

A Mechanistic Account of Attention Sinks in GPT-2: One Circuit, Broader Implications for Mitigation

Yuval Ran-Milo* and Hila Ofek* and Shahar Mendel

Tel Aviv University

{yuvalmilo, hilaofek1, shaharmendel}@mail.tau.ac.il

Abstract

Transformers commonly exhibit an attention sink: disproportionately high attention to the first position. We study this behavior in GPT-2-style models with learned query biases and absolute positional embeddings. Combining structural analysis with causal interventions, validated across natural-language, mathematical, and code inputs, we find that the sink arises from the interaction among (i) a learned query bias, (ii) the first-layer MLP transformation of the positional encoding, and (iii) structure in the key projection. Crucially, each component we identify is individually dispensable: architectures omitting each of them robustly exhibit sinks. This indicates that attention sinks may arise through distinct circuits across architectures. These findings inform mitigation of sinks, and motivate broader investigation into why sinks emerge.

1 Introduction

Transformers (Vaswani et al., 2017) routinely display an *attention sink*: a persistent tendency to allocate disproportionate attention mass to early (often first) positions independent of semantic content (Xiao et al., 2024; Gu et al., 2025). The effect has been observed across training stages and hyperparameters (Gu et al., 2025; Guo et al., 2024), across model families and datasets (Xiao et al., 2024), and under diverse positional encodings—including ALiBi (Press et al., 2022), RoPE (Su et al., 2023), and even no explicit positional encodings (Gu et al., 2025). Similar sink-like patterns have also been reported in large multimodal models and vision transformers (Kang et al., 2025; Wang et al., 2025; Feng and Sun, 2025).¹

The practical stakes are significant. Attention sinks can reduce effective context use and lower accuracy (Yu et al., 2024a; Guo et al., 2024), aggravate numerical error and hinder quantization (Sun et al., 2024; Lin et al., 2024), obscure interpretability by dominating attention maps (Guo et al., 2024),

and complicate streaming and KV-cache strategies (Xiao et al., 2024). Analogous effects in vision and multimodal settings waste representational capacity on irrelevant tokens and can be exploited to induce hallucinations (Kang et al., 2025; Wang et al., 2025; Feng and Sun, 2025). Understanding when sinks arise and how to control them is therefore directly relevant for model performance and interpretability.

We study the sink mechanistically in GPT-2-style Transformers with learned query biases and absolute positional embeddings (Radford et al., 2019). Combining structural analysis with targeted causal interventions validated across natural-language, mathematical, and code inputs, we tie the first-token sink to three interacting components: learned query bias, first-layer MLP transformation of the positional encoding, and structure in the key projection.² We establish causality by showing that disrupting any component weakens, removes, or relocates the sink. When alternative explanations exist, we ablate them, isolating the causal circuit.

A direct consequence of our analysis is that sinks must arise through distinct mechanisms across architectures: each component we identify is individually dispensable in models that still exhibit sinks. This implies that while attention sinks are robust as a phenomenon, they are not governed by a single universal mechanism. Consequently, many post-training mitigations may need to be architecture-specific, and many pre-training architectural interventions may not generalize (see section 4). This highlights the need to understand why the sink emerges as a prerequisite for developing robust pre-training strategies.

2 Preliminaries

2.1 Attention mechanism

We use D for model dimension, H for number of heads, D_h for per-head dimension, and L for number of layers. We index layers by l , heads by h , and sequence positions by i or j . The representation at layer l , position i is

*Equal contribution.

¹However, some architectures have been reported to have little to no sink (Qiu et al., 2025; Endy et al., 2025).

²Our findings connect to and extend a body of work on the origins of attention sinks; see Appendix B for a discussion.

