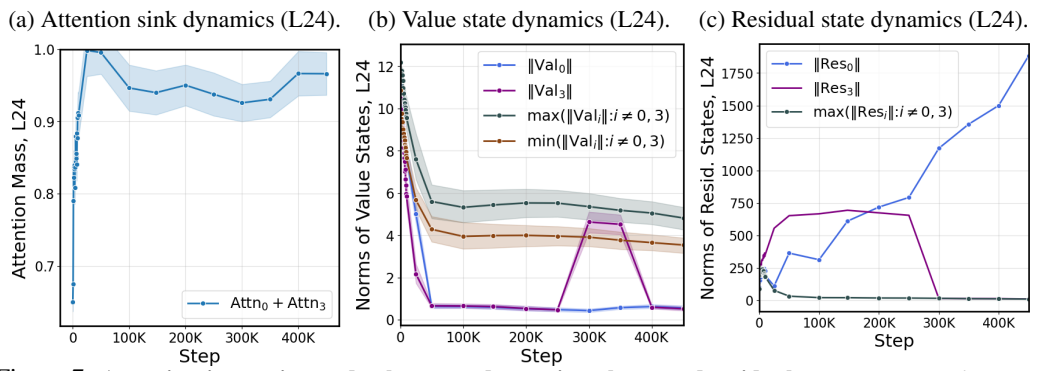


Figure 7: **Attention-increasing and value state-decreasing phase, and residual state norms.** *Left (a):* We plot the total attention mass on extreme tokens 0 and 3 at Layer 24 and averaged over all attention heads, during OLMo training. We observe that it increases rapidly and then maintains its value in  $[0.9, 1]$  for the rest of training, which is in line with our predictions. *Middle (b):* We plot the norm of each token’s value state at Layer 24 during training, averaged over all heads. We observe that the value states of all tokens shrink initially and then converge, while the value states of the extreme tokens shrink to much lower than all other tokens. *Right (c):* We plot the norm of each token’s residual state at Layer 24 during training. We observe that the residual state of token 0 increases linearly in magnitude during training.

293 **4 Conclusion**

294 In this work, we investigated the *extreme-token phenomena*, namely *attention sinks*, *value state drains*,  
 295 and *residual state peaks*. We analyzed a simple evocative model called the Bigram-Backcopy task,  
 296 and theoretically and empirically showed that it exhibited the same extreme-token phenomena as in

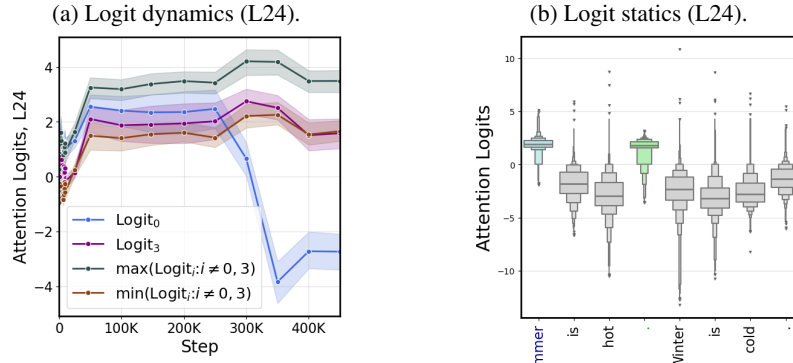


Figure 8: **Stable phase.** *Left (a):* We plot the normalized attention logits of all tokens’ query states against token 0’s key state during training. We observe that the logits of all non-extreme tokens’ query states against token 0’s key state in OLMo’s Layer 24 are stable for a large fraction of the training run, after an initialization period. This echoes the stable phase prediction made in the BB model in Section 2. Note that this prediction makes no guarantees about the logit corresponding to the zeroth query token and zeroth key token, which will be set to 1 by the softmax and so its behavior is irrelevant for prediction. Also note that we use normalization, similar to Section 2, to make all terms comparable; namely we have  $\text{logit}_i = \langle \text{Qry}_i, \text{Key}_0 \rangle - \text{mean}_j(\langle \text{Qry}_i, \text{Key}_j \rangle)$ . *Right (b):* For this experiment, we generate 128 randomly sampled test tokens with IDs from 100 to 50000 in the OLMo tokenizer. We append each token separately to the test phrase “Summer is warm. Winter is cold.”, creating 128 different samples, which we feed to the LLM to record the model behavior. We plot the distribution of (un-normalized) dot products  $\langle \text{Qry}_{\text{test}}, \text{Key}_j \rangle$  across all heads at Layer 24 and all test tokens. We observe that logits of all regular tokens have very similar distributions, and the distributions of the logits corresponding to extreme tokens 0 and 3 are also similar. This confirms the hypothesis that at the end of training, attention heads converge to the stable phase, with similar logits on extreme tokens.

297 LLMs. Based on the Bigram-Backcopy task, we made several detailed predictions about the behavior  
 298 of extreme-token phenomena in LLMs. In particular, we identified the *active-dormant mechanism* for  
 299 attention heads in both the BB model and LLMs, of which attention sinks and value state drains are  
 300 indicators, and a *mutual reinforcement mechanism* by which these phenomena are induced during  
 301 pretraining. Using intuition about these mechanisms, we applied minor interventions to the model  
 302 architecture and optimization procedure which disabled extreme-token phenomena within the BB  
 303 model. Overall, our work uncovers the causes of extreme-token phenomena and points to possible  
 304 pathways to eliminate them during LLM training.

305 We believe the most compelling direction for future work in this area is as follows. Specifically,  
 306 one could build more performant and scalable interventions which would eliminate extreme-token  
 307 phenomena and observe the effect on training dynamics and the finished model. This would make it  
 308 easier to understand whether extreme token phenomena are necessary to build a powerful transformer-  
 309 based LLM, whether they are merely helpful, or whether they are completely incidental to the  
 310 particular architecture and optimization algorithms used by the community.

## 311 References

- 312 [1] Kwangjun Ahn et al. “Linear attention is (maybe) all you need (to understand transformer  
313 optimization)”. In: *arXiv preprint arXiv:2310.01082* (2023).
- 314 [2] Kwangjun Ahn et al. “Transformers learn to implement preconditioned gradient descent for  
315 in-context learning”. In: *Advances in Neural Information Processing Systems* 36 (2024).
- 316 [3] Zeyuan Allen-Zhu and Yuanzhi Li. “Physics of language models: Part 1, context-free grammar”.  
317 In: *arXiv preprint arXiv:2305.13673* (2023).
- 318 [4] Stella Biderman et al. “Pythia: A suite for analyzing large language models across training and  
319 scaling”. In: *International Conference on Machine Learning*. PMLR, 2023, pp. 2397–2430.
- 320 [5] Alberto Bietti et al. “Birth of a transformer: A memory viewpoint”. In: *Advances in Neural  
321 Information Processing Systems* 36 (2024).
- 322 [6] François Charton. “What is my math transformer doing?—Three results on interpretability and  
323 generalization”. In: *arXiv preprint arXiv:2211.00170* (2022).
- 324 [7] Liang Chen et al. “An Image is Worth 1/2 Tokens After Layer 2: Plug-and-Play Inference  
325 Acceleration for Large Vision-Language Models”. In: *arXiv preprint arXiv:2403.06764* (2024).
- 326 [8] Together Computer. *RedPajama: An Open Source Recipe to Reproduce LLaMA training  
327 dataset*. 2023. URL: <https://github.com/togethercomputer/RedPajama-Data>.
- 328 [9] Timothée Darcet et al. “Vision transformers need registers”. In: *arXiv preprint  
329 arXiv:2309.16588* (2023).
- 330 [10] Puneesh Deora et al. “On the optimization and generalization of multi-head attention”. In:  
331 *arXiv preprint arXiv:2310.12680* (2023).
- 332 [11] Tim Dettmers et al. “Gpt3. int8 (): 8-bit matrix multiplication for transformers at scale”. In:  
333 *Advances in Neural Information Processing Systems* 35 (2022), pp. 30318–30332.
- 334 [12] Abhimanyu Dubey et al. “The llama 3 herd of models”. In: *arXiv preprint arXiv:2407.21783*  
335 (2024).
- 336 [13] Nelson Elhage, Robert Lasenby, and Christopher Olah. “Privileged bases in the transformer  
337 residual stream”. In: *Transformer Circuits Thread* (2023).
- 338 [14] Nelson Elhage et al. “A mathematical framework for transformer circuits”. In: *Transformer  
339 Circuits Thread* 1 (2021), p. 1.
- 340 [15] Jiahai Feng and Jacob Steinhardt. “How do language models bind entities in context?” In:  
341 *arXiv preprint arXiv:2310.17191* (2023).
- 342 [16] Yao Fu. “How Do Language Models put Attention Weights over Long Context?” In: *Yao Fu’s  
343 Notion* (2024). URL: <https://yaofu.notion.site/How-Do-Language-Models-put-Attention-Weights-over-Long-Context-10250219d5ce42e8b465087c383a034e?pvs=4>.
- 344 [17] Mor Geva et al. “Dissecting recall of factual associations in auto-regressive language models”.  
345 In: *arXiv preprint arXiv:2304.14767* (2023).
- 346 [18] Dirk Groeneveld et al. “Olmo: Accelerating the science of language models”. In: *arXiv preprint  
347 arXiv:2402.00838* (2024).
- 348 [19] Tianyu Guo et al. “How do transformers learn in-context beyond simple functions? a case  
349 study on learning with representations”. In: *arXiv preprint arXiv:2310.10616* (2023).
- 350 [20] Zhiyu Guo, Hidetaka Kamigaito, and Taro Watanabe. “Attention Score is not All You Need for  
351 Token Importance Indicator in KV Cache Reduction: Value Also Matters”. In: *arXiv preprint  
352 arXiv:2406.12335* (2024).
- 353 [21] Wes Gurnee et al. “Universal neurons in gpt2 language models”. In: *arXiv preprint  
354 arXiv:2401.12181* (2024).
- 355 [22] Chi Han et al. “Lm-infinite: Simple on-the-fly length generalization for large language models”.  
356 In: *arXiv preprint arXiv:2308.16137* (2023).
- 357 [23] Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- 358 [24] Yu Huang, Yuan Cheng, and Yingbin Liang. “In-context convergence of transformers”. In:  
359 *arXiv preprint arXiv:2310.05249* (2023).
- 360 [25] Albert Q Jiang et al. “Mistral 7B”. In: *arXiv preprint arXiv:2310.06825* (2023).
- 361 [26] Juno Kim, Tai Nakamaki, and Taiji Suzuki. “Transformers are Minimax Optimal Nonparametric  
362 In-Context Learners”. In: *arXiv preprint arXiv:2408.12186* (2024).
- 363
- 364

