# ACTIVE-DORMANT ATTENTION HEADS: MECHANISTICALLY DEMYSTIFYING EXTREME-TOKEN PHENOMENA IN LLMS

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

025

026

027

028 029 030

031

037

038

039

040

041

042

043

Paper under double-blind review

### ABSTRACT

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

### 1 INTRODUCTION

Recent analyses of transformer-based open-source large language models (LLMs), such as GPT-2
(Radford et al., 2019), Llama-2 (Touvron et al., 2023), Llama-3 (Dubey et al., 2024), Mixtral (Jiang
et al., 2023), Pythia (Biderman et al., 2023), and OLMo (Groeneveld et al., 2024), have revealed
several intriguing phenomena:

- Attention sinks (Xiao et al., 2023): In many attention heads, the initial token consistently attracts a large proportion of attention weights. In certain LLMs, other special tokens, such as the delimiter token, also draw significant attention. We refer to these as *sink tokens*.
- Value state drains (Guo et al., 2024): The value states of sink tokens are consistently much smaller than those of other tokens.
  - **Residual state peaks** (Sun et al., 2024): The intermediate representations of sink tokens, excluding those from the first and last layers, exhibit a significantly larger norm than other tokens.

044 These phenomena often appear simultaneously, and we collectively refer to them as the extreme-045 token phenomena. Figure 1 illustrates these phenomena using a fixed prompt: " $\langle s \rangle$  Summer is warm. Winter is cold." in Llama-3.1-8B-Base, where the first token,  $\langle s \rangle$ , the beginning-of-sentence 046 token, serves as the sink token. We note that the first token does not have to be  $\langle s \rangle$  to function 047 as a sink token, as in GPT-2, where other tokens, being the initial token, can also serve this role. 048 Furthermore, in models like Llama-2, a delimiter token can also act as the sink token. Despite 049 the consistency of these observations, no prior work has provided a satisfying explanation for the mechanisms behind these phenomena. As a tentative explanation, Xiao et al. (2023) suggested that 051 models tend to dump unnecessary attention values to specific tokens. 052

053 This work aims to demystify the extreme-token phenomena in LLMs. We show that the extremetoken phenomena are manifestations of the *active-dormant mechanism* of attention heads. We sup-



Figure 1: Extreme-token phenomena in Llama 3.1-8B-Base. We evaluate the sentence "(s) 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  $\langle s \rangle$  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  $\langle s \rangle$  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  $\langle s \rangle$  token is much smaller than those of other tokens on average.

077

078 079

081

082

084

085

090

091

093

094

095

096

067

068

069

071

073

port this claim through studies on simplified transformer architectures and tasks, a dynamical theory of simplified models, and experiments on pre-trained LLMs. Our contributions are as follows:

- 1. In Section 2, we train one-to-three-layer transformers on the *Bigram-Backcopy* (BB) task, which also exhibits extreme-token phenomena similar to those observed in LLMs. We show that attention sinks and value-state drains are driven by the active-dormant mechanism mechanism. Both theoretically and empirically, we demonstrate that the mutual reinforcement dynamics underpin the extreme-token phenomena: attention sinks and value-state drains reinforce each other, leading to a stable phase where all query tokens produce identical attention logits for the keys of extreme tokens. Empirical evidence further shows that residual state peaks result from the interaction between this mutual reinforcement mechanism and Adam.
- 2. In Section 3, we demonstrate the *active-dormant mechanism* mechanism in LLMs by identifying an interpretable active-dormant head (Layer 16, Head 25 in Llama 2-7B-Base (Touvron et al., 2023)), confirmed through causal intervention analyses. We also discover circuits in LLMs related to extreme tokens that partially align with models trained on the BB task. Examining the dynamics of OLMo-7B-0424 (Groeneveld et al., 2024), we observe the same mutual reinforcement mechanism and stable phase, consistent with predictions from the BB task. 092
  - 3. Through causal interventions, we isolate the extreme-token phenomena to architecture and optimization strategy. Specifically, we show that replacing SoftMax with ReLU activations in attention heads can eliminate extreme-token phenomena in the BB task, and switching from Adam to SGD removes the residual-state peak phenomenon in the BB task. Our work demonstrates potential classes of modifications to mitigate extreme-token phenomena in LLMs.
- 097 098
- 099 100

101

1.1 NOTATION

102

103 We denote the SoftMax attention layer with a causal mask as attn, the MLP layer as mlp, and the 104 transformer block as TF. The query, key, value states, and residuals of a token v are represented as 105  $Qry_{u}$ , Key<sub>u</sub>, Val<sub>v</sub>, and Res<sub>v</sub>, respectively, with the specific layer and head indicated in context. We use  $\langle s \rangle$  to refer to the "Beginning of Sequence" (bos) token. Throughout the paper, we employ zero-106 indexing (i.e., attention head and layer indices start from 0 rather than 1) for consistency between 107 code and writing.



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  $\langle s \rangle$  token has a tiny norm.

124

125

126

127

128

129

130

## 2 THE BIGRAM-BACKCOPY TASK

The Bigram-Backcopy task consists of two sub-tasks: *Bigram-transition* and *Backcopy*. Each input sequence begins with a  $\langle s \rangle$  token, followed by tokens sampled according to a pre-determined bigram transition probability P. When some special trigger tokens are encountered, instead of sampling, the preceding token is copied to the next position. Following Bietti et al. (2024), we select the transition P and the vocabulary  $\mathcal{V}$  with  $|\mathcal{V}| = V = 64$  based on the estimated character-level bigram distribution from the tiny *Shakespeare* dataset. In all experiments, the set of trigger tokens  $\mathcal{T}$  is fixed and consists of the  $|\mathcal{T}| = 3$  most frequent tokens in the unigram distribution. Thus, the non-trigger token set,  $\mathcal{V} \setminus \mathcal{T}$ , comprises 61 tokens.

- 131 132
- 133 134

### 2.1 ONE-LAYER TRANSFORMER SHOWS ATTENTION SINKS AND VALUE-STATE DRAINS.

135 On the Bigram-Backcopy task, we pre-train a standard one-layer transformer with only one soft-136 max attn head and one mlp layer. Unless otherwise specified, the model is trained with Adam 137 for 10,000 steps. We relegate the training details in Appendix C. Figure 2b shows that the trained 138 transformer also exhibits the attention sink phenomenon, where the  $\langle s \rangle$  token captures a significant 139 proportion of the attention weights. More importantly, the attention weights reveal interpretable pat-140 terns: all non-trigger tokens exhibit attention sinks, while the attention for trigger tokens is concen-141 trated on their preceding positions. Furthermore, Figure 2c reveals a value state drain phenomenon 142 similar to LLMs, indicating that on non-trigger tokens, the attn head adds a minimal value to the residual stream. 143

144

145 The active-dormant mechanism of the attention head: Inspired by the observed interpretable 146 attention weight patterns, we propose the *active-dormant mechanism*. For any given token, an at-147 tention head is considered *active* if it contributes significantly to the residual state, and *dormant* if 148 its contribution is minimal. As illustrated in Figure 2b, trained on the BB task, the attention head is 149 active on trigger tokens and dormant on non-trigger tokens.