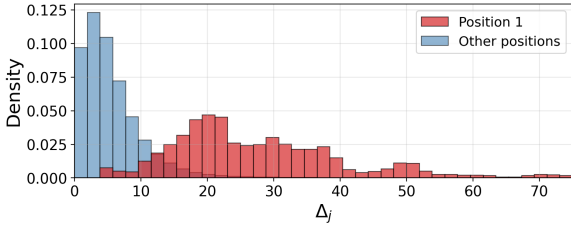


Figure 1: **The source-agnostic shift $\Delta_{j,h}^{(l)}$ at position 1 is systematically larger than all other positions.** Two overlapping density histograms of the source-agnostic shift $\Delta_{j,h}^{(l)} = b_{Q,h}^{(l)\top} W_{k,h}^{(l)\top} x_j^{(l)\top}$ (normalized so $\min_j \Delta_{j,h}^{(l)} = 0$; see section 3.2.2), pooled across all heads h and layers $l \in [4, 11]$ (see section 3.2.1). The distribution of scores for the first position ($j = 1$, red) is clearly separated from and systematically larger than the distribution for all other positions ($j > 1$, blue). The x -axis is truncated for clarity; see full histogram in Appendix A.1.

$x_i^{(l)} \in \mathbb{R}^D$. Within each head h , queries, keys, and values are $q_{i,h}^{(l)} = x_i^{(l)} W_{q,h}^{(l)} + b_{Q,h}^{(l)}$, $k_{j,h}^{(l)} = x_j^{(l)} W_{k,h}^{(l)} + b_{K,h}^{(l)}$, and $v_{j,h}^{(l)} = x_j^{(l)} W_{v,h}^{(l)} + b_{V,h}^{(l)}$, where $W_{q,h}^{(l)}, W_{k,h}^{(l)}, W_{v,h}^{(l)} \in \mathbb{R}^{D \times D_h}$ and biases $b_{Q,h}^{(l)}, b_{K,h}^{(l)}, b_{V,h}^{(l)} \in \mathbb{R}^{D_h}$. The pre-softmax attention score from position i to position j in head h is $s_{i \rightarrow j,h}^{(l)} = q_{i,h}^{(l)} (k_{j,h}^{(l)})^\top$. The attention weights are $\alpha_{i \rightarrow j,h}^{(l)} = \text{softmax}_j(s_{i \rightarrow j,h}^{(l)} / \sqrt{D_h})$ over valid positions $j \leq i$, and the output for position i is $o_{i,h}^{(l)} = \sum_j \alpha_{i \rightarrow j,h}^{(l)} v_{j,h}^{(l)}$; head outputs are concatenated and projected by $W_o^{(l)} \in \mathbb{R}^{D \times D}$.³

2.2 Positional encoding

Attention layers are not inherently order-aware; apart from masking-induced structure, they are invariant to input permutations. To address this, Transformers incorporate positional information through various schemes (Su et al. (2023), Press et al. (2022)). GPT-2 uses learned absolute positional encodings: a set of trainable vectors $\{p_i\}_{i=1}^N \subset \mathbb{R}^D$, where i is the token position and N is the maximal sequence length. These are added to token embeddings e_i : $x_i^{(0)} = e_i + p_i$.

2.2.1 Effective positional encoding (EPE)

We define the *effective positional encoding* (EPE) for position i as $\text{EPE}_i = \text{MLP}^{(1)}(p_i) + p_i \in \mathbb{R}^D$, where $\text{MLP}^{(1)}$ is the first layer’s feed-forward network. We call this “effective” because it roughly captures the net positional signal that the first layer’s MLP writes into the residual stream⁴: as

³We adopt the usual transformer convention $D = HD_h$.

⁴This equals *exactly* the positional contribution if we simplify the first layer to only a linear MLP (removing its nonlinearity), omitting attention and removing normalization.