- 365 [27] Ruikang Liu et al. “IntactKV: Improving Large Language Model Quantization by Keeping  
366 Pivot Tokens Intact”. In: *arXiv preprint arXiv:2403.01241* (2024).
- 367 [28] Ziming Liu et al. “Towards understanding grokking: An effective theory of representation  
368 learning”. In: *Advances in Neural Information Processing Systems* 35 (2022), pp. 34651–  
369 34663.
- 370 [29] Kevin Meng et al. “Locating and editing factual associations in GPT”. In: *Advances in Neural  
371 Information Processing Systems* 35 (2022), pp. 17359–17372.
- 372 [30] Neel Nanda et al. “Progress measures for grokking via mechanistic interpretability”. In: *arXiv  
373 preprint arXiv:2301.05217* (2023).
- 374 [31] Eshaan Nichani, Alex Damian, and Jason D Lee. “How transformers learn causal structure  
375 with gradient descent”. In: *arXiv preprint arXiv:2402.14735* (2024).
- 376 [32] Catherine Olsson et al. “In-context learning and induction heads”. In: *arXiv preprint  
377 arXiv:2209.11895* (2022).
- 378 [33] Alec Radford et al. “Language models are unsupervised multitask learners”. In: *OpenAI blog  
379* 1.8 (2019), p. 9.
- 380 [34] Gautam Reddy. “The mechanistic basis of data dependence and abrupt learning in an in-context  
381 classification task”. In: *The Twelfth International Conference on Learning Representations*.  
382 2023.
- 383 [35] Luca Soldaini et al. “Dolma: An open corpus of three trillion tokens for language model  
384 pretraining research”. In: *arXiv preprint arXiv:2402.00159* (2024).
- 385 [36] Seungwoo Son et al. “Prefixing Attention Sinks can Mitigate Activation Outliers for Large  
386 Language Model Quantization”. In: *arXiv preprint arXiv:2406.12016* (2024).
- 387 [37] Mingjie Sun et al. “Massive Activations in Large Language Models”. In: *arXiv preprint  
388 arXiv:2402.17762* (2024).
- 389 [38] Yuandong Tian et al. “Joma: Demystifying multilayer transformers via joint dynamics of mlp  
390 and attention”. In: *arXiv preprint arXiv:2310.00535* (2023).
- 391 [39] Yuandong Tian et al. “Scan and snap: Understanding training dynamics and token composition  
392 in 1-layer transformer”. In: *Advances in Neural Information Processing Systems* 36 (2023),  
393 pp. 71911–71947.
- 394 [40] Eric Todd et al. “Function vectors in large language models”. In: *arXiv preprint  
395 arXiv:2310.15213* (2023).
- 396 [41] Hugo Touvron et al. “Llama 2: Open foundation and fine-tuned chat models”. In: *arXiv preprint  
397 arXiv:2307.09288* (2023).
- 398 [42] Roman Vershynin. *High-dimensional probability: An introduction with applications in data  
399 science*. Vol. 47. Cambridge university press, 2018.
- 400 [43] Kevin Wang et al. “Interpretability in the wild: a circuit for indirect object identification in  
401 gpt-2 small”. In: *arXiv preprint arXiv:2211.00593* (2022).
- 402 [44] Jingfeng Wu et al. “How Many Pretraining Tasks Are Needed for In-Context Learning of  
403 Linear Regression?” In: *arXiv preprint arXiv:2310.08391* (2023).
- 404 [45] Guangxuan Xiao et al. “Efficient Streaming Language Models with Attention Sinks”. In: *arXiv  
405 preprint arXiv:2309.17453* (2023).
- 406 [46] Zhongzhi Yu et al. “Unveiling and Harnessing Hidden Attention Sinks: Enhancing Large  
407 Language Models without Training through Attention Calibration”. In: *arXiv preprint  
408 arXiv:2406.15765* (2024).
- 409 [47] Shuangfei Zhai et al. “Stabilizing transformer training by preventing attention entropy col-  
410 lapse”. In: *International Conference on Machine Learning*. PMLR, 2023, pp. 40770–40803.
- 411 [48] Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. “Trained transformers learn linear models  
412 in-context”. In: *arXiv preprint arXiv:2306.09927* (2023).
- 413 [49] Ruiqi Zhang, Jingfeng Wu, and Peter L Bartlett. “In-context learning of a linear Transformer  
414 block: benefits of the MLP component and one-step GD initialization”. In: *arXiv preprint  
415 arXiv:2402.14951* (2024).
- 416 [50] Yi Zhang et al. “Unveiling transformers with lego: a synthetic reasoning task”. In: *arXiv  
417 preprint arXiv:2206.04301* (2022).
- 418 [51] Zeyuan Allen Zhu and Yuanzhi Li. “Physics of language models: Part 3.1, knowledge storage  
419 and extraction”. In: *arXiv preprint arXiv:2309.14316* (2023).

420 **A Related works**

421 Several studies independently identified the “attention sink” phenomenon in language models and  
 422 vision transformers, where attention weights were found to be concentrated on a few tokens [45,  
 423 9, 22, 47, 13, 11]. Recent research has provided more detailed characterizations of this attention  
 424 pattern and the attention sink phenomenon [16, 37]. Sun et al. [37] attributed the attention sink to  
 425 the massive activation of the hidden representations of the corresponding tokens. Both Sun et al.  
 426 [37] and Zhai et al. [47] discussed methods for mitigating the attention sink by modifying the model  
 427 and training recipes. Additionally, recent studies have leveraged the attention sink phenomenon to  
 428 develop improved quantization and more efficient inference algorithms [27, 7, 46, 36].

429 The dynamics of transformers are studied under various simplifications, including linear attention  
 430 structures [48, 2], reparametrizations [38], NTK [10], often in the setting of in-context linear  
 431 regressions [1, 44, 49] and structured sequence [5, 31, 39]. Notably, Zhang, Frei, and Bartlett [48]  
 432 proves that a one-layer linear attention head trained with gradient descent converges to a model that  
 433 implements the in-context linear regression algorithm. [24, 26] extend this to non-linear settings. [5]  
 434 shows the fast learning of bigram memorization and the slow development of in-context abilities.  
 435 [39] shows the scan and snap dynamics in reparametrized one-layer transformers. [34] simplifies the  
 436 structure of the induction head, showing the connection between the sharp transitions of in-context  
 437 learning dynamics and the nested nonlinearities of multi-layer operations.

438 Mechanistic interpretability is a growing field focused on understanding the internal mechanisms of  
 439 language models in solving specific tasks [14, 17, 29, 30, 32, 5, 43, 15, 40]. This includes mechanisms  
 440 like the induction head and function vector for in-context learning [14, 32, 40, 5], the binding ID  
 441 mechanism for binding tasks [15], association-storage mechanisms for factual identification tasks  
 442 [29], and a complete circuit for indirect object identification tasks [43]. The task addressed in this  
 443 paper is closely related to [5], which explored synthetic tasks where tokens are generated from either  
 444 global or context-specific bigram distributions. Several other studies have also used synthetic tasks to  
 445 investigate neural network mechanisms [6, 28, 30, 3, 51, 19, 50].

446 We note that Gurnee et al. [21] proposed Attention Deactivation Neurons, a concept similar to  
 447 Dormant Attention Heads. Gurnee et al. [21] hypothesized that when such a head attends to the first  
 448 token, it indicates that the head is deactivated and has minimal effect.

449 **B Proofs**

450 Since we drop the trigger tokens in the loss function, we neglect  $\mathcal{T}$  throughout the proof for notational  
 451 convenience, assuming that  $\mathcal{V}$  consists of only non-trigger tokens. We provide new notations which  
 452 are frequently used in the proofs. Define the full bigram transition probability.

$$\mathbf{P} = \begin{pmatrix} p_{11} & \cdots & p_{1V} \\ \vdots & \ddots & \vdots \\ p_{V1} & \cdots & p_{VV} \end{pmatrix} = \begin{pmatrix} \mathbf{p}_1^\top \\ \vdots \\ \mathbf{p}_V^\top \end{pmatrix}. \quad (2)$$

453 Given token  $v$ , define the predicted probability, which is the logit output passed through the softmax  
 454 activation

$$\mathbf{q}_v = \text{SoftMax}(\text{TF}([\langle \text{bos} \rangle; v_{1:n-1}; v]_n)). \quad (3)$$

455 Similarly, define the full output probability matrix.

$$\mathbf{Q} = \begin{pmatrix} q_{11} & \cdots & q_{1V} \\ \vdots & \ddots & \vdots \\ q_{V1} & \cdots & q_{VV} \end{pmatrix} = \begin{pmatrix} \mathbf{q}_1^\top \\ \vdots \\ \mathbf{q}_V^\top \end{pmatrix}. \quad (4)$$

456 Given any vector  $\mathbf{u} = [u_1; \dots; u_d]$ , define the corresponding diagonal matrix as

$$\text{diag}(\mathbf{u}) = \begin{pmatrix} u_1 & 0 & \cdots & 0 \\ \vdots & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & u_d \end{pmatrix}.$$

457 Define

$$\mathbf{G}_v^{\mathbf{Q}} = \text{diag}(\mathbf{q}_v) - \mathbf{q}_v \mathbf{q}_v^\top \quad \mathbf{G}_v^{\mathbf{P}} = \text{diag}(\mathbf{p}_v) - \mathbf{p}_v \mathbf{p}_v^\top.$$

458 Denote  $\mathbf{z} = \mathbf{W} \cdot \boldsymbol{\beta} - \mathbf{W} \circ \boldsymbol{\xi}$ . We present a technical lemma.