150 Figure 3a demonstrates that the mlp layer is responsible for the Bigram task whereas the attn head 151 takes care of the Backcopy task. When the mlp layer is zeroed out, the backcopy loss remains signif-152 icantly better than a random guess, but the bigram loss degrades to near-random levels. Conversely, when the attn layer is zeroed out, the backcopy loss becomes worse than a random guess, while 153 the bigram loss remains unaffected. This suggests that on trigger tokens, the attn head is active 154 and handles the backcopy task, whereas on non-trigger tokens, the attn head is dormant, allow-155 ing the mlp layer to handle the Bigram task. We summarize the active-dormant mechanism of the 156 attn head in Claim 1. 157

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

160 161

• **Dormant phase**: On non-trigger tokens, the attn head puts dominant weights to the  $\langle s \rangle$  token, adding minimal value to the residual stream, having little impact on the model's output.



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  $\langle s \rangle$  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 logit<sub>(s)</sub> curve denotes the mean of attention logits from all given non-trigger query tokens v on the  $\langle s \rangle$  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  $\langle s \rangle$  tokens.

• 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.

The growth of attention logits on the  $\langle s \rangle$  token and the decrease in the norm of its value state. 200 Figure 3b displays the training dynamics of excess risks, attention weights, attention logits, and value 201 state norms for the  $\langle s \rangle$  token. All values are rescaled to highlight the trends. The backcopy excess 202 risk and the bigram excess risk both drop to zero within the first 1000 steps. As the backcopy risk 203 decreases, the attention weights on the  $\langle s \rangle$  token increase, suggesting a relationship between the for-204 mation of attention sinks and the functional development of the attention heads. For each token  $v_n$  at 205 position n in the prompt, we compute  $\text{logit}_{(s)} = \text{mean}_n[\langle \text{Qry}_{v_n}, \text{Key}_{(s)} \rangle - \text{mean}_i(\langle \text{Qry}_{v_n}, \text{Key}_{v_i}) \rangle],$ 206 which serves as a progress measure for attention sinks. Even after the attention weights on the 207  $\langle s \rangle$  token is nearly 1, logit<sub>(s)</sub> continues to increase. Simultaneously, the norm of the value state of 208 the  $\langle s \rangle$  token continues to decrease to a small value.

209

181

183

186

187

188

189

190 191

192

193 194

196

197

199

210 211

2.2 ANALYSIS OF A MINIMALLY-SUFFICIENT TRANSFORMER ARCHITECTURE

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  $\langle s \rangle$  token, and  $\mathcal{T}$  denote the set of all trigger tokens. Given any  $v \in \mathcal{V}$ , we denote  $p_{vk} = \mathsf{P}(k|v)$  to be the next token Markov transition probability, and  $\mathbf{p}_v = [p_{v1}, \ldots, p_{vV}]$  be the row vector in the simplex. We assume that the tokens are embedded into V-dimensional space 216 using one-hot encoding, and for notation simplicity, we abuse v to stand for its one-hot encoding 217 vector  $e_v \in \mathbb{R}^V$  which is a row vector. The predicted probability of the n+1 token is given by 218 SoftMax $(TF([\langle s \rangle; v_{1:n-1}; v])_n)$ , where transformer architecture is given by  $TF(\cdot) = \mathtt{attn}(\cdot) +$  $mlp(\cdot)$ . Here  $attn(\cdot) = SoftMax(mask(Qry(\cdot)Key(\cdot)^{\top}))Val(\cdot)$  and (Qry, Key, Val) are linear maps from  $\mathbb{R}^V \to \mathbb{R}^V$ . Since the mlp layer could handle the Bigram task, we assume that mlp 219 220 outputs the Markov transition probabilities  $\mathbf{p}_v$  on non-trigger tokens v and zero on trigger tokens. 221 For the attn head, we assume that the attention logits on the  $\langle s \rangle$  key-token are  $(\alpha_{v_1}; \ldots; \alpha_{v_n})$ , 222 the attention logits on any trigger query-token are  $(0, \ldots, \lambda, 0)$  where the second last coordinate 223 is  $\lambda$ , and assume other logits are zero. Assume that the value state of  $\langle s \rangle$  is  $\beta \in \mathbb{R}^V$ , and the 224 value state of each non-trigger token v is a one-hot encoding vector  $e_v$  multiplied by  $\xi_v \ge 0$ . 225 Figure 4 illustrates this simplified transformer architecture. These assumptions are summarized in 226 the following equations. 227

229 230

231 232

246

258

259

260

261

262

264 265

266

267

268

269

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

$$\mathtt{Val}(v) = \xi_v \boldsymbol{e}_v \quad ext{with} \ \xi_v = 0 \ ext{for} \ v \in \mathcal{T}, \ ext{and} \ \xi_v \geq 0 \ ext{for} \ v \in \mathcal{V} \setminus \mathcal{T}.$$

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

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 transformer (mlp, Qry, Key, Val) such that Eq. (1) holds. The transformer gives ground truth transition of the BB model if min<sub> $v \in V$ </sub>  $\alpha_v \to \infty$ , min<sub> $v \in V$ </sub>  $\xi_v \to \infty$ ,  $\lambda \to \infty$ , and  $\beta = 0$ .

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

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

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

$$\operatorname{loss}_{v}(\alpha_{v},\boldsymbol{\beta}) = \sum_{k=1}^{V} p_{vk} \Big\{ \log \Big[ \sum_{i=1}^{V} p_{vi} \exp \Big( \frac{e^{\alpha_{v}} \beta_{i} + W_{i} \xi_{i}}{e^{\alpha_{v}} + W} \Big) \Big] - \frac{e^{\alpha_{v}} \beta_{k} + W_{k} \xi_{k}}{e^{\alpha_{v}} + W} - \log p_{vk} \Big\}.$$

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

$$\mathsf{loss}(oldsymbol{lpha},oldsymbol{eta}) = \sum_{v \in \mathcal{V} ackslash \mathcal{T}} \pi_v \mathsf{loss}_v(lpha_v,oldsymbol{eta}).$$

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

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

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

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

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

<sup>&</sup>lt;sup>1</sup>We note that Reddy (2023) makes similar simplification in analyzing induction heads.

271

272

273

274

283

308

310

311

312

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

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

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

275 The proof of Theorem 3 is provided in Appendix B.2. We give three key remarks: (1) As  $\alpha_v \to \infty$ , 276 a Taylor expansion of the gradient  $\partial \log/\partial \alpha_v$  suggests that  $d\alpha_v/dt \propto \exp(-2\alpha_v)$ , which leads 277 to the logarithmic growth of  $\alpha_v$ . Similar logarithmic growth exists in the literature under different 278 setups (Tian et al., 2023a; Han et al., 2023). (2) For a fixed  $\alpha = \alpha \mathbf{1}$ , under additional assumptions 279 on the initial value  $\beta(0)$ , we can prove a linear convergence for  $\beta$ . (3) The stable phase described 280 in Theorem 3 seems to imply that the system could be stable without attention sinks, as it does not 281 require  $\alpha$  to be large. However, in practice, models trained on the BB task tend to converge to a 282 stable phase where  $\alpha$  is relatively large.

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

290 **Experimental verification.** Revisiting Figure 3b, which shows the dynamics of a full transformer 291 model trained with Adam, we observe that both  $logit_{(s)}$  and  $||Val_{(s)}||_2$  exhibit growth rates con-292 sistent with Theorem 3. The logit<sub>(s)</sub> is equivalent to  $\alpha$  in this context, as all other attention logits 293 are assumed to be zero under the setup of Theorem 3. When plotted on a logarithmic scale, the 294 logit<sub>(s)</sub> curve grows approximately linearly between 1,000 and 10,000 steps, then accelerates before 295 stabilizing around 100,000 steps. Meanwhile, the norm of the value state decreases monotonically. 296 The simultaneous increase in attention weights and decrease in value-state norms suggest that these 297 phases occur together during the training process. To further validate Theorem 3, we construct a sim-298 plified model that aligns with Equ. (1), and train the parameters  $(\alpha \in \mathbb{R}^V, \beta \in \mathbb{R}^V, \xi \in \mathbb{R}^V, \lambda \in \mathbb{R})$ 299 with Adam. The resulting training curves are similar to those of a one-layer transformer, also ex-300 hibiting the mutual reinforcement mechanism.