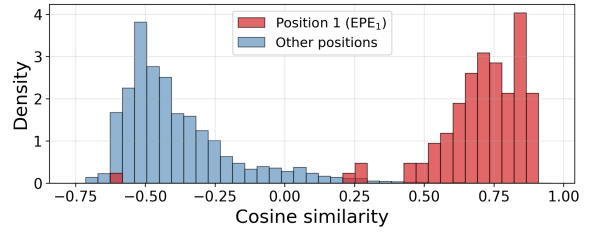


Figure 2: **The first-position key projection $\text{EPE}_1 W_{k,h}^{(l)}$ is uniquely aligned with the query bias.** Two overlapping density histograms of the per-head cosine similarity between $b_{Q,h}^{(l)}$ and $\text{EPE}_j W_{k,h}^{(l)}$, pooled across all heads h and layers $l \in [4, 11]$. The first-position distribution ($j=1$, red) is concentrated at high positive values, indicating strong alignment with the query bias, while the distribution for all other positions ($j>1$, blue) is negatively aligned.

we verify in section 3.2.3, adding EPE_i to the first layer’s output produces roughly the same result as adding p_i before the first layer and then applying the first layer’s MLP.

3 Methodology and Results

We first describe the mechanism underlying attention sinks in models with learned query biases and absolute positional encodings in section 3.1. We then provide empirical evidence through four analyses in sections 3.2.2 to 3.2.5 and confirm causality through targeted interventions in section 3.2.6.

3.1 The Mechanism behind Attention Sinks

Expanding the attention score $s_{i \rightarrow j,h}^{(l)}$ (layer l , head h , source position i , target position j) from section 2.1 by substituting the expressions for $q_{i,h}^{(l)}$ and $k_{j,h}^{(l)}$ gives $s_{i \rightarrow j,h}^{(l)} = (x_i^{(l)} W_{q,h}^{(l)}) (x_j^{(l)} W_{k,h}^{(l)})^\top + (x_i^{(l)} W_{q,h}^{(l)}) b_{K,h}^{(l)\top} + b_{Q,h}^{(l)} (x_j^{(l)} W_{k,h}^{(l)})^\top + b_{Q,h}^{(l)} b_{K,h}^{(l)\top}$. The third term, $\Delta_{j,h}^{(l)} \triangleq b_{Q,h}^{(l)} W_{k,h}^{(l)\top} x_j^{(l)\top}$, is a token-specific, *source-agnostic* shift: it raises or lowers the score toward target position j identically for every source position i . We find that this term for the first token, $\Delta_{1,h}^{(l)}$, is conspicuously large across heads in most deep layers, creating a strong prior to attend to position 1.

To understand why $\Delta_{1,h}^{(l)}$ is so large, recall that the effect of adding p_1 to the first token and passing it through the first layer is roughly equivalent to adding EPE_1 directly to the residual stream at position 1 (as verified in section 3.2.3). EPE_1 has very large absolute values on a small set of coordinates—the phenomenon of *massive activations* (Sun et al., 2024)—which are exactly those coordinates where $b_{Q,h}^{(l)} W_{k,h}^{(l)\top}$ has the largest magnitude in almost all layers. This co-adaptation allows EPE_1 to dramatically amplify $\Delta_{1,h}^{(l)}$ across heads, yielding an

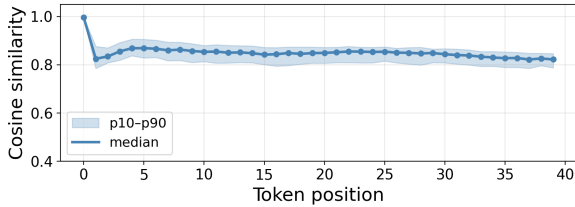


Figure 3: **EPE closely tracks the net positional signal at every position.** Per-position cosine similarity between EPE_i and the net added positional signal N_i (see section 3.2.3). The shaded band spans the 10th–90th percentile across our dataset; the line marks the median. The similarity exceeds 0.82 at every position and reaches > 0.99 at position 1, confirming that EPE_1 (the key driver of the sink) is an accurate proxy for the actual positional contribution.

attention sink at position 1.