459 **Lemma 5.** *The matrices  $\mathbf{G}_v^{\mathbf{P}}$  and  $\mathbf{G}_v^{\mathbf{Q}}$  are positive semi-definite for any  $v$ .*

460 *Proof.* Since we have that  $\sum_{k=1}^V p_{vk} = 1$  and  $\sum_{k=1}^V q_{vk} = 1$  for any  $v$ ,

$$\begin{aligned} (\mathbf{G}_v^{\mathbf{P}})_{ii} &= p_i - p_i^2 = p_i \left( \sum_{k \neq i} p_k \right) \geq \sum_{k \neq i} |(\mathbf{G}_v^{\mathbf{P}})_{ik}| \\ (\mathbf{G}_v^{\mathbf{Q}})_{ii} &= q_i - q_i^2 = q_i \left( \sum_{k \neq i} q_k \right) \geq \sum_{k \neq i} |(\mathbf{G}_v^{\mathbf{Q}})_{ik}|. \end{aligned}$$

461 This shows that both  $\mathbf{G}_v^{\mathbf{P}}$  and  $\mathbf{G}_v^{\mathbf{Q}}$  are diagonally dominant matrices. By Corollary 6.2.27 in Horn  
462 and Johnson [23], they are positive semi-definite.  $\square$

### 463 B.1 Proof of Theorem 2

464 We denote the hidden dimension as  $d$  and the sequence length as  $N$ . We begin with the assumption  
465 regarding the transformer’s positional embedding:

466 **Assumption A.** *For any token  $v$  and position  $i$ , assume that the encoding combined with the positional  
467 embedding ensures that  $\{\text{ebd}(v_i)\}$  is linearly independent.*

468 Assumption A requires that  $d \geq VN$ . Given the fact that there are  $O(\exp(d))$  approximately linearly  
469 independent vectors for large  $d$  [42], it is possible to apply approximation theory to avoid Assumption  
470 A. However, since Assumption A pertains only to the construction of  $\lambda$  for trigger tokens and is  
471 unrelated to Theorem 3, we adopt it to simplify the proof of Theorem 2.

472 *Proof.* Consider vectors  $\mathbf{u}_i \in \mathbb{R}^d$ ,  $i \in [N]$  such that  $\mathbf{u}_i^\top \mathbf{u}_j = 0$ ,  $i \neq j$ , and  $\mathbf{u}_i^\top \text{ebd}(v_j)$  for any  
473  $v \in \mathcal{V}$  and  $i, j \in [N]$ . Adopting Assumption A, there exists a matrix  $\text{Qry}$  such that

$$\begin{aligned} \text{Qry}(\text{ebd}(v_i)) &= \lambda \mathbf{u}_{i-1} \quad \text{for } v_i \in \mathcal{T}, \quad i > 1, \\ \text{Qry}(\text{ebd}(v_i)) &= \alpha_{v_i} \mathbf{u}_0 \quad \text{for } v_i \in \mathcal{V} \setminus \mathcal{T}, \quad i > 0. \end{aligned} \tag{5}$$

474 Define the corresponding key matrix.

$$\begin{aligned} \text{Key}(\text{ebd}(v_i)) &= \mathbf{u}_i \quad \text{for } v_i \in \mathcal{V}, \quad i > 0, \\ \text{Key}(\text{ebd}(\langle \text{bos} \rangle)) &= \mathbf{u}_0. \end{aligned} \tag{6}$$

475 There exists a value matrix  $\text{Val}$  such that

$$\begin{aligned} \text{Val}(\text{ebd}(v_i)) &= 0 \quad \text{for } v_i \in \mathcal{T}, \quad i > 1, \\ \text{Val}(\text{ebd}(v_i)) &= \xi_{v_i} \mathbf{u}_i \quad \text{for } v_i \in \mathcal{V} \setminus \mathcal{T}, \quad i > 0, \\ \text{Val}(\text{ebd}(\langle \text{bos} \rangle)) &= \boldsymbol{\beta}. \end{aligned} \tag{7}$$

476 Further define the matrix  $\mathbf{M}$  that satisfies

$$\begin{aligned} \mathbf{M}(\text{ebd}(v_i)) &= \log \mathbf{p}_{v_i} \cdot \mathbf{1}\{v_i \notin \mathcal{T}\} \quad \text{for } v_i \in \mathcal{V}, \quad i \in [N], \\ \mathbf{M}(\mathbf{u}_i) &= \mathbf{e}_i \quad \text{for } i \in [N]. \end{aligned} \tag{8}$$

477 Setting  $\text{mlp}(\cdot) = \text{ReLU}(\mathbf{M}(\cdot))$ , we can then verify that the residual connection gives that  
478  $\text{TF}(\langle \text{bos} \rangle; v_{1:n-1}; v_n) = \text{mlp}(\text{ebd}(v_n) + \text{attn}(\text{ebd}(v_n)))$ , which is equivalent to the simplified  
479 model.

480 When  $\min_{v \in \mathcal{V}} \alpha_v \rightarrow \infty$ ,  $\min_{v \in \mathcal{V}} \xi_v \rightarrow \infty$ ,  $\lambda \rightarrow \infty$ , and  $\boldsymbol{\beta} = \mathbf{0}$ , if  $v_n \in \mathcal{T}$ ,  
481  $\text{SoftMax}[\text{TF}(\langle \text{bos} \rangle; v_{1:n-1}; v_n)] = \delta_{v_{n-1}}$ . If  $v_n \in \mathcal{V} \setminus \mathcal{T}$ ,  $\text{SoftMax}[\text{TF}(\langle \text{bos} \rangle; v_{1:n-1}; v_n)] =$   
482  $\mathbf{p}_{v_n}$ . All next-token probabilities match those in the data-generating procedure, aligning with the  
483 oracle algorithm.  $\square$

484 **B.2 The stable phase in Theorem 3**

485 Lemma 6 computes the gradient of  $\mathbf{Q}$ .

486 **Lemma 6.** *We have*

$$\begin{aligned}\frac{\partial q_{ik}}{\partial \alpha_v} &= \frac{\mathbf{1}\{i=v\}q_{ik}e^{\alpha_i}}{(e^{\alpha_i}+W)^2} \left[ W\beta_k - W_k\xi_k - \sum_{j=1}^V q_{ij}(W\beta_j - W_j\xi_j) \right], \\ \frac{\partial q_{ik}}{\partial \beta_v} &= \frac{e^{\alpha_i}}{e^{\alpha_i}+W} [q_{ik}\mathbf{1}\{k=v\} - q_{ik}q_{iv}].\end{aligned}$$

487 *Furthermore,*

$$\sum_{v=1}^V \frac{\partial q_{ik}}{\partial \alpha_v} = 0, \quad \sum_{v=1}^V \frac{\partial q_{ik}}{\partial \beta_v} = 0.$$

488 *Proof.* We repeatedly use the following two facts:

$$\begin{aligned}\frac{\partial \left\{ \exp \left[ \frac{W_k\xi_k + e^{\alpha_i}\beta_k}{e^{\alpha_i}+W} \right] \right\}}{\partial \alpha_v} &= \frac{e^{\alpha_v}(W\alpha_k - W_k\xi_k)}{(e^{\alpha_i}+W)^2} \exp \left[ \frac{W_k\xi_k + e^{\alpha_i}\beta_k}{e^{\alpha_i}+W} \right], \\ \frac{\partial \left\{ \exp \left[ \frac{W_k\xi_k + e^{\alpha_i}\beta_k}{e^{\alpha_i}+W} \right] \right\}}{\partial \beta_v} &= \frac{\mathbf{1}\{i=v\}e^{\alpha_i}}{e^{\alpha_i}+W} \exp \left[ \frac{W_k\xi_k + e^{\alpha_i}\beta_k}{e^{\alpha_i}+W} \right].\end{aligned}$$

489 When  $i \neq v$ ,  $q_{ik}$  does not include  $\alpha_v$ , making the gradients as zero. When  $i = v$ , we have

$$\begin{aligned}\frac{\partial q_{vk}}{\partial \alpha_v} &= q_{vk}e^{\alpha_v} \left[ \frac{W\beta_k - W_k\xi_k}{(e^{\alpha_v}+W)^2} \right] - \frac{q_{vk} \sum_{i=1}^V p_{vi}e^{\alpha_v} \left[ \frac{W\beta_i - W_i\xi_i}{(e^{\alpha_v}+W)^2} \right] \exp \left[ \frac{W_i\xi_i + e^{\alpha_v}\beta_i}{e^{\alpha_v}+W} \right]}{\sum_{i=1}^V p_{vi} \exp \left[ \frac{W_i\xi_i + e^{\alpha_v}\beta_i}{e^{\alpha_v}+W} \right]} \\ &= \frac{e^{\alpha_v}}{(e^{\alpha_v}+W)^2} \left\{ q_{vk}[W\beta_k - W_k\xi_k] - q_{vk} \sum_{j=1}^V q_{vj}^\top (W\alpha_j - W_j\xi_j) \right\},\end{aligned}$$

490 and

$$\begin{aligned}\frac{\partial q_{ik}}{\partial \beta_v} &= \left[ \frac{e^{\alpha_i}}{e^{\alpha_i}+W} \right] q_{ik}\mathbf{1}\{k=v\} - \frac{\left[ \frac{e^{\alpha_i}}{e^{\alpha_i}+W} \right] p_{iv} \exp \left[ \frac{W_v\xi_v + e^{\alpha_i}\beta_v}{e^{\alpha_i}+W} \right] p_{iv} \exp \left[ \frac{W_k\xi_k + e^{\alpha_i}\beta_k}{e^{\alpha_i}+W} \right]}{\left( \sum_{j=1}^V p_{vj} \exp \left[ \frac{W_j\xi_j + e^{\alpha_i}\beta_j}{e^{\alpha_i}+W} \right] \right)^2} \\ &= \left[ \frac{e^{\alpha_i}}{e^{\alpha_i}+W} \right] [q_{ik}\mathbf{1}\{k=v\} - q_{ik}q_{iv}].\end{aligned}$$