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

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

- 1. The SoftMax mechanism pushes the attention weights to the value state drains, reinforcing attention sinks.
- 2. The attention sinks on the value state drains further pushes down the value state, reinforcing value state drains.
- The mutual reinforcement stabilizes at the phase when all tokens have identical large attention logits on the value state drains. Finally, due to the causal mask, the training dynamics favor the  $\langle s \rangle$  token to become an extreme token.

316 We expect that the formation of extreme tokens in LLMs follows a similar mutual reinforcement 317 mechanism. Indeed, although Theorem 3 focuses on a specific BB task with a simplified architecture 318 and loss function, the same principles can be applied to more general scenarios. Specifically, for an 319 attention head attn, we assume that  $(LLM \setminus attn)(v) = \log p_v$ , meaning that the LLM, even if 320 we zeroed out attn, can still output an accurate next token prediction. Furthermore, we assume 321  $Val(v) = \xi_v e_v$ , indicating that adding the value state from any previous tokens performs a specific function. Under these assumptions, we expect the same theoretical results to apply to LLMs. In 322 Section 3, we will explore the formation of attention sinks and value state drains along the training 323 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.

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

335

336

337

338

339 340 341

342

343

344

345

346

347

THE EMERGENCE OF RESIDUAL STATE PEAKS 2.3

352 The residual state peaks require a three-layer structure. No residual state peaks appear in a 353 one-layer transformer trained on the BB task. We train various models on the BB task and track 354 the  $\langle s \rangle$  token's residual state norms after layer 0. We relegate the experimental results to Appendix C. We find that a three-layer transformer is enough to produce residual state peaks. If we allow 355 to skip some mlp or attn layers, the "attn+mlp+attn+mlp" combination becomes the 356 simplest model that produces residual state peaks (Figure 10). Circuit analysis also reveals that 357 LLMs typically add a large vector in the first layer and cancel it in the last layer. We propose that 358 the add-then-cancel mechanism is essential for residual state peaks and requires at least three layers. 359

360

Residual state peak reinforces attention sinks and value state drains in trained models. Figure 361 5b presents the intervention results on the "attn+mlp+attn+mlp+mlp" model. We recenter the 362  $\|\text{Res}_{(s)}\|_2$  by subtracting the average norm of other tokens from the  $\langle s \rangle$  token norm. The logit<sub>(s)</sub> 363 and  $\|Val_{(s)}\|$  are computed in layer 1 following the same ways as in Figure 3b. When layer 0 is 364 zeroed out, the residual norm returns to normal, attention logits decrease, and the value state norm rises. It verifies that the residual state peak contributes to the attention sink and value state drain 366 phenomenon in the trained transformer. 367

368 **Replacing Adam by SGD removes the linear growth of residual state norm.** Figure 5c shows 369 the  $\langle s \rangle$ 's residual state norms at the output of layer 0 of three-layer transformers with different 370 configurations. Adam leads to a linear increase in residual norms. In contrast, with SGD, attention sinks persist, but residual state peaks vanish. The ReLU attention, which lacks the active-dormant 372 mechanism, shows no residual state peaks.

- 373
- 374 375

371

#### EXTENDING PREDICTIONS OF THE BB MODEL TO LLMS 3

376

In this section, we examine extreme-token phenomena in open-source pre-trained LLMs. In Sec-377 tion 3.1, we analyze the static behavior of these phenomena in Llama 2-7B-Base (Touvron et al., 2023), confirming that certain attention heads in LLMs exhibit both active and dormant phases. No-tably, we identify a specific head that is active on GitHub samples but dormant on Wikipedia samples, illustrating the *active-dormant mechanism*. In Section 3.2, we explore the dynamic behavior of extreme-token phenomena during the pre-training process of OLMo-7B (Groeneveld et al., 2024).
We show that the attention logits, value state norms, and residual state norms of the sink token(s) in OLMo mirror their behavior in the simpler BB model. Specifically, the simultaneous formation of attention sinks and value state drains gives evidence for the *mutual reinforcement mechanism*.

385 386

387

388

389 390

391

392

3.1 ACTIVE-DORMANT MECHANISM IN LLMS

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

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

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

398 Accordingly, we strive to find instances of heads in pretrained LLMs which satisfy this principle, i.e., 399 which are dormant on some domains and active on others. In Figure 6, we show a particular attention head – Layer 16 Head 25 of Llama 2-7B-Base (Touvron et al., 2023) — which has an extremely clear 400 active-dormant distinction across two distinct contexts (e.g., tokens from RedPajama (Computer, 401 2023) drawn from the GitHub subset versus the Wikipedia subset). While there are many such 402 attention heads which are context-dependent — we provide some in Appendix D — we demonstrate 403 this one because the conditions under which it is active are simple and interpretable, while others 404 have more involved or complex criteria to become active. We observe that this attention head is 405 dormant (i.e., an attention sink) on samples from Wikipedia, which more closely resemble prose, 406 and active (i.e., not an attention sink) on samples from Github, which more closely resemble code. 407 We also observe that this attention head, in general, contributes significantly to the performance of 408 the model on code sequences, but has negligible impact on the performance of the model on prose 409 sequences (Figure 6b). This is a further justification, from a practical perspective, of why this head is sometimes dormant and sometimes active — in some contexts we can ablate it from the model 410 entirely with no effect, but in other contexts ablating the head leads to huge performance drops. We 411 include more detail in Appendix E, where we extract a circuit for extreme-token phenomena in order 412 to analyze the dormant-active mechanism and its interaction with the semantics of the input tokens. 413

414 415

416

417

418 419

420

421

422

3.2 TRAINING DYNAMICS OF EXTREME-TOKEN PHENOMENA IN LLMS

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

The attention heads go through a attention-increasing and value-state-shrinking phase. They then go into a stable phase, with identical attention logits on the  $\langle s \rangle$  token. Meanwhile, the residual state norm of the  $\langle s \rangle$  token linearly increases during pre-training.

We confirm these predictions below. To observe the training dynamics of a large-scale LLM, we use the setup of OLMo-7B-0424 (Groeneveld et al., 2024) (henceforth just referred to as OLMo), who have open-sourced weights at several steps during their training run. For our analysis, we inspect OLMo at a variety of training steps: every 500 steps throughout the first 10,000 steps, then 25,000 steps, then 50,000 steps, then every 50,000 steps until 449,000 steps (which is roughly the end of their training). Again, we use the input "Summer is warm. Winter is cold.".<sup>2</sup> Notice that in this prompt, token 3, namely ".", is not very semantically meaningful; it becomes a sink token along with token 0 (c.f. Section 3.1, Appendix E, Appendix F.2).

<sup>&</sup>lt;sup>2</sup>Note that OLMo does not have a  $\langle s \rangle$  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 (Computer, 2023) 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.