3.2 Empirical Validation

We validate our proposed mechanism through four complementary analyses, followed by causal interventions that confirm the necessity of each component described in section 3.1. The model and evaluation dataset are described in section 3.2.1. In section 3.2.2 we show that $\Delta_{1,h}^{(l)}$ is conspicuously large relative to other positions across layers and heads. In section 3.2.3 we verify that the EPE accurately captures the net positional signal written by the first layer’s MLP. We then investigate the underlying cause of the large shift and show in section 3.2.4 that $\text{EPE}_1 W_{k,h}^{(l)}$ is strongly aligned with $b_{Q,h}^{(l)}$ across heads and layers. In section 3.2.5 we establish that EPE_1 exhibits massive activations precisely at coordinates where the bias projection $b_{Q,h}^{(l)} W_{k,h}^{(l)\top}$ has high magnitude. Finally, in section 3.2.6 we use causal interventions to verify that disrupting any component weakens or removes the sink while transplanting it transfers the sink to a new position.

3.2.1 Setup

All experiments are conducted on the 124M-parameter GPT-2 (Radford et al., 2019)⁵. We evaluate on a dataset of 300 examples: 100 randomly sampled from each of SST-2 (Socher et al., 2013) (natural language), GSM8K (Cobbe et al., 2021) (mathematical reasoning), and HumanEval (Chen et al., 2021) (code generation). Only examples with at least 40 tokens are kept; all selected examples are then truncated to exactly 40 tokens so that every input has the same sequence length. This standardization ensures that no dataset disproportionately influences aggregate statistics due to varying input

⁵Pretrained checkpoint from the Hugging Face Hub: <https://huggingface.co/openai-community/gpt2>; License: MIT.

lengths; we find that results are qualitatively equivalent for other length choices. All per-head analyses are reported over layers 4–11 (of 12), excluding the first three layers, and the last layer, where attention sinks are known to be weaker (Yu et al., 2024b). Code and data are available at <https://github.com/YuvMilo/MechanisticAccountofSinks>.

3.2.2 Source-Agnostic Shift Analysis

First, we verify that $\Delta_{1,h}^{(l)}$ is anomalously large relative to other positions. We compute $\Delta_{j,h}^{(l)}$ for every head h and position j across our evaluation dataset. Because softmax attention is invariant to adding a constant across all target positions, we re-center Δ by subtracting $\min_j \Delta_{j,h}^{(l)}$ from all positions for each (sentence, layer, head) slice, setting the smallest value to zero; this does not affect the model’s output and makes magnitudes comparable across slices. Figure 1 shows two overlapping density histograms of $\Delta_{j,h}^{(l)}$ across all heads h and layers $l \in [4, 11]$: one for the first position ($j=1$, red) and one for all other positions ($j>1$, blue). The two distributions are clearly separated, with the first-position distribution shifted far to the right, confirming that $\Delta_{1,h}^{(l)}$ is indeed anomalously large and creates a strong, consistent prior to attend to position 1. Later, in section 3.2.6, we show that nullifying b_Q alone—effectively zeroing out this term in every attention score—is sufficient to drastically diminish the sink.

3.2.3 EPE Captures the Net Positional Contribution

Our mechanism relies on EPE_i faithfully representing the net positional signal written into the residual stream by the first layer’s MLP. To verify this, we define the *net added positional signal* $N_i := R_i - O_i$, where $R_i = x_i^{(0)} + \text{MLP}^{(1)}(x_i^{(0)})$ is the MLP output given an input with positional encoding, and $O_i = e_i + \text{MLP}^{(1)}(e_i)$ is the output for the same token embedding without positional encoding ($e_i = x_i^{(0)} - p_i$). N_i is thus the incremental change to the residual stream attributable to the positional encoding. We then measure the cosine similarity between EPE_i and N_i across our dataset (see section 3.2.1) and all positions. Figure 3 shows that the median similarity exceeds 0.82 at every position. Crucially, for position 1 (the sink position), the cosine similarity exceeds 0.99.