491 We can verify that

$$\begin{aligned}\sum_{v=1}^V \frac{\partial q_{ik}}{\partial \alpha_v} &= \frac{e^{\alpha_v}}{(e^{\alpha_v}+W)^2} \sum_{v=1}^V \left\{ q_{vk}[W\beta_k - W_k\xi_k] - q_{vk} \sum_{j=1}^V q_{vj}^\top (W\alpha_j - W_j\xi_j) \right\} \\ &= \frac{e^{\alpha_v}}{(e^{\alpha_v}+W)^2} \left\{ \sum_{v=1}^V q_{vk}[W\beta_k - W_k\xi_k] - \sum_{j=1}^V q_{vj}^\top (W\alpha_j - W_j\xi_j) \right\} \\ &= 0,\end{aligned}$$

492 and

$$\begin{aligned}\sum_{v=1}^V \frac{\partial q_{ik}}{\partial \beta_v} &= \left[ \frac{e^{\alpha_i}}{e^{\alpha_i}+W} \right] \sum_{v=1}^V [q_{ik}\mathbf{1}\{k=v\} - q_{ik}q_{iv}] \\ &= \left[ \frac{e^{\alpha_i}}{e^{\alpha_i}+W} \right] [q_{iv} - q_{iv}] \\ &= 0.\end{aligned}$$

493 This finishes the proof of Lemma 6. □

494 Proposition 7 computes the gradient of loss with respect to  $\alpha$  and  $\beta$ , giving the gradient flow.

495 **Proposition 7.** *The gradient flow of optimizing  $\text{loss}(\alpha, \beta)$  is given by*

$$\begin{aligned}\dot{\alpha}_v(t) &= \frac{\pi_v e^{\alpha_v}}{(e^{\alpha_v} + W)^2} \sum_{i=1}^V (p_{vi} - q_{vi})(W\beta_i - W_i\xi_i), \\ \dot{\beta}_v(t) &= \sum_{k=1}^V \left\{ \frac{\pi_k e^{\alpha_k} [p_{kv} - q_{kv}]}{e^{\alpha_k} + W} \right\}.\end{aligned}$$

496 *Proof.* The gradient flow gives that

$$\dot{\alpha}_v(t) = -\frac{\partial \text{loss}(\alpha, \beta)}{\partial \alpha_v}, \quad \text{and} \quad \dot{\beta}_v(t) = -\frac{\partial \text{loss}(\alpha, \beta)}{\partial \beta_v}.$$

497 Taking the derivative of  $\text{loss}(\alpha, \beta)$  gives that

$$\begin{aligned}\frac{\partial \text{loss}(\alpha, \beta)}{\partial \alpha_v} &= \pi_v \sum_{k=1}^V p_{vk} \cdot \frac{-1}{q_{vi}} \cdot \frac{\partial q_{vi}}{\partial \alpha_v} \\ &= \frac{\pi_v e^{\alpha_v}}{(e^{\alpha_v} + W)^2} \left\{ \sum_{i=1}^V q_{vi} [W\beta_i - W_i\xi_i] - \sum_{k=1}^V p_{vk} [W\beta_k - W_k\xi_k] \right\} \\ &= \frac{\pi_v e^{\alpha_v}}{(e^{\alpha_v} + W)^2} \sum_{k=1}^V \left\{ [q_{vk} - p_{vk}] [W\beta_k - W_k\xi_k] \right\}.\end{aligned}$$

498 Similarly, we have that

$$\begin{aligned}\frac{\partial \text{loss}(\alpha, \beta)}{\partial \beta_v} &= \sum_{j=1}^V \pi_j \sum_{k=1}^V p_{jk} \left\{ \frac{e^{\alpha_j} q_{jv}}{e^{\alpha_j} + W} - \frac{e^{\alpha_j} \mathbf{1}\{k=v\}}{e^{\alpha_j} + W} \right\} \\ &= \sum_{j=1}^V \left\{ \frac{\pi_j e^{\alpha_j} [q_{jv} - p_{jv}]}{e^{\alpha_j} + W} \right\}.\end{aligned}$$

499 This proves Proposition 7. □

500 **Theorem 8** (Restatement the stable phase part in Theorem 3). *Consider the gradient flow of optimizing*  
501  *$\text{loss}(\alpha, \beta)$ . The gradient flow has sink stationary points*

$$\alpha^* = \alpha \mathbf{1}, \quad \beta^* = c \cdot \mathbf{1} - e^{-\alpha} \cdot \mathbf{W} \circ \xi.$$

502 *Proof.* When  $\alpha = \alpha^*$  and  $\beta = \beta^*$ ,

$$\begin{aligned}q_{vi} &= \frac{p_{vi} \exp \left[ \frac{W_i \xi_i + e^{\alpha} \beta_i}{e^{\alpha} + W} \right]}{\sum_{k=1}^V p_{vk} \exp \left[ \frac{W_k \xi_k + e^{\alpha} \beta_k}{e^{\alpha} + W} \right]} \\ &= \frac{p_{vi} \exp \left[ \frac{c}{e^{\alpha} + W} \right]}{\sum_{k=1}^V p_{vk} \exp \left[ \frac{c}{e^{\alpha} + W} \right]} \\ &= p_{vi}.\end{aligned}$$

503 Take  $q_{vi}$ 's into  $\partial \text{loss}(\alpha, \beta) / \partial \alpha$  and  $\partial \text{loss}(\alpha, \beta) / \partial \beta$ .

$$\begin{aligned}\frac{\partial \text{loss}(\alpha, \beta)}{\partial \alpha_v} \Big|_{\alpha^*, \beta^*} &= \frac{\pi_v e^{\alpha_v}}{(e^{\alpha_v} + W)^2} \sum_{k=1}^V \left\{ (q_{vk} - p_{vk}) [W\beta_k - W_k\xi_k] \right\} = 0, \\ \frac{\partial \text{loss}(\alpha, \beta)}{\partial \beta_v} \Big|_{\alpha^*, \beta^*} &= \sum_{k=1}^V \left\{ \frac{\pi_k e^{\alpha_k} [q_{kv} - p_{kv}]}{e^{\alpha_k} + W} \right\} = 0.\end{aligned}$$



504 This shows that the given points are stationary points. We further compute the second-order derivative  
 505 using Lemma 6.

$$\begin{aligned}
 \frac{\partial^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \alpha_i \partial \alpha_v} \Big|_{\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*} &= \mathbf{1}\{v = i\} \cdot \frac{\pi_v e^\alpha}{(e^\alpha + W)^2} \sum_{k=1}^V \left\{ \frac{\partial q_{ik}}{\partial \alpha_v} [W \beta_k - W_k \xi_k] \right\} \\
 &= \mathbf{1}\{v = i\} \cdot \frac{-\pi_v e^{2\alpha}}{(e^\alpha + W)^4} \left\{ \sum_{k=1}^V q_{ik} (e^{-\alpha} W + W_k)^2 \xi_k^2 - \left[ \sum_{k=1}^V q_{ik} (e^{-\alpha} W + W_k) \xi_k \right]^2 \right\}, \\
 &= \mathbf{1}\{v = i\} \cdot \frac{-\pi_v e^{2\alpha}}{(e^\alpha + W)^4} \left\{ \sum_{k=1}^V p_{ik} (e^{-\alpha} W + W_k)^2 \xi_k^2 - \left[ \sum_{k=1}^V p_{ik} (e^{-\alpha} W + W_k) \xi_k \right]^2 \right\}.
 \end{aligned}$$

506 where in the second line, we take  $\beta_k^* = c - e^{-\alpha} \xi_k$  and use that  $\sum_{k=1}^V \partial q_{ik} / \partial \alpha_v = 0$ . In the last line,  
 507 we take  $\mathbf{Q} = \mathbf{P}$ . Similarly, we compute the gradients with respect to  $\alpha_i$  and  $\beta_v$ .

$$\begin{aligned}
 \frac{\partial^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \alpha_i \partial \beta_v} \Big|_{\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*} &= \frac{\pi_i e^\alpha}{(e^\alpha + W)^2} \sum_{k=1}^V \left\{ \frac{\partial q_{ik}}{\partial \beta_v} [W \beta_k - W_k \xi_k] \right\} \\
 &= \frac{p_{iv} \pi_i e^{2\alpha}}{(e^\alpha + W)^3} \left\{ - (e^{-\alpha} W + W_k) \xi_k + \sum_{k=1}^V p_{ik} (e^{-\alpha} W + W_k) \xi_k \right\}.
 \end{aligned}$$

508 With the same manner, we compute the gradients with respect to  $\beta_i$  and  $\beta_v$ .

$$\begin{aligned}
 \frac{\partial^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \beta_i \partial \beta_v} \Big|_{\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*} &= \sum_{k=1}^V \left\{ \frac{\partial q_{ki}}{\partial \beta_v} \frac{\pi_k e^\alpha}{e^\alpha + W} \right\} \\
 &= \frac{e^{2\alpha}}{(e^\alpha + W)^2} \sum_{k=1}^V [\mathbf{1}\{v = i\} p_{kv} - p_{ki} p_{kv}].
 \end{aligned}$$

509 Define  $\mathbf{z} = [z_1; \dots; z_V]$  so that  $z_k = -(e^{-\alpha} W + W_k) \xi_k$ . Combining above computations gives that

$$\text{Hessian}(\text{loss}(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)) = \begin{pmatrix} \nabla_{\boldsymbol{\alpha}}^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) & \nabla_{\boldsymbol{\alpha}} \nabla_{\boldsymbol{\beta}} \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) \\ \nabla_{\boldsymbol{\beta}} \nabla_{\boldsymbol{\alpha}} \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) & \nabla_{\boldsymbol{\beta}}^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) \end{pmatrix},$$

510 with

$$\begin{aligned}
 \nabla_{\boldsymbol{\alpha}}^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{e^{2\alpha}}{(e^\alpha + W)^4} \text{diag} \left\{ \pi \circ [\mathbf{z}^\top \mathbf{G}_1^{\mathbf{P}} \mathbf{z}; \dots; \mathbf{G}_V^{\mathbf{P}} \mathbf{z}] \right\}, \\
 \nabla_{\boldsymbol{\alpha}} \nabla_{\boldsymbol{\beta}} \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{e^{2\alpha}}{(e^\alpha + W)^3} \text{diag} \left\{ \pi \right\} [\mathbf{z}^\top \mathbf{G}_1^{\mathbf{P}}; \dots; \mathbf{z}^\top \mathbf{G}_V^{\mathbf{P}}], \\
 \nabla_{\boldsymbol{\beta}}^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) &= \frac{e^{2\alpha}}{(e^\alpha + W)^2} \sum_{k=1}^V \pi_k \mathbf{G}_k^{\mathbf{P}}.
 \end{aligned}$$