In Figure 7, we confirm that attention heads go through an attention-increasing and value-stateshrinking phase, and that the residual state norm of the  $\langle s \rangle$  token increases linearly during pretraining. We show that, at Layer 24 of OLMo, the average attention on extreme tokens (token 0 and token 3) increases rapidly at the beginning of training and converges to a constant, while the value state norms of extreme tokens decrease rapidly. Also, the residual states of extreme tokens also increase linearly, while the rest quickly converge. In Figure 8 we show that attention heads converge to a stable phase, and that all logits corresponding to the first token's value states (i.e., all tokens' value of  $logit_0$ , except possibly the value of  $logit_0$  corresponding to token 0 itself) have 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. 



Figure 8: Stable phase. Left (a): We plot the normalized attention logits of all tokens' query states against to-499 ken 0's key state during training. We observe that the logits of all non-extreme tokens' query states against token 500 0's key state in OLMo's Layer 24 are stable for a large fraction of the training run, after an initialization period. 501 This echoes the stable phase prediction made in the BB model in Section 2. Note that this prediction makes no 502 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  $logit_i = \langle Qry_i, Key_0 \rangle - mean_i(\langle Qry_i, Key_i \rangle)$ . Right 504 (b): For this experiment, we generate 128 randomly sampled test tokens with IDs from 100 to 50000 in the 505 OLMo tokenizer. We append each token separately to the test phrase "Summer is warm. Winter is cold.", cre-506 ating 128 different samples, which we feed to the LLM to record the model behavior. We plot the distribution 507 of (un-normalized) dot products  $\langle Qry_{test}, Key_j \rangle$  across all heads at Layer 24 and all test tokens. We observe 508 that logits of all regular tokens have very similar distributions, and the distributions of the logits corresponding 509 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. 510

### 4 CONCLUSION

513 514

In this work, we investigated the extreme-token phenomena, namely attention sinks, value state 515 drains, and residual state peaks. We analyzed a simple evocative model called the Bigram-Backcopy 516 task, and theoretically and empirically showed that it exhibited the same extreme-token phenomena 517 as in LLMs. Based on the Bigram-Backcopy task, we made several detailed predictions about the 518 behavior of extreme-token phenomena in LLMs. In particular, we identified the active-dormant 519 mechanism for attention heads in both the BB model and LLMs, of which attention sinks and value 520 state drains are indicators, and a *mutual reinforcement mechanism* by which these phenomena are induced during pretraining. Using intuition about these mechanisms, we applied minor interventions 521 to the model architecture and optimization procedure which disabled extreme-token phenomena 522 within the BB model. Overall, our work uncovers the causes of extreme-token phenomena and 523 points to possible pathways to eliminate them during LLM training. 524

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

- 531
- 532
- 53
- 534
- 533
- 527

538

# 540 ETHICS STATEMENT

541	
542	This paper contributes towards the analysis of large language models. This paper does not add any
543	ethical concerns beyond the usual ethics associated with use and analysis of large language models.
544	
545	
546	
547	
548	
549	
550	
551	
552	
555	
555	
556	
557	
558	
559	
560	
561	
562	
563	
564	
565	
566	
567	
568	
569	
570	
571	
572	
573	
574	
575	
576	
577	
578	
579	
580	
500	
583	
584	
585	
586	
587	
588	
589	
590	
591	
592	
593	

# 594 REFERENCES

602

609

621

633

- Kwangjun Ahn, Xiang Cheng, Minhak Song, Chulhee Yun, Ali Jadbabaie, and Suvrit Sra. Lin ear attention is (maybe) all you need (to understand transformer optimization). arXiv preprint
   arXiv:2310.01082, 2023.
- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Zeyuan Allen-Zhu and Yuanzhi Li. Physics of language models: Part 1, context-free grammar. *arXiv preprint arXiv:2305.13673*, 2023.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pp. 2397–2430. PMLR, 2023.
- Alberto Bietti, Vivien Cabannes, Diane Bouchacourt, Herve Jegou, and Leon Bottou. Birth of a transformer: A memory viewpoint. Advances in Neural Information Processing Systems, 36, 2024.
- François Charton. What is my math transformer doing?-three results on interpretability and gener alization. *arXiv preprint arXiv:2211.00170*, 2022.
- Liang Chen, Haozhe Zhao, Tianyu Liu, Shuai Bai, Junyang Lin, Chang Zhou, and Baobao Chang.
  An image is worth 1/2 tokens after layer 2: Plug-and-play inference acceleration for large vision-language models. *arXiv preprint arXiv:2403.06764*, 2024.
- Together Computer. Redpajama: An open source recipe to reproduce llama training dataset, 2023.
   URL https://github.com/togethercomputer/RedPajama-Data.
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need
   registers. *arXiv preprint arXiv:2309.16588*, 2023.
- Puneesh Deora, Rouzbeh Ghaderi, Hossein Taheri, and Christos Thrampoulidis. On the optimization and generalization of multi-head attention. *arXiv preprint arXiv:2310.12680*, 2023.
- Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Gpt3. int8 (): 8-bit matrix
   multiplication for transformers at scale. *Advances in Neural Information Processing Systems*, 35: 30318–30332, 2022.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The Ilama 3 herd of models.
   *arXiv preprint arXiv:2407.21783*, 2024.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann,
   Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, et al. A mathematical framework for
   transformer circuits. *Transformer Circuits Thread*, 1:1, 2021.
- <sup>637</sup> Nelson Elhage, Robert Lasenby, and Christopher Olah. Privileged bases in the transformer residual stream. *Transformer Circuits Thread*, 2023.
- Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context? *arXiv preprint arXiv:2310.17191*, 2023.

# Yao Fu. How do language models put attention weights over long context? Yao Fu's Notion, 2024. URL https://yaofu.notion.site/ How-Do-Language-Models-put-Attention-Weights-over-Long-Context-10250219d5ce42e8b46 pvs=4.

647 Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir Globerson. Dissecting recall of factual associations in auto-regressive language models. *arXiv preprint arXiv:2304.14767*, 2023.