Furthermore, we conducted an additional experiment comparing the cosine similarity of EPE_1 and the added positional signal of the entire first layer (including normalization and first layer attention) in a similar fashion, and found the cosine similarity median to be 0.76 (90th and 10th percentiles are

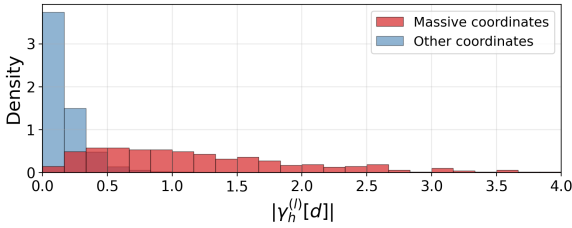


Figure 4: EPE_1 is large exactly where the bias projection is large. Two overlapping density histograms of $|\gamma_h^{(l)}[d]|$, across all heads h and layers $l \in [4, 11]$: massive-activation coordinates of EPE_1 (red) versus all other coordinates (blue). The x -axis is truncated for clarity; see full histogram in Appendix A.3.

0.8 and 0.67 respectively).

3.2.4 EPE-Bias Projection Alignment

Having established the magnitude of $\Delta_{1,h}^{(l)}$, we investigate its underlying cause. Since $x_1^{(l)}$ contains both token and positional information, it remains to disentangle which of the two is responsible for the large $\Delta_{1,h}^{(l)}$. To that end, we compute the cosine similarity between $b_{Q,h}^{(l)}$ and $EPE_j W_{k,h}^{(l)}$ for every head h , position j , and layer l . Figure 2 shows two overlapping density histograms pooled over layers 4–11 and all heads: the first-position distribution ($j=1$, red) is tightly concentrated at high positive values, indicating that the first-position key projection $EPE_1 W_{k,h}^{(l)}$ is strongly aligned with the query bias, while the distribution for all other positions ($j>1$, blue) is negatively aligned. Since EPE_1 is the only position positively aligned with the query bias, its addition to the residual stream (see section 3.2.3) directly enlarges $\Delta_{1,h}^{(l)}$, attributing the outsized shift to the positional signal.

3.2.5 Coordinate-Level Structural Analysis

Since $\Delta_{j,h}^{(l)} = b_{Q,h}^{(l)} W_{k,h}^{(l)\top} x_j^{(l)\top}$, the large $\Delta_{1,h}^{(l)}$ requires $b_{Q,h}^{(l)} W_{k,h}^{(l)\top}$ to have high magnitude at the same coordinates where EPE_1 is large. We verify this by comparing, for each head h , the bias projection magnitude $\gamma_h^{(l)}[d] := |b_{Q,h}^{(l)} \cdot W_{k,h}^{(l)}[:, d]|$ at the massive-activation coordinates of EPE_1 (identified in Appendix A.2) against all other coordinates. Figure 4 shows two overlapping density histograms of $\gamma_h^{(l)}[d]$ across all heads h and layers $l \in [4, 11]$: massive-activation coordinates (red) are systematically larger than all other coordinates (blue), confirming the coordinate-level co-adaptation that amplifies $\Delta_{1,h}^{(l)}$.

3.2.6 Causal Interventions

To establish causality beyond correlation, we perform ten targeted interventions during forward

	Intervention	BOS Attn	% Base
(a)	Baseline	.563 \pm .004	100.0
(b)	Nullify b_Q	.251 \pm .003	44.7
(c)	Remove First PE	.017 \pm .001	3.0
(d)	Swap EPE	.035 \pm .002	6.2
(e)	Swap PE	.557 \pm .006	99.0
(f)	Nullify BOS Token	.561 \pm .004	99.7
(g)	No MLP	.283 \pm .009	50.3
(h)	No PE	.099 \pm .015	17.5
(i)	Zero Top-3 W_k	.367 \pm .004	65.2
(j)	Zero Random W_k	.563 \pm .004	100.0

Table 1: **Each component of the sink circuit is necessary.** BOS-attention metric (mean $f \pm 2$ SE) across our dataset, averaged over all heads and layers 4–11. “% Base” is the fraction of the baseline (a) score retained. Disrupting any elements of the b_Q – EPE_1 – W_k pathway (b–d, g–i) substantially reduces the sink; controls (e, f, j) do not.