511 At last, we diagonalize the Hessian matrix and get that

$$\text{Diag-Hessian}(\text{loss}(\boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)) = \begin{pmatrix} \nabla_{\boldsymbol{\alpha}}^2 \text{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) & 0 \\ 0 & \frac{e^{2\alpha}}{(e^\alpha + W)^2} \mathbf{H} \end{pmatrix},$$

512 where the  $\mathbf{H}$  is given by

$$\mathbf{H} = \sum_{k=1}^V \pi_k \left( \mathbf{G}_k^{\mathbf{P}} - (\mathbf{z}^\top \mathbf{G}_k^{\mathbf{P}} \mathbf{z})^{-1} \mathbf{G}_k^{\mathbf{P}} \mathbf{z} \mathbf{z}^\top \mathbf{G}_k^{\mathbf{P}} \right).$$

513 To prove that  $\mathbf{H}$  is positive semi-definite, consider any vector  $\boldsymbol{\eta}$  with  $\|\boldsymbol{\eta}\|_2 = 1$ .

$$\boldsymbol{\eta}^\top \mathbf{H} \boldsymbol{\eta} = \sum_{k=1}^V \pi_k \left( \boldsymbol{\eta}^\top \mathbf{G}_k^{\mathbf{P}} \boldsymbol{\eta} - \frac{\boldsymbol{\eta}^\top \mathbf{G}_k^{\mathbf{P}} \mathbf{z} \mathbf{z}^\top \mathbf{G}_k^{\mathbf{P}} \boldsymbol{\eta}}{\mathbf{z}^\top \mathbf{G}_k^{\mathbf{P}} \mathbf{z}} \right).$$

514 Since  $\mathbf{G}_k^{\mathbf{P}}$ 's are positive semi-definite, the Cauchy inequality gives that

$$z^\top \mathbf{G}_k^{\mathbf{P}} \eta \leq \sqrt{z^\top \mathbf{G}_k^{\mathbf{P}} z \eta^\top \mathbf{G}_k^{\mathbf{P}} \eta}.$$

515 As a result, we have that

$$\eta^\top \mathbf{H} \eta \geq \sum_{k=1}^V \pi_k \left( \eta^\top \mathbf{G}_k^{\mathbf{P}} \eta - \frac{z^\top \mathbf{G}_k^{\mathbf{P}} z \eta^\top \mathbf{G}_k^{\mathbf{P}} \eta}{z^\top \mathbf{G}_k^{\mathbf{P}} z} \right) = 0.$$

516 This shows that  $\mathbf{H}$  is positive semi-definite. Therefore,  $\text{Hessian}(\text{loss}(\alpha^*, \beta^*))$  is positive semi-definite.  
517 This proves Theorem 8.  $\square$

518 We prove Theorem 8 through direct computation. Due to the non-linearity, it's unclear whether  
519 other stationary points exist. However, we observe that all of our simulations converge to the given  
520 stationary points.

### 521 B.3 Attention sinks in Theorem 3

522 **Theorem 9** (Restatement of the attention sink part in Theorem 3). *Fixing  $\beta = c \cdot \mathbf{1}$ , with any initial*  
523 *value, there exists  $\mathbf{r}(t)$  with bounded norm such that*

$$\alpha(t) = \frac{1}{2} \log t \cdot \mathbf{1} + \mathbf{r}(t).$$

524 *Proof.* We separately analyze each entry of  $\alpha$ . Focusing on  $\alpha_v$ , to simplify the notation, we introduce  
525 a random variable  $\varphi$  such that  $\mathbb{P}(\varphi = W_k \xi_k) = p_{vk}$ . Define

$$u = e^{\alpha_v}.$$

526 Therefore, using Lemma 7, we get that

$$\frac{du}{dt} = \frac{\pi_v e^{2\alpha_v}}{(e^{\alpha_v} + W)^2} \sum_{i=1}^V (q_{vi} - p_{vi})(W \beta_i - W_i \xi_i).$$

527 We take in  $\beta = c$  and expand the expression of  $du/dt$ . This gives us

$$\begin{aligned} \frac{du}{dt} &= \frac{\pi_v u^2}{(u + W)^2} \frac{\sum_{k=1}^V p_{vk} e^{W_k \xi_k / (u+W)} W_k \xi_k - \sum_{k=1}^V p_{vk} e^{W_k \xi_k / (u+W)} \sum_{k=1}^V W_k \xi_k}{\sum_{k=1}^V p_{vk} e^{W_k \xi_k / (u+W)}} \\ &= \frac{\pi_v u^2}{(u + W)^2} \frac{\text{Cov}(e^{\frac{\varphi}{u+W}}, \varphi)}{\mathbb{E} e^{\frac{\varphi}{u+W}}}. \end{aligned}$$

528 Since both  $e^{x/(u+W)}$  and  $x$  are monotonically increasing with respect to  $x$ ,  $u$  is monotonically  
529 increasing. This means that

$$\frac{u(t)^2}{[u(t) + W]^2} \geq \frac{u(0)^2}{[u(0) + W]^2}, \quad \mathbb{E} e^{\frac{\varphi}{u(t)+W}} \leq \mathbb{E} e^{\frac{\varphi}{u(0)+W}}.$$

530 Meanwhile, if we consider the first and second order approximation of  $e^{\varphi/(u+W)}$ ,

$$e^{\frac{\varphi}{u+W}} = 1 + \frac{\theta_1(\varphi)\varphi}{u+W}, \quad e^{\frac{\varphi}{u+W}} = 1 + \frac{\varphi}{u+W} + \theta_2(\varphi) \left[ \frac{\varphi}{u+W} \right]^2.$$

531 Both  $\theta_1(\varphi)$  and  $\theta_2(\varphi)$  are monotonically increasing functions of  $\varphi$ . We also have the bound

$$\theta(\varphi) \leq \frac{e^{\frac{\max \varphi}{u(0)+W}} - 1}{\frac{\max \varphi}{u(0)+W} - 1} = C_\theta.$$

532 Therefore, we get two more inequalities

$$\text{Cov}(\theta_1(\varphi)\varphi, \varphi) \leq C_\theta \text{Var}(\varphi), \quad \text{Cov}(\theta_2(\varphi)\varphi^2, \varphi) \geq 0.$$

533 With all the preparatory works down, we give upper and lower bounds for  $du/dt$ . We first upper-  
 534 bound  $du/dt$ .

$$\begin{aligned}\frac{du}{dt} &\leq \pi_v \text{Cov}(e^{\frac{\varphi}{u+W}}, \varphi) \\ &= \pi_v \text{Cov}\left(1 + \frac{\theta_1(\varphi)\varphi}{u+W}, \varphi\right) \\ &\leq \frac{\pi_v C_\theta \text{Var}(\varphi)}{u}.\end{aligned}$$

535 By solving the corresponding ODE, we get that

$$\frac{1}{2}u^2 \leq \sqrt{C_\theta \text{Var}(\varphi)t} + C.$$

536 To give a lower bound, we have that

$$\begin{aligned}\frac{du}{dt} &\geq \frac{u(0)^2}{[u(0)+W]^2} \frac{\pi_v \text{Cov}(e^{\frac{\varphi}{u+W}}, \varphi)}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \\ &\geq \frac{u(0)^2}{[u(0)+W]^2} \frac{\pi_v}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \text{Cov}\left(1 + \frac{\varphi}{u+W} + \theta_2(\varphi)\left[\frac{\varphi}{u+W}\right]^2, \varphi\right) \\ &\geq \frac{u(0)^2}{[u(0)+W]^2} \frac{\pi_v}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \frac{\text{Var}(\varphi)}{u+W} \\ &\geq \frac{u(0)^2}{[u(0)+W]^2} \frac{\pi_v}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \cdot \frac{u(0)}{u(0)+W} \cdot \frac{\text{Var}(\varphi)}{u} \\ &= \tilde{C}_\theta \frac{1}{u}.\end{aligned}$$

537 Therefore,  $u \geq \sqrt{\tilde{C}_\theta t} + \tilde{C}$ . In conclusion,

$$y_v = \log u = \frac{1}{2} \log t + r_v,$$

538 with  $r_v$  bounded. □

#### 539 **B.4 Value state drains in Theorem 3**

540 **Theorem 10** (Restatement of Theorem 3). *Fixing  $\alpha = y\mathbf{1}$ ,  $\beta = c\mathbf{1} - e^{-\alpha}\mathbf{W} \circ \xi$  with  $c \in \mathbb{R}$ . Define*  
 541  $\bar{\beta}(t) = V^{-1} \sum_{i=1}^V \beta_i(t)$ . *Then the gradient flow of  $\beta(t)$  converges:*

$$\beta(t) \rightarrow \beta^* = \bar{\beta}(0)\mathbf{1} - e^{-\alpha}\mathbf{W} \circ \xi.$$

542 *Proof.* Theorem 8 has already verified that  $\beta = c\mathbf{1} - e^{-\alpha}\mathbf{W} \circ \xi$  are stationary points of loss. In the  
 543 proof of Theorem 8, we have derived  $\nabla_{\beta}^2 \text{loss}(\alpha, \beta)$ .

$$\nabla_{\beta}^2 \text{loss}(\alpha, \beta) = \sum_{k=1}^V \pi_k \mathbf{G}_k^{\mathbf{Q}}.$$

544 Lemma 5 indicates that it is positive semi-definite. Therefore, all stationary points attain the minimum  
 545 of  $\text{loss}(\alpha, \beta)$ . Suppose  $\beta^*$  is a stationary point, we therefore get that  $q_{vk} = p_{vk}$  for any  $v, k$ . This  
 546 implies that  $e^y \beta_k^* + W_k \xi_k$  are constants across  $k$ . We can solve  $\beta^*$  and get that  $\beta^* = c\mathbf{1} - e^{-\alpha}\mathbf{W} \circ \xi$ .  
 547 The convexity of the  $\text{loss}(\alpha, \beta)$  guarantees that  $\beta$  always converges to a stationary point  $\beta^*$ .