658

665

667

682

683

684 685

686

687

688

- 648 Dirk Groeneveld, Iz Beltagy, Pete Walsh, Akshita Bhagia, Rodney Kinney, Oyvind Tafjord, 649 Ananya Harsh Jha, Hamish Ivison, Ian Magnusson, Yizhong Wang, et al. Olmo: Accelerating the 650 science of language models. arXiv preprint arXiv:2402.00838, 2024.
- Tianyu Guo, Wei Hu, Song Mei, Huan Wang, Caiming Xiong, Silvio Savarese, and Yu Bai. How do 652 transformers learn in-context beyond simple functions? a case study on learning with representa-653 tions. arXiv preprint arXiv:2310.10616, 2023. 654
- 655 Zhiyu Guo, Hidetaka Kamigaito, and Taro Watanabe. Attention score is not all you need for token 656 importance indicator in ky cache reduction: Value also matters. arXiv preprint arXiv:2406.12335, 657 2024.
- Wes Gurnee, Theo Horsley, Zifan Carl Guo, Tara Rezaei Kheirkhah, Qinyi Sun, Will Hathaway, 659 Neel Nanda, and Dimitris Bertsimas. Universal neurons in gpt2 language models. arXiv preprint 660 arXiv:2401.12181, 2024. 661
- 662 Chi Han, Qifan Wang, Wenhan Xiong, Yu Chen, Heng Ji, and Sinong Wang. Lm-infinite: Simple 663 on-the-fly length generalization for large language models. arXiv preprint arXiv:2308.16137, 664 2023.
- Roger A Horn and Charles R Johnson. Matrix analysis. Cambridge university press, 2012. 666
- Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. arXiv preprint 668 arXiv:2310.05249, 2023. 669
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 670 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 671 Mistral 7b. arXiv preprint arXiv:2310.06825, 2023. 672
- 673 Juno Kim, Tai Nakamaki, and Taiji Suzuki. Transformers are minimax optimal nonparametric in-674 context learners. arXiv preprint arXiv:2408.12186, 2024. 675
- Ruikang Liu, Haoli Bai, Haokun Lin, Yuening Li, Han Gao, Zhengzhuo Xu, Lu Hou, Jun Yao, 676 and Chun Yuan. Intactky: Improving large language model quantization by keeping pivot tokens 677 intact. arXiv preprint arXiv:2403.01241, 2024. 678
- 679 Ziming Liu, Ouail Kitouni, Niklas S Nolte, Eric Michaud, Max Tegmark, and Mike Williams. To-680 wards understanding grokking: An effective theory of representation learning. Advances in Neu-681 ral Information Processing Systems, 35:34651–34663, 2022.
  - Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. Advances in Neural Information Processing Systems, 35:17359–17372, 2022.
  - Neel Nanda, Lawrence Chan, Tom Lieberum, Jess Smith, and Jacob Steinhardt. Progress measures for grokking via mechanistic interpretability. arXiv preprint arXiv:2301.05217, 2023.
  - Eshaan Nichani, Alex Damian, and Jason D Lee. How transformers learn causal structure with gradient descent. arXiv preprint arXiv:2402.14735, 2024.
- 690 Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, 691 Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, et al. In-context learning and induction 692 heads. arXiv preprint arXiv:2209.11895, 2022. 693
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language 694 models are unsupervised multitask learners. OpenAI blog, 1(8):9, 2019. 695
- 696 Gautam Reddy. The mechanistic basis of data dependence and abrupt learning in an in-context 697 classification task. In The Twelfth International Conference on Learning Representations, 2023. 698
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, 699 Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of 700 three trillion tokens for language model pretraining research. arXiv preprint arXiv:2402.00159, 701 2024.

702 703 704	Seungwoo Son, Wonpyo Park, Woohyun Han, Kyuyeun Kim, and Jaeho Lee. Prefixing atten- tion sinks can mitigate activation outliers for large language model quantization. <i>arXiv preprint</i> <i>arXiv:2406.12016</i> , 2024.
705 706 707	Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. Massive activations in large language models. <i>arXiv preprint arXiv:2402.17762</i> , 2024.
708 709 710 711	Yuandong Tian, Yiping Wang, Beidi Chen, and Simon S Du. Scan and snap: Understanding train- ing dynamics and token composition in 1-layer transformer. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 36:71911–71947, 2023a.
712 713 714 715	Yuandong Tian, Yiping Wang, Zhenyu Zhang, Beidi Chen, and Simon Du. Joma: Demystifying mul- tilayer transformers via joint dynamics of mlp and attention. <i>arXiv preprint arXiv:2310.00535</i> , 2023b.
716 717 718	Eric Todd, Millicent L Li, Arnab Sen Sharma, Aaron Mueller, Byron C Wallace, and David Bau. Function vectors in large language models. <i>arXiv preprint arXiv:2310.15213</i> , 2023.
719 720 721	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko- lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda- tion and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> , 2023.
722 723 724	Roman Vershynin. <i>High-dimensional probability: An introduction with applications in data science</i> , volume 47. Cambridge university press, 2018.
725 726 727 728	Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Inter- pretability in the wild: a circuit for indirect object identification in gpt-2 small. <i>arXiv preprint</i> <i>arXiv:2211.00593</i> , 2022.
729 730 731 732	Jingfeng Wu, Difan Zou, Zixiang Chen, Vladimir Braverman, Quanquan Gu, and Peter L Bartlett. How many pretraining tasks are needed for in-context learning of linear regression? <i>arXiv</i> preprint arXiv:2310.08391, 2023.
733 734 735	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming language models with attention sinks. <i>arXiv preprint arXiv:2309.17453</i> , 2023.
736 737 738	Zhongzhi Yu, Zheng Wang, Yonggan Fu, Huihong Shi, Khalid Shaikh, and Yingyan Celine Lin. Unveiling and harnessing hidden attention sinks: Enhancing large language models without training through attention calibration. <i>arXiv preprint arXiv:2406.15765</i> , 2024.
739 740 741 742 743	Shuangfei Zhai, Tatiana Likhomanenko, Etai Littwin, Dan Busbridge, Jason Ramapuram, Yizhe Zhang, Jiatao Gu, and Joshua M Susskind. Stabilizing transformer training by preventing attention entropy collapse. In <i>International Conference on Machine Learning</i> , pp. 40770–40803. PMLR, 2023.
744 745 746	Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. arXiv preprint arXiv:2306.09927, 2023.
747 748 749	Ruiqi Zhang, Jingfeng Wu, and Peter L Bartlett. In-context learning of a linear transformer block: benefits of the mlp component and one-step gd initialization. <i>arXiv preprint arXiv:2402.14951</i> , 2024.
750 751 752 753	Yi Zhang, Arturs Backurs, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, and Tal Wagner. Unveiling transformers with lego: a synthetic reasoning task. <i>arXiv preprint arXiv:2206.04301</i> , 2022.
754 755	Zeyuan Allen Zhu and Yuanzhi Li. Physics of language models: Part 3.1, knowledge storage and extraction. <i>arXiv preprint arXiv:2309.14316</i> , 2023.

# 756 A RELATED WORKS

758 Several studies independently identified the "attention sink" phenomenon in language models and 759 vision transformers, where attention weights were found to be concentrated on a few tokens (Xiao 760 et al., 2023; Darcet et al., 2023; Han et al., 2023; Zhai et al., 2023; Elhage et al., 2023; Dettmers 761 et al., 2022). Recent research has provided more detailed characterizations of this attention pattern and the attention sink phenomenon (Fu, 2024; Sun et al., 2024). Sun et al. (2024) attributed the 762 attention sink to the massive activation of the hidden representations of the corresponding tokens. Both Sun et al. (2024) and Zhai et al. (2023) discussed methods for mitigating the attention sink by 764 modifying the model and training recipes. Additionally, recent studies have leveraged the attention 765 sink phenomenon to develop improved quantization and more efficient inference algorithms (Liu 766 et al., 2024; Chen et al., 2024; Yu et al., 2024; Son et al., 2024). 767

768 The dynamics of transformers are studied under various simplifications, including linear attention structures (Zhang et al., 2023; Ahn et al., 2024), reparametrizations (Tian et al., 2023b), NTK (Deora 769 et al., 2023), often in the setting of in-context linear regressions (Ahn et al., 2023; Wu et al., 2023; 770 Zhang et al., 2024) and structured sequence (Bietti et al., 2024; Nichani et al., 2024; Tian et al., 771 2023a). Notably, Zhang et al. (2023) proves that a one-layer linear attention head trained with 772 gradient descent converges to a model that implements the in-context linear regression algorithm. 773 Huang et al. (2023); Kim et al. (2024) extend this to non-linear settings. Bietti et al. (2024) shows 774 the fast learning of bigram memorization and the slow development of in-context abilities. Tian et al. 775 (2023a) shows the scan and snap dynamics in reparametrized one-layer transformers. Reddy (2023) 776 simplifies the structure of the induction head, showing the connection between the sharp transitions 777 of in-context learning dynamics and the nested nonlinearities of multi-layer operations.