passes, testing the necessity of components by removing them and sufficiency by transplanting them. Where alternative explanations exist, we include control interventions to ablate them. We quantify each intervention’s effect using the **BOS-attention metric**: the average attention weight assigned to position 1 by tokens in the second half of the sequence,⁶ averaged over all heads and layers 4–11 (see section 3.2.1). Results are in table 1. Figure 5 illustrates the effect of each intervention on the attention map for a single representative sentence.⁷

- **Baseline.** No intervention; the attention sink is clearly present.
- **Nullify b_Q .** We set $b_Q=0$ in every layer; the sink substantially diminishes, showing that b_Q is necessary for the large first-token contribution.
- **Remove First PE.** We replace PE_1 with PE_2 at position 1; the sink nearly vanishes. Together with the BOS-nullification control (f) below, this shows that the positional encoding, not the BOS token identity, drives the sink.
- **Swap EPE.** We perform a swap of the EPE component between positions 1 and 2⁸ in the first layer’s MLP output (technical details in section B.1). The sink nearly vanishes, demonstrating that the MLP-transformed positional signal controls sink location.
- **Swap PE.** Same swap as above but using the raw positional embeddings instead of the EPE (details in section B.1). The sink persists, showing that downstream layers rely specifically on the MLP-processed signal, not the raw positional embedding.

⁶We restrict to the second half so that early positions, whose attention is concentrated on position 1 simply because few alternatives are available, do not inflate the metric.

⁷“It was the best of times, it was the worst of times, it was the age of wisdom.” (from *A Tale of Two Cities* by Charles Dickens, 1859).

⁸Results are qualitatively the same when swapping position 1 and any other position, not necessarily 2.

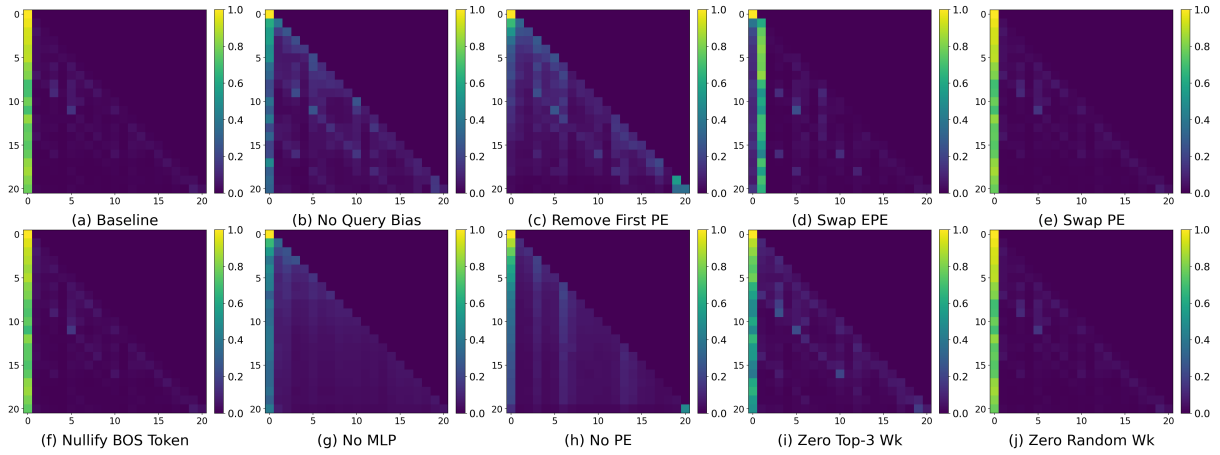


Figure 5: **Disrupting any component of the b_Q -EPE₁- W_k pathway significantly diminishes the sink.** Head-averaged attention maps (layers 4–11) under all ten interventions for a single representative sentence. Each panel shows attention weights from query positions (rows) to key positions (columns). The first-position sink (bright first column) persists when the pathway is left intact (a, e, f, j), but is significantly diminished when any component is disrupted (b, c, d, g, h, i). See table 1 for aggregate statistics across our dataset.