548 To find the value of  $c$  in  $\beta^*$ , note that  $\sum_{v=1}^V \dot{\beta}_v(t) = 0$ . We get that  $\bar{\beta}^* = \bar{\beta}(0)$ . Therefore,  
 549  $\beta^* = \beta^* = \bar{\beta}(0)\mathbf{1} - e^{-\alpha}\mathbf{W} \circ \xi$ . □

550 **Remark 11.** *If we assume that  $p_{vk} > 0$  for any  $v, k$  and suppose that the initial value  $\beta(0)$  is close*  
 551 *enough to  $\beta^*$ , it is possible to prove the fast convergence of  $\beta(t)$  to  $\beta^*$ .*

$$\|\beta(t) - \beta^*\|_2^2 \leq \delta e^{-\mu t}.$$

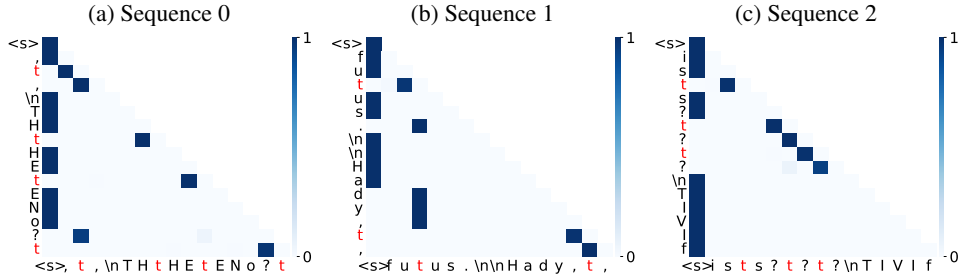


Figure 9: Attention plots of the one-layer transformer trained on the Bigram-Backcopy task.

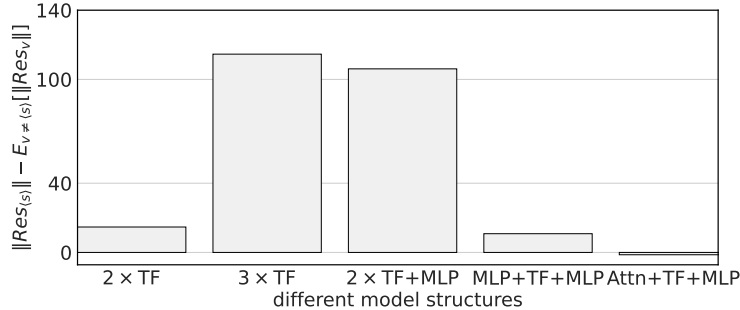


Figure 10: Minimal structures to elicit residual state peaks. We use  $A + B + C$  to indicate the model with structure  $A$ ,  $B$ ,  $C$  in layers 0, 1, and 2, respectively.

## 552 C Ablations

553 **Experimental details.** We train transformers with positional embedding, pre-layer norm, SoftMax  
 554 activation in `attn`, and ReLU activation in `mLp`. We use Adam with constant learning rate 0.0003,  
 555  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ ,  $\varepsilon = 10^{-8}$ , and a weight decay of 0.01. We choose a learning rate of 0.03 for  
 556 the SGD. In each training step, we resample from the BB task with a batch size of  $B = 512$  and  
 557 sequence length  $N = 256$ . Unless otherwise specified, the model is trained for 10,000 steps. Results  
 558 are consistent across different random seeds.

559 **More attention plots** : Figure 9 presents more attention-weight heat maps of the one-layer trans-  
 560 former model trained on the BB task. All attention maps show the attention sink phenomenon.  
 561 Interestingly, the trigger tokens serve as attention sinks in some inputs.

### 562 C.1 Ablations of different model structures trained on the Bigram-Backcopy task.

563 **Exploring the minimal structure for massive norms.** Figure 10 presents the difference of residual  
 564 norms between the `<bos>` token and others ( $\|Res_{<bos>}\| - \mathbb{E}_{v \neq <bos>}[\|Res_v\|]$ ), with different combi-  
 565 nations of model structures. The  $3 \times TF$  and  $2 \times TF + mLp$  are two outliers, showing clear evidence  
 566 of residual state peaks.

567 **Attention plots, value state norms, and residual norms for a three-layer transformer trained on**  
 568 **BB task.** Figures 11, 12, and 13 show the extreme token phenomena in a three-layer transformer.  
 569 The residual state peaks show different phenomena from those in LLMs, with the last layer output  
 570 increasing the residual norms of non-`<bos>` tokens. Figure 1 demonstrates that the residual state  
 571 norms of `<bos>` drop match the magnitudes of other tokens at the last layer.

572 **Statics and dynamics of the simplified model in Theorem 3.** With the simplified model structure  
 573 in Figure 4, we pre-train the model using Adam with learning rate 0.03. Figure 14 and 15 show  
 574 results that match both the theory and the observations of the one-layer transformer.

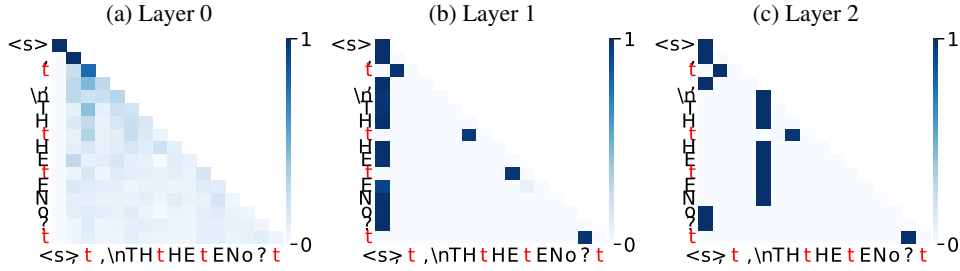


Figure 11: Value state norms of three-layer transformer trained on the BB task

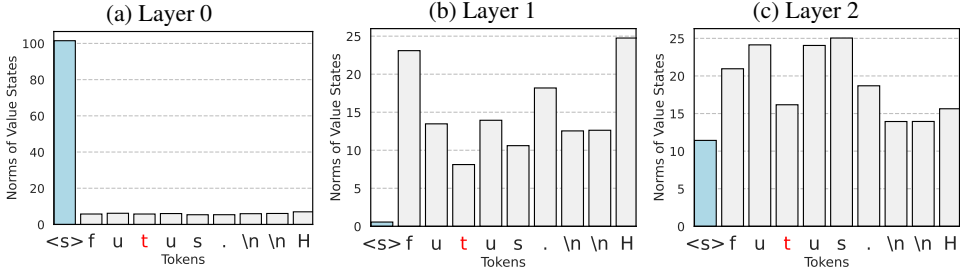


Figure 12: Value state norms of three-layer transformer trained on the BB task

## 575 C.2 Variations of the Bigram-Backcopy task

576 **Bigram-Backcopy task without the <bos> token.** We train a one-layer transformer on the BB  
 577 task without the <bos> token. Figure 16 shows that the <bos> token is perhaps not the extreme token.  
 578 Instead, trigger tokens and delimiter tokens seem to become extreme tokens. The results indicate  
 579 that initial tokens may not be the only candidates for the extreme token, partially explaining why  
 580 delimiter tokens could also be extreme tokens in LLMs.

581 **The Bigram-Skip-one (BS) task.** We make slight modifications to the Bigram-Backcopy task.  
 582 On trigger tokens, instead of copying the preceding token, we sample from the bigram-probability  
 583 of the preceding token  $P(\cdot | \text{Second-to-last token})$ . We train a one-layer transformer on it using the  
 584 same configuration as the BB task. Figure 17 shows that extreme token phenomena are mitigated.  
 585 The reason is that trained under BS, both the value states  $Va_{1,v}$  and the token embedding  $ebd_v$  give  
 586 the logit of the bigram transition probability. Therefore, other than having attention sink on the  
 587 <bos> token, self-attention becomes a new possibility to achieve the active-dormant mechanism.

## 588 D More Attention Heads in Dormant and Active Phase

589 In this section, we present two more dormant- and active- phase heads in Llama 2-7B-Base, in  
 590 Figures 18 and 19, which are more difficult to interpret than Layer 16 Head 25, but go dormant on  
 591 some inputs and active on others.

## 592 E Fine-Grained Static Mechanisms for Extreme-Token Phenomena

593 In this section, we will identify more fine-grained static mechanisms for extreme-token phenomena  
 594 in Llama 3.1-8B-Base. To do this, we identify circuits for the origin of attention sinks and small  
 595 value states. Then, using ablation studies, we study the origin of massive norms. Again, we use the  
 596 generic test phrase “<bos> Summer is warm. Winter is cold.”

597 **Attention sinks and global contextual semantics.** There are many attention sinks at layer 0, and  
 598 the <bos> token is always the sink token (see Figure 20). From now on until the end of this section,  
 599 we restrict our attention to Head 31 of Layer 0, which is an attention sink. These attention sinks are  
 600 caused by two linear-algebraic factors, demonstrated in Figure 21.

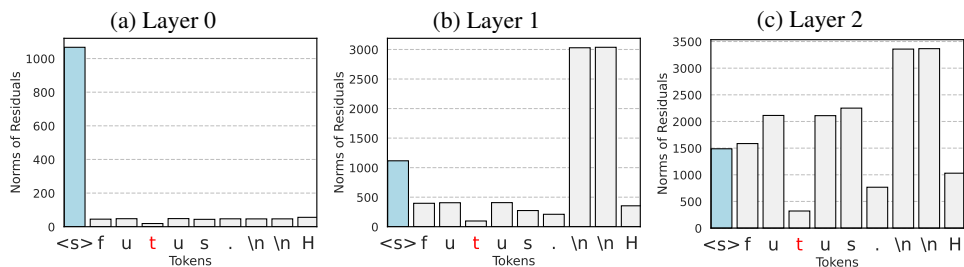


Figure 13: Residual state norms of three-layer transformer trained on the BB task

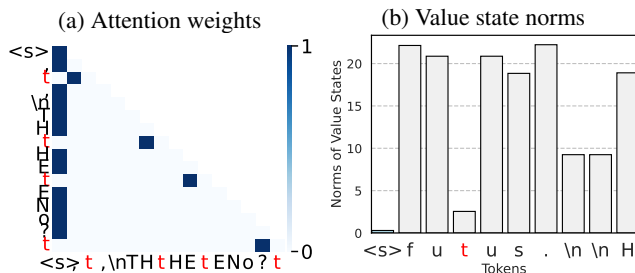


Figure 14: The simplified model structure trained on the BB task.