778 Mechanistic interpretability is a growing field focused on understanding the internal mechanisms of 779 language models in solving specific tasks (Elhage et al., 2021; Geva et al., 2023; Meng et al., 2022; Nanda et al., 2023; Olsson et al., 2022; Bietti et al., 2024; Wang et al., 2022; Feng & Steinhardt, 781 2023; Todd et al., 2023). This includes mechanisms like the induction head and function vector 782 for in-context learning (Elhage et al., 2021; Olsson et al., 2022; Todd et al., 2023; Bietti et al., 783 2024), the binding ID mechanism for binding tasks (Feng & Steinhardt, 2023), association-storage 784 mechanisms for factual identification tasks (Meng et al., 2022), and a complete circuit for indirect 785 object identification tasks (Wang et al., 2022). The task addressed in this paper is closely related to Bietti et al. (2024), which explored synthetic tasks where tokens are generated from either global 786 or context-specific bigram distributions. Several other studies have also used synthetic tasks to 787 investigate neural network mechanisms (Charton, 2022; Liu et al., 2022; Nanda et al., 2023; Allen-788 Zhu & Li, 2023; Zhu & Li, 2023; Guo et al., 2023; Zhang et al., 2022). 789

We note that Gurnee et al. (2024) proposed Attention Deactivation Neurons, a concept similar to
Dormant Attention Heads. Gurnee et al. (2024) hypothesized that when such a head attends to the
first token, it indicates that the head is deactivated and has minimal effect.

## B PROOFS

793 794

795

804

805

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

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

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

$$\boldsymbol{q}_{v} = \mathsf{SoftMax}(\mathrm{TF}([\langle \mathsf{s} \rangle; v_{1:n-1}; v])_{n}). \tag{3}$$

806 Similarly, define the full output probability matrix.

807 808 809  $\mathbf{Q} = \begin{pmatrix} q_{11} & \cdots & q_{1V} \\ \vdots & \ddots & \vdots \\ q_{V1} & \cdots & q_{VV} \end{pmatrix} = \begin{pmatrix} \boldsymbol{q}_1^\top \\ \vdots \\ \boldsymbol{q}_V^\top \end{pmatrix}.$ (4) Given any vector  $\boldsymbol{u} = [u_1; \ldots; u_d]$ , define the corresponding diagonal matrix as 

$$\operatorname{diag}(\boldsymbol{u}) = \begin{pmatrix} u_1 & 0 & \dots & 0\\ \vdots & \ddots & & \vdots\\ \vdots & & \ddots & \vdots\\ 0 & \dots & 0 & u_d \end{pmatrix}$$

Define 

$$\mathbf{G}_{v}^{\mathbf{Q}} = \operatorname{diag}(\boldsymbol{q}_{v}) - \boldsymbol{q}_{v}\boldsymbol{q}_{v}^{\top} \quad \mathbf{G}_{v}^{\mathbf{Q}} = \operatorname{diag}(\boldsymbol{p}_{v}) - \boldsymbol{p}_{v}\boldsymbol{p}_{v}^{\top}$$

Denote  $z = W \cdot \beta - W \circ \xi$ . We present a technical lemma.

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

821  
822 *Proof.* Since we have that 
$$\sum_{k=1}^{V} p_{vk} = 1$$
 and  $\sum_{k=1}^{V} q_{vk} = 1$  for any  $v$ ,  
823  
824  $(\mathbf{G}_{v}^{\mathbf{P}})_{ii} = p_{i} - p_{i}^{2} = p_{i}(\sum p_{k}) \ge \sum |(\mathbf{G}_{v}^{\mathbf{P}})_{ik}|$ 

$$(\mathbf{G}_{v}^{\mathbf{P}})_{ii} = p_{i} - p_{i}^{2} = p_{i}(\sum_{k \neq i} p_{k}) \ge \sum_{k \neq i} |(\mathbf{G}_{v}^{\mathbf{P}})_{ik}|$$
$$(\mathbf{G}_{v}^{\mathbf{Q}})_{ii} = q_{i} - q_{i}^{2} = q_{i}(\sum_{k \neq i} q_{k}) \ge \sum_{k \neq i} |(\mathbf{G}_{v}^{\mathbf{Q}})_{ik}|.$$

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 & Johnson (2012), they are positive semi-definite. 

#### B.1 PROOF OF THEOREM 2

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

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

Assumption A requires that  $d \ge VN$ . Given the fact that there are  $O(\exp(d))$  approximately linearly independent vectors for large d (Vershynin, 2018), it is possible to apply approximation theory to avoid Assumption A. However, since Assumption A pertains only to the construction of  $\lambda$  for trigger tokens and is unrelated to Theorem 3, we adopt it to simplify the proof of Theorem 2.

*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  $\mathbf{u}_i^\top ebd(v_j)$  for any  $v \in \mathcal{V}$  and  $i, j \in [N]$ . Adopting Assumption A, there exists a matrix Qry such that

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

Define the corresponding key matrix. 

Q

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

There exists a value matrix Val such that 

$$Val(ebd(v_i)) = 0 \quad \text{for } v_i \in \mathcal{T}, \quad i > 1,$$
  

$$Val(ebd(v_i)) = \xi_{v_i} \mathbf{u}_i \quad \text{for } v_i \in \mathcal{V} \setminus \mathcal{T}, \quad i > 0,$$
  

$$Val(ebd(\langle s \rangle)) = \boldsymbol{\beta}.$$
(7)

Further define the matrix M that satisfies 

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

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

When  $\min_{v \in \mathcal{V}} \alpha_v \to \infty$ ,  $\min_{v \in \mathcal{V}} \xi_v \to \infty$ ,  $\lambda \to \infty$ , and  $\beta = 0$ , if  $v_n \in \mathcal{T}$ , SoftMax $[TF([\langle s \rangle; v_{1:n-1}; v_n])] = \delta_{v_{n-1}}$ . If  $v_n \in \mathcal{V} \setminus \mathcal{T}$ , SoftMax $[TF([\langle s \rangle; v_{1:n-1}; v_n])] = p_{v_n}$ . All next-token probabilities match those in the data-generating procedure, aligning with the oracle algorithm. 

### B.2 THE STABLE PHASE IN THEOREM 3

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

Lemma 6. We have 

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

Furthermore,

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

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

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

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

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

and

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

$$= \left[\frac{e^{\alpha_i}}{e^{\alpha_i} + W}\right] [q_{ik} \mathbf{1}\{k = v\} - q_{ik} q_{iv}].$$

We can verify that

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

and

911  
912  
913  
913  
914  
915  
917  
918  

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

914  
915  
916  

$$= \left[\frac{e^{\alpha_i}}{e^{\alpha_i} + W}\right] [q_{iv} - q_{iv}]$$

$$= 0.$$

This finishes the proof of Lemma 6.