- **Nullify BOS Token.** Zeroing the BOS token embedding before adding positional signals leaves the sink intact, ruling out BOS token identity as a driver.
- **No MLP.** Skipping all MLP blocks substantially reduces the sink,⁹ consistent with the first-layer MLP’s role in amplifying the positional signal into the EPE.
- **No PE.** Removing all positional embeddings greatly diminishes the sink, confirming that positional information is the primary driver.
- **Zero Top-3 W_k .** Zeroing the three W_k columns at the massive-activation coordinates of EPE₁ reduces the sink, confirming that those coordinates are the mechanism through which EPE₁ amplifies $\Delta_{1,h}^{(l)}$.
- **Zero Random W_k .** Zeroing three random W_k columns has nearly no effect, confirming that the result above is specific to the massive-activation coordinates of EPE₁.

4 Conclusions

Through fine-grained structural analysis and targeted causal interventions, validated across natural-language, mathematical, and code inputs, we identified the circuit behind attention sinks in GPT-2: the interaction of a learned query bias (b_Q), the first-layer MLP’s transformation of the positional encoding, and coordinate-level structure in the key projection (W_k). Disrupting any single component significantly weakens the sink, while alternative explanations are ablated.

⁹We disable all MLPs rather than only the first because we observe that later MLPs also process the positional encoding in a qualitatively similar manner, though more weakly.

Crucially, each component of this circuit is individually dispensable: many architectures omit b_Q entirely yet still display strong attention sinks (Gu et al., 2025; Xiao et al., 2024), sinks persist even in Transformers without any MLP blocks (Hong and Lee, 2025), and sinks persist in Transformers not using any positional encodings (Gu et al., 2025). The sink therefore arises through distinct circuits across architectures¹⁰—suggesting it reflects a structural computational necessity that optimization reliably learns using whatever architectural components are available.¹¹

This implies that while attention sinks are robust as a phenomenon, they are not governed by a single universal circuit. Consequently, many pre-training architectural modifications aimed at preventing sinks may prove insufficient, as optimization can potentially reconstruct an equivalent sink using alternative components. Likewise, many post-training mitigations may need to be architecture-specific: targeting a particular component will succeed only if it forms part of the active sink-implementing pathway in that specific model. We hope our causal analysis serves as a first step toward designing architecture-aware mitigations, and motivates further investigation into why the sink emerges as a prerequisite for robust pre-training strategies.

¹⁰Other architectures may construct similar mechanisms using different tools (e.g., a constant-feature slice of W_Q could substitute for an absent b_Q). Our claim is that the *specific* parameter-level pathway we identify does not carry over.

¹¹This view is supported by works showing that sinks are provably necessary for certain attention computations (Ran-Milo, 2026), and that they serve a functional role in mitigating over-mixing (Barbero et al., 2025).

5 Limitations

5.1 Scope across architectures and scales

Our analyses focus on a GPT-2–style model with learned query biases and absolute positional encodings. The broader Transformer ecosystem includes architectures that omit such biases or use alternative positional schemes (e.g., RoPE, ALiBi). While we do show that the specific components constituting the GPT-2 circuit are absent or altered in those settings, we leave the investigation of what mechanisms do form in those models to future work. In addition, GPT-2 is small by contemporary standards; with scale, the mechanism could strengthen, fragment into multiple pathways, or be replaced by different circuits.

5.2 Learning dynamics

We provide a post-hoc, static analysis of a trained checkpoint. We do not track when the circuit emerges during pre-training, which gradients give rise to it, or whether intermediate snapshots exhibit qualitatively different pathways. We believe our static analysis could inform future work researching the emergence of the attention sink mechanism.

5.3 Secondary contributors

Our interventions reduce the sink substantially but do not eliminate it entirely: nullifying b_Q or zeroing the massive-activation rows of W_k leaves a residual sink. This indicates that secondary contributors beyond the three main components we isolate also play a role. Identifying them is a natural extension of our work which we leave to future work.

5.4 Mechanism vs. function

Our contribution is mechanistic: we explain *how* an attention sink can be implemented in the studied architecture. We do not claim a definitive functional rationale for *why* such a sink is beneficial or harmful across tasks. Establishing the downstream utility or cost of the sink, and the conditions under which it is selected by optimization, is left for future work.