1. The key state of the <bos> token has small dot product with all other tokens.
2. The query states of all tokens are nearly orthogonal to the key states of all tokens except the <bos> token.

These two facts combine to ensure that the key state of the <bos> token is picked out by each query state, causing the attention sink. Since these query and key states are produced without any cross-token interaction, the alignment of different states is caused purely by the token’s global importance or meaning imparted via pretraining. The <bos> token has no semantic meaning in the context of prose tokens, so its key state is not aligned with key states of meaningful prose tokens. Also, delimiter tokens, oft considered secondary attention sinks (c.f. Appendix F.2), have the most aligned key states to the key state of the <bos> token, and are also the tokens with the least semantic meaning in the prose context. Thus, we identify that, at least in this restricted example, query state and key state alignment depends heavily on the contextual semantics of the token.

**Value state drains.** The value states of the <bos> token at Layer 0 Head 31 are already near zero, as demonstrated in Figure 22. While the delimiter tokens, which are less semantically meaningful in the prose context, have smaller value states than the rest, they are not as small as the value state of the <bos> token which is guaranteed to not have any semantics.

**Residual state peaks.** Residual state peaks are caused by the first two layers’ MLPs. In particular, we perform several ablations, comparing between the residual state norms in a later layer (24) of an un-edited forward pass versus forward passes where we force the output of either multiple layers, a single layer, an attention block, or an MLP to be zero (and hence remove its contribution from the residual stream). This intervention showed that ablating *either* Layer 0’s or Layer 1’s MLP is sufficient to remove the residual state peak. In particular, the second-largest token at Layer 24 in *each* ablation (including the original setup) has norm between 29 and 38, so the interventions ensure that all tokens have similar size.

## F Assorted Caveats

### F.1 Multiple Attention Sinks vs. One Attention Sink

As we have seen, attention heads in the BB task (Section 2), Llama 2-7B-Base (Section 3.1), and OLMo (Section 3.2) exhibit multiple attention sinks. That is, when heads in these models are

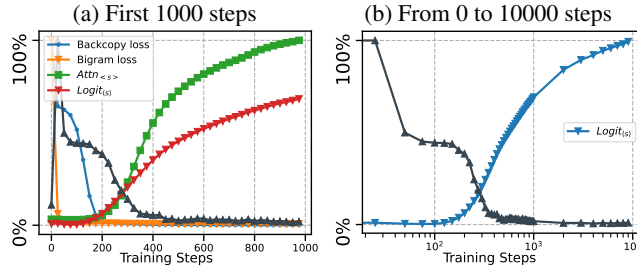


Figure 15: The dynamics of the simplified model structure trained on the BB task. *Left (a)*: The training curves match the one-layer transformer. *Right (b)*: The logit curve is close to the logarithmic growth predicted in Theorem 3.

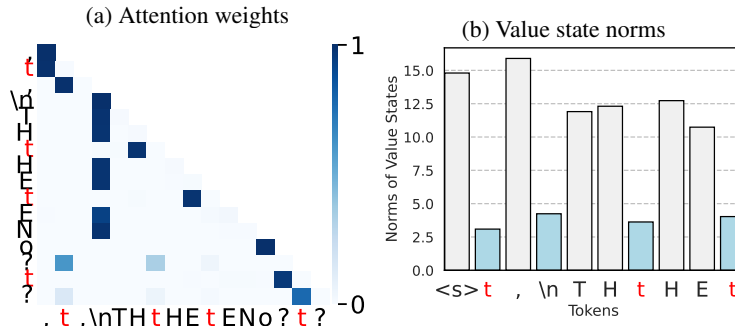


Figure 16: Attention weights and value state norms of a one-layer transformer trained on the BB task without the `<bos>` token.

629 dormant, they tend to have two attention sinks. For the LLMs in this group, at least on prose data, the  
 630 `<bos>` token as well as the first delimiter token (e.g., representing `.` or `;`) are sink tokens. Meanwhile,  
 631 Llama-3.1-8B-Base (Section 3) only ever has one attention sink on prose data, and the `<bos>` token  
 632 is always the sink token. Here, we offer a possible explanation of this phenomenon. For the BB  
 633 task, multiple sink tokens are necessary to solve the task. For LLMs, we believe this distinction may  
 634 be explained by the relative proportion of coding data, in which delimiters have a greater semantic  
 635 meaning than prose, within the training set. For instance, OLMo was trained on DOLMA [35], which  
 636 has around 411B coding tokens. Meanwhile, Llama 2 used at most  $(2T \times 0.08 =) 0.16T$  coding  
 637 tokens. Finally, Llama 3.1 used around  $(15.6T \times 0.17 =) 2.6T$  coding tokens [12]. On top of the raw  
 638 count being larger, coding tokens are a larger proportion of the whole pre-training dataset for Llama  
 639 3.1 compared to other model families. Thus, during training, the presence of delimiters would not be  
 640 considered unhelpful towards next-token prediction, since such delimiters carry plenty of semantics  
 641 in a wide variety of cases. Our earlier hypothesis in Section 3.1 proposes that only tokens which lack  
 642 semantics in almost all cases are made to be sink tokens. This could be a reason for the distinction.

## 643 F.2 The Role of a Fixed `<bos>` Token in the Active-Dormant Mechanism

644 Some models, such as OLMo, are not trained with a `<bos>` token. Despite this, the first token of  
 645 the input still frequently develops into a sink token. We can study the effect of positional encoding  
 646 of the tokens on the attention sink phenomenon by shuffling the tokens before inputting them into  
 647 the transformer, and observing how and why attention sinks form. If we do this with the phrase  
 648 “Summer is warm. Winter is cold.” with OLMo, we observe that at Layer 24, there are many attention  
 649 sink heads where the first token and first delimiter token share attention mass, even if the sentence  
 650 is jumbled up and makes no grammatical sense. This points towards the observation that without  
 651 a `<bos>` token, the attention sink formation uses both positional data and, to a greater degree, the  
 652 semantic data of each token. We leave studying this effect in greater detail to future work.

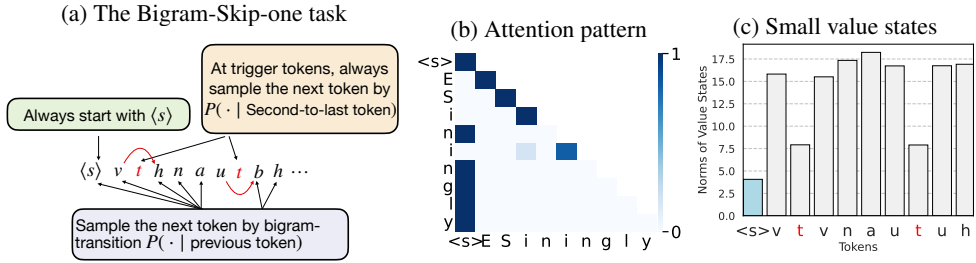


Figure 17: **Experiments on the Bigram-Skip-one task.** All phenomena are close to those in the BB task, but with diagonal attention sinks and relatively larger  $\|\text{val}_{\langle \text{bos} \rangle}\|$  compared with Figure 2.