Proposition 7 computes the gradient of loss with respect to  $\alpha$  and  $\beta$ , giving the gradient flow. Proposition 7. The gradient flow of optimizing loss( $\alpha, \beta$ ) is given by

$$\dot{\alpha}_{v}(t) = \frac{\pi_{v}e^{\alpha_{v}}}{(e^{\alpha_{v}} + W)^{2}} \sum_{i=1}^{V} (p_{vi} - q_{vi})(W\beta_{i} - W_{i}\xi_{i}),$$
$$\dot{\alpha}_{v}(t) = \sum_{i=1}^{V} \left(\pi_{k}e^{\alpha_{k}}[p_{kv} - q_{kv}]\right)$$

$$\dot{\beta}_{v}(t) = \sum_{k=1} \left\{ \frac{\pi_{k} e^{-\kappa_{k}} [p_{kv} - q_{kv}]}{e^{\alpha_{k} + W}} \right\}.$$

*Proof.* The gradient flow gives that

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

Taking the derivative of  $loss(\alpha, \beta)$  gives that

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

Similarly, we have that

$$\frac{\partial \mathsf{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \beta_v} = \sum_{j=1}^V \pi_j \sum_{k=1}^V p_{jk} \Big\{ \frac{e^{\alpha_j} q_{jv}}{e^{\alpha_j} + W} - \frac{e^{\alpha_j} \mathbf{1}\{k = v\}}{e^{\alpha_j} + W} \Big\}$$
$$= \sum_{j=1}^V \Big\{ \frac{\pi_j e^{\alpha_j} [q_{jv} - p_{jv}]}{e^{\alpha_j} + W} \Big\}.$$

950 This proves Proposition 7.

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

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

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

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

$$= p_{vi}.$$

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

$$\frac{\partial \mathsf{loss}(\boldsymbol{\alpha},\boldsymbol{\beta})}{\partial \alpha_v}\Big|_{\boldsymbol{\alpha}^\star,\boldsymbol{\beta}^\star} = \frac{\pi_v e^{\alpha_v}}{(e^{\alpha_v}+W)^2} \sum_{k=1}^V \Big\{ (q_{vk}-p_{vk}) [W\beta_k - W_k \xi_k] \Big\} = 0,$$

970  
971 
$$\frac{\partial \mathsf{loss}(\boldsymbol{\alpha},\boldsymbol{\beta})}{\partial \beta_v}\Big|_{\boldsymbol{\alpha}^\star,\boldsymbol{\beta}^\star} = \sum_{k=1}^V \left\{ \frac{\pi_k e^{\alpha_k} [q_{kv} - p_{kv}]}{e^{\alpha_k} + W} \right\} = 0.$$

=

This shows that the given points are stationary points. We further compute the second-order deriva-tive using Lemma 6. 

$$\frac{\partial^2 \mathsf{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta})}{\partial \alpha_i \partial \alpha_v} \Big|_{\boldsymbol{\alpha}^\star, \boldsymbol{\beta}^\star} = \mathbf{1} \{ v = i \} \cdot \frac{\pi_v e^\alpha}{(e^\alpha + W)^2} \sum_{k=1}^V \Big\{ \frac{\partial q_{ik}}{\partial \alpha_v} [W \beta_k - W_k \xi_k] \Big\}$$

$$\mathbf{1}\{v=i\} \cdot \frac{-\pi_v e^{2\alpha}}{(e^{\alpha}+W)^4} \Big\{ \sum_{k=1}^V q_{ik} (e^{-\alpha}W + W_k)^2 \xi_k^2 - \Big[ \sum_{k=1}^V q_{ik} (e^{-\alpha}W + W_k) \xi_k \Big]^2 \Big\},$$

$$= \mathbf{1}\{v=i\} \cdot \frac{-\pi_v e^{2\alpha}}{(e^{\alpha}+W)^4} \Big\{ \sum_{k=1}^V p_{ik} (e^{-\alpha}W + W_k)^2 \xi_k^2 - \Big[ \sum_{k=1}^V p_{ik} (e^{-\alpha}W + W_k) \xi_k \Big]^2 \Big\}.$$

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

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

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

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

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

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

with 

$$\nabla_{\boldsymbol{\alpha}}^{2} \mathsf{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{e^{2\alpha}}{(e^{\alpha} + W)^{4}} \operatorname{diag} \left\{ \pi \circ [\mathbf{z}^{\top} \mathbf{G}_{1}^{\mathbf{P}} \mathbf{z}; \dots; \mathbf{G}_{V}^{\mathbf{P}} \mathbf{z}] \right\},\$$
$$\nabla_{\boldsymbol{\alpha}} \nabla_{\boldsymbol{\beta}} \mathsf{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{e^{2\alpha}}{(\alpha + W)^{3}} \operatorname{diag} \left\{ \pi \right\} [\mathbf{z}^{\top} \mathbf{G}_{1}^{\mathbf{P}}; \dots; \mathbf{z}^{\top} \mathbf{G}_{V}^{\mathbf{P}}],$$

$$\nabla_{\boldsymbol{\alpha}} \nabla_{\boldsymbol{\beta}} \mathsf{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{(e^{\alpha} + W)^3} \operatorname{diag} \left\{ \pi \right\} [\mathbf{z}^{\top} \mathbf{G}_1^{\mathbf{P}}; \dots; \mathbf{z}^{\top} \mathbf{G}_V^{\mathbf{P}}].$$

$$\nabla_{\boldsymbol{\beta}}^{2} \mathsf{loss}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{e^{2\alpha}}{(e^{\alpha} + W)^{2}} \sum_{k=1}^{V} \pi_{k} \mathbf{G}_{k}^{\mathbf{P}}.$$

At last, we diagonalize the Hessian matrix and get that

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

where the H is given by 

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

To prove that **H** is positive semi-definite, consider any vector  $\eta$  with  $\|\eta\|_2 = 1$ .  $oldsymbol{\eta}^{ op} \mathbf{H}oldsymbol{\eta} = \sum_{k=1}^V \pi_k \Big( oldsymbol{\eta}^{ op} \mathbf{G}_k^{\mathbf{P}} oldsymbol{\eta} - rac{oldsymbol{\eta}^{ op} \mathbf{G}_k^{\mathbf{P}} oldsymbol{z} oldsymbol{\pi}^{\mathbf{P}} oldsymbol{g}_k^{\mathbf{P}} oldsymbol{\eta}}{oldsymbol{z}^{ op} \mathbf{G}_k^{\mathbf{P}} oldsymbol{z}} \Big).$ 

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

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

1030 As a result, we have that

$$\boldsymbol{\eta}^{\top} \mathbf{H} \boldsymbol{\eta} \geq \sum_{k=1}^{V} \pi_{k} \Big( \boldsymbol{\eta}^{\top} \mathbf{G}_{k}^{\mathbf{P}} \boldsymbol{\eta} - \frac{\boldsymbol{z}^{\top} \mathbf{G}_{k}^{\mathbf{P}} \boldsymbol{z} \boldsymbol{\eta}^{\top} \mathbf{G}_{k}^{\mathbf{P}} \boldsymbol{\eta}}{\boldsymbol{z}^{\top} \mathbf{G}_{k}^{\mathbf{P}} \boldsymbol{z}} \Big) = 0.$$

<sup>1035</sup> This shows that **H** is positive semi-define. Therefore,  $\text{Hessian}(\text{loss}(\alpha^*, \beta^*))$  is positive semi-define. This proves Theorem 8.

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

1042 B.3 Attention sinks in Theorem 3

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

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

1050 Proof. We separately analyze each entry of  $\alpha$ . Focusing on  $\alpha_v$ , to simplify the notation, we introduce a random variable  $\varphi$  such that  $\mathbb{P}(\varphi = W_k \xi_k) = p_{vk}$ . Define

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

1053 Therefore, using Lemma 7, we get that

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

1057 We take in  $\beta = c$  and expand the expression of du/dt. This gives us 

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

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

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

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

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

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

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

1078 Therefore, we get two more inequalities 1079

$$\operatorname{Cov}(\theta_1(\varphi)\varphi,\varphi) \le C_{\theta}\operatorname{Var}(\varphi), \quad \operatorname{Cov}(\theta_2(\varphi)\varphi^2,\varphi) \ge 0.$$

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

$$\frac{\mathrm{d}u}{\mathrm{d}t} \leq \pi_v \operatorname{Cov}(e^{\frac{\varphi}{u+W}}, \varphi) \\
= \pi_v \operatorname{Cov}(1 + \frac{\theta_1(\varphi)\varphi}{u+W}, \varphi) \\
= \frac{\pi_v \operatorname{Cov}(1 + \frac{\theta_1(\varphi)\varphi}{u+W}, \varphi) \\
\leq \frac{\pi_v C_\theta \operatorname{Var}(\varphi)}{u}.$$

By solving the corresponding ODE, we get that 

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

To give a lower bound, we have that

$$\begin{array}{ll} \begin{array}{ll} 1094 \\ 1095 \\ 1096 \\ 1096 \\ 1097 \\ 1098 \\ 1098 \\ 1098 \\ 1098 \\ 1098 \\ 1098 \\ 1099 \\ 1100 \\ 1100 \\ 1100 \\ 1101 \\ 1102 \\ 1102 \\ 1103 \\ 1104 \\ 1104 \\ 1105 \end{array} \\ \begin{array}{ll} \begin{array}{ll} \displaystyle \frac{\mathrm{d} u}{\mathrm{d} t} \geq \frac{u(0)^2}{[u(0)+W]^2} \frac{\pi_v \operatorname{Cov}(e^{\frac{\omega}{u}+W},\varphi)}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \operatorname{Cov}(1+\frac{\varphi}{u+W}+\theta_2(\varphi)\left[\frac{\varphi}{u+W}\right]^2,\varphi) \\ \displaystyle \frac{\omega(0)^2}{[u(0)+W]^2} \frac{\pi_v}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \frac{\operatorname{Var}(\varphi)}{u+W} \\ \displaystyle \frac{\mathrm{Var}(\varphi)}{u} \\ \displaystyle \frac{\mathrm{d} u(0)^2}{[u(0)+W]^2} \frac{\pi_v}{\mathbb{E}e^{\frac{\varphi}{u(0)+W}}} \cdot \frac{u(0)}{u(0)+W} \cdot \frac{\operatorname{Var}(\varphi)}{u} \\ \displaystyle = \tilde{C}_{\theta} \frac{1}{u}. \end{array}$$

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

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

with  $r_v$  bounded. 

### B.4 VALUE STATE DRAINS IN THEOREM 3

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

$$\boldsymbol{\beta}(t) \rightarrow \boldsymbol{\beta}^{\star} = \overline{\boldsymbol{\beta}}(0)\mathbf{1} - e^{-\alpha} \boldsymbol{W} \circ \boldsymbol{\xi}$$

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

$$abla_{oldsymbol{eta}}^2 \mathsf{loss}(oldsymbol{lpha},oldsymbol{eta}) = \sum_{k=1}^V \pi_k \mathbf{G}_k^{\mathbf{Q}}.$$

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

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

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

$$\|\boldsymbol{\beta}(t) - \boldsymbol{\beta}^{\star}\|_2^2 \le \delta e^{-\mu t}.$$



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

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

1172 1173 C.1 ABLATIONS OF DIFFERENT MODEL STRUCTURES TRAINED ON THE BIGRAM-BACKCOPY TASK.

**Exploring the minimal structure for massive norms.** Figure 10 presents the difference of residual norms between the  $\langle s \rangle$  token and others ( $||\text{Res}_{\langle s \rangle}|| - \mathbb{E}_{v \neq \langle s \rangle}[||\text{Res}_{v}||]$ ), with different combinations of model structures. The 3 × TF and 2 × TF + mlp are two outliers, showing clear evidence of residual state peaks.

1179

1185

1180Attention plots, value state norms, and residual norms for a three-layer transformer trained on1181BB task.Figures 11, 12, and 13 show the extreme token phenomena in a three-layer transformer.1182The residual state peaks show different phenomena from those in LLMs, with the last layer output1183increasing the residual norms of non- $\langle s \rangle$  tokens. Figure 1 demonstrates that the residual state norms1184of  $\langle s \rangle$  drop match the magnitudes of other tokens at the last layer.

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



1236 In this section, we will identify more fine-grained static mechanisms for extreme-token phenomena 1237 in Llama 3.1-8B-Base. To do this, we identify circuits for the origin of attention sinks and small 1238 value states. Then, using ablation studies, we study the origin of massive norms. Again, we use the 1239 generic test phrase " $\langle s \rangle$  Summer is warm. Winter is cold."

- 1240
- **Attention sinks and global contextual semantics.** There are many attention sinks at layer 0, and the  $\langle s \rangle$  token is always the sink token (see Figure 20). From now on until the end of this section,





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  $\langle s \rangle$  token.

- <sup>1323</sup> F Assorted Caveats
- 1325 1326

1310

1311

1313

1315 1316 1317

1318

1319

### F.1 MULTIPLE ATTENTION SINKS VS. ONE ATTENTION SINK

1327 As we have seen, attention heads in the BB task (Section 2), Llama 2-7B-Base (Section 3.1), and 1328 OLMo (Section 3.2) exhibit multiple attention sinks. That is, when heads in these models are dormant, they tend to have two attention sinks. For the LLMs in this group, at least on prose data, the 1330  $\langle s \rangle$  token as well as the first delimiter token (e.g., representing . or ;) are sink tokens. Meanwhile, 1331 Llama-3.1-8B-Base (Section 3) only ever has one attention sink on prose data, and the  $\langle s \rangle$  token is always the sink token. Here, we offer a possible explanation of this phenomenon. For the BB 1332 task, multiple sink tokens are necessary to solve the task. For LLMs, we believe this distinction may 1333 be explained by the relative proportion of coding data, in which delimiters have a greater semantic 1334 meaning than prose, within the training set. For instance, OLMo was trained on DOLMA (Sol-1335 daini et al., 2024), which has around 411B coding tokens. Meanwhile, Llama 2 used at most ( $2T \times$ 1336 0.08 = 0.16T coding tokens. Finally, Llama 3.1 used around  $(15.6T \times 0.17 =) 2.6T$  coding tokens 1337 (Dubey et al., 2024). On top of the raw count being larger, coding tokens are a larger proportion of 1338 the whole pre-training dataset for Llama 3.1 compared to other model families. Thus, during train-1339 ing, the presence of delimiters would not be considered unhelpful towards next-token prediction, 1340 since such delimiters carry plenty of semantics in a wide variety of cases. Our earlier hypothesis in 1341 Section 3.1 proposes that only tokens which lack semantics in almost all cases are made to be sink tokens. This could be a reason for the distinction. 1342

- 1343
- 1344 1345

F.2 The Role of a Fixed  $\langle s \rangle$  Token in the Active-Dormant Mechanism

Some models, such as OLMo, are not trained with a  $\langle s \rangle$  token. Despite this, the first token of the input still frequently develops into a sink token. We can study the effect of positional encoding of the tokens on the attention sink phenomenon by shuffling the tokens before inputting them into the transformer, and observing how and why attention sinks form. If we do this with the phrase "Summer is warm. Winter is cold." with OLMo, we observe that at Layer 24, there are many



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



1454 1455 Figure 21: Alignment between query states and key states at Layer 0 Head 31 of Liama 3.1-bB-base. we 1456 observe that the key state of  $\langle s \rangle$  is orthogonal to all other key states, and heavily aligned with all query states. 1456 Meanwhile, all semantically meaningful (i.e., not delimiter) tokens have aligned key states.



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  $\langle s \rangle$  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.