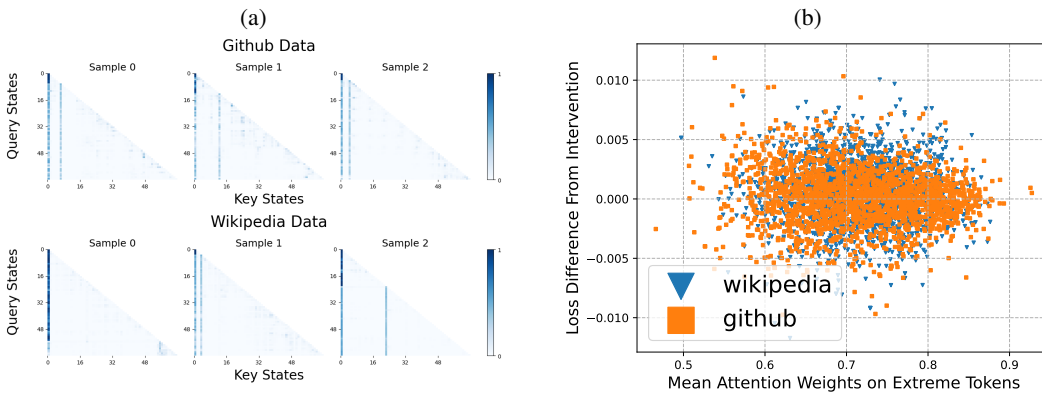


Figure 18: **Layer 16 Head 20 of Llama 2-7B-Base.**

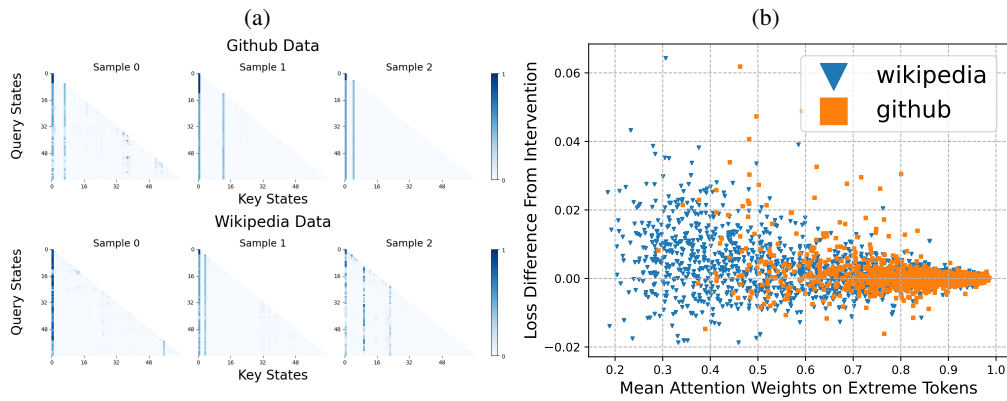


Figure 19: **Layer 16 Head 28 of Llama 2-7B-Base.**



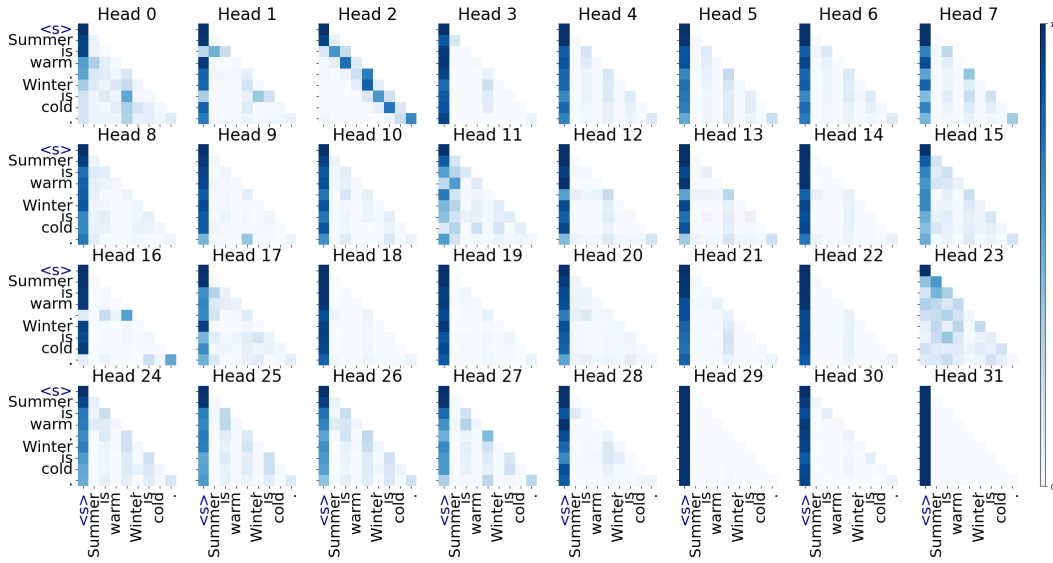


Figure 20: A visualization of attention heads at Layer 0 of Llama 3.1-8B-Base. Notice that many heads have the attention sink property, even at Layer 0 without any cross-token interaction. As usual, the test phrase is “Summer is warm. Winter is cold.” The most clear attention sink is Head 31.

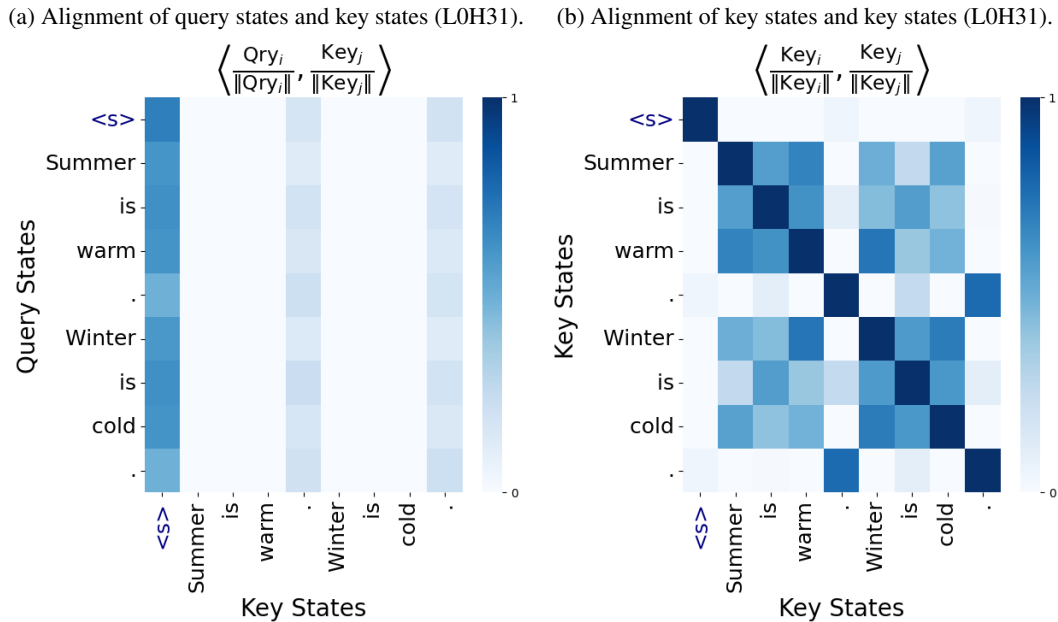


Figure 21: Alignment between query states and key states at Layer 0 Head 31 of Llama 3.1-8B-Base. We observe that the key state of <bos> is orthogonal to all other key states, and heavily aligned with all query states. Meanwhile, all semantically meaningful (i.e., not delimiter) tokens have aligned key states.

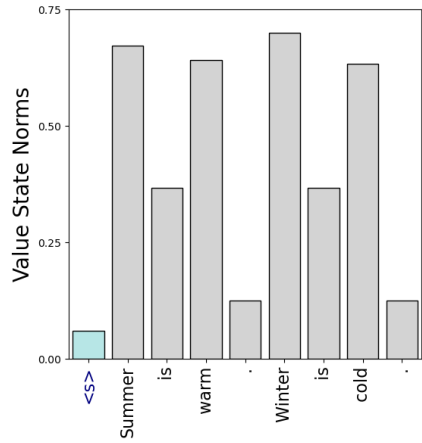


Figure 22: **Value state drains at Layer 0 Head 31 of Llama 3.1-8B-Base.** We observe that the value state associated with <bos> is already much smaller than every other semantically meaningful token, and still smaller than the delimiter tokens in the same sentence.

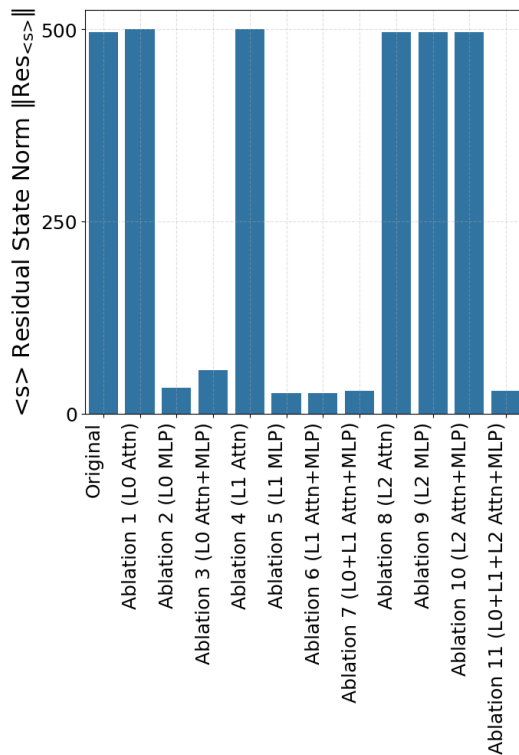


Figure 23: **Ablation study on the cause of the residual state peak in Llama 3.1-8B-Base.** We perform a series of ablations to understand which components of the network promote the residual state peaks. We find that ablating either the zeroth or first layer’s MLP is sufficient to remove the residual state peak phenomenon, while no other layer-level ablation can do it.

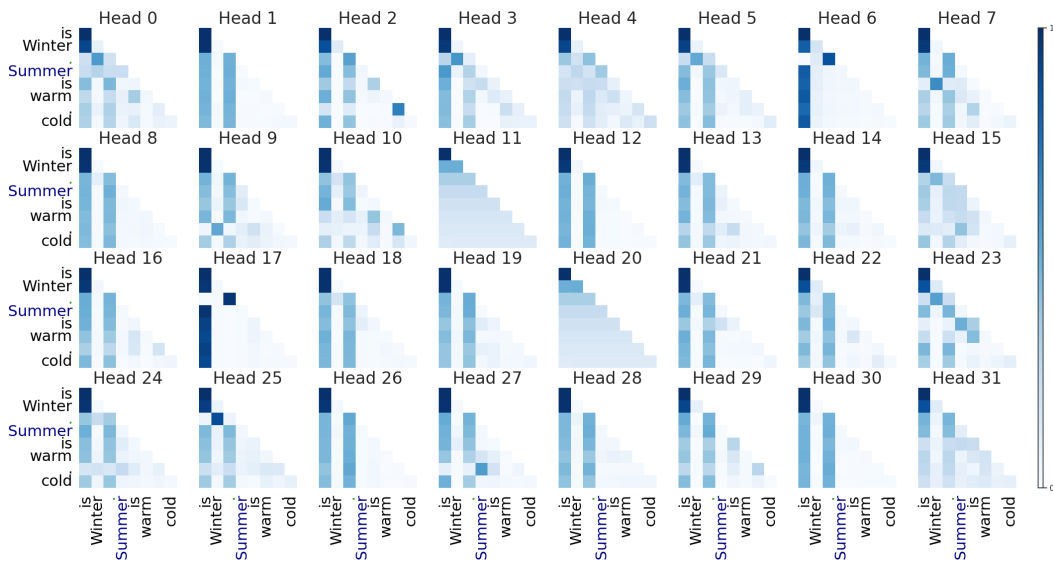


Figure 24: **Attention sinks with shuffled input in Layer 24 of OLMo.** In order to understand the impact of positional encodings when there is no `<bos>` token, we shuffle the input of the test string “Summer is warm. Winter is cold.” in OLMo. We observe that there is still an attention sink on token 0, despite it being a random token that does not usually start sentences or phrases (since it is uncapitalized). This shows that the positional embedding, say via RoPE, has a large impact on the formation of attention sinks — when the semantics of each token have switched positions, the attention sink still forms on the zeroth token.