Acknowledgments

I thank Amit Elhelo and Daniela Gottesman for illuminating discussions. Special thanks to my advisor Nadav Cohen for his guidance and mentorship. We used AI assistance for writing and code development. This work was supported by the European Research Council (ERC) grant NN4C 101164614, a Google Research Scholar Award, a Google Research Gift, Meta, the Yandex Initiative in Machine Learning, the Israel Science Foundation (ISF) grant 1780/21, the Tel Aviv University Center for AI and

Data Science, the Adelis Research Fund for Artificial Intelligence, Len Blavatnik and the Blavatnik Family Foundation, and Amnon and Anat Shashua.

References

- Federico Barbero, Álvaro Arroyo, Xiangming Gu, Christos Perivolaropoulos, Michael Bronstein, Petar Veličković, and Razvan Pascanu. 2025. [Why do llms attend to the first token?](#) *Preprint*, arXiv:2504.02732.
- Nicola Cancedda. 2024. [Spectral filters, dark signals, and attention sinks.](#) *Preprint*, arXiv:2402.09221.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, and 39 others. 2021. [Evaluating large language models trained on code.](#) *Preprint*, arXiv:2107.03374.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems.](#) *Preprint*, arXiv:2110.14168.
- Enrique Queipo de Llano, Álvaro Arroyo, Federico Barbero, Xiaowen Dong, Michael Bronstein, Yann LeCun, and Ravid Shwartz-Ziv. 2026. [Attention sinks and compression valleys in llms are two sides of the same coin.](#) *Preprint*, arXiv:2510.06477.
- Jack Dial. 2025. [The curious case of the BOS token.](#)
- Nir Endy, Idan Daniel Grosbard, Yuval Ran-Milo, Yonatan Slutzky, Itay Tshuva, and Raja Giryes. 2025. [Mamba knockout for unraveling factual information flow.](#) *Preprint*, arXiv:2505.24244.
- Wenfeng Feng and Guoying Sun. 2025. [Edit: Enhancing vision transformers by mitigating attention sink through an encoder-decoder architecture.](#) *Preprint*, arXiv:2504.06738.
- Xiangming Gu, Tianyu Pang, Chao Du, Qian Liu, Fengzhuo Zhang, Cunxiao Du, Ye Wang, and Min Lin. 2025. [When attention sink emerges in language models: An empirical view.](#) *Preprint*, arXiv:2410.10781.
- Tianyu Guo, Druv Pai, Yu Bai, Jiantao Jiao, Michael I. Jordan, and Song Mei. 2024. [Active-dormant attention heads: Mechanistically demystifying extreme-token phenomena in llms.](#) *Preprint*, arXiv:2410.13835.
- Jonghyun Hong and Sungyoon Lee. 2025. [Variance sensitivity induces attention entropy collapse and instability in transformers.](#) In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 8360–8378, Suzhou, China. Association for Computational Linguistics.

Seil Kang, Jinyeong Kim, Junhyeok Kim, and Seong Jae Hwang. 2025. [See what you are told: Visual attention sink in large multimodal models](#). *ArXiv*, abs/2503.03321.

Haokun Lin, Haobo Xu, Yichen Wu, Jingzhi Cui, Yingtao Zhang, Linzhan Mou, Linqi Song, Zhenan Sun, and Ying Wei. 2024. [Duquant: Distributing outliers via dual transformation makes stronger quantized llms](#). *Preprint*, arXiv:2406.01721.

Ofir Press, Noah A. Smith, and Mike Lewis. 2022. [Train short, test long: Attention with linear biases enables input length extrapolation](#). *Preprint*, arXiv:2108.12409.

Zihan Qiu, Zeyu Huang, Kaiyue Wen, Peng Jin, Bo Zheng, Yuxin Zhou, Haofeng Huang, Zekun Wang, Xiao Li, Huaqing Zhang, Yang Xu, Haoran Lian, Siqi Zhang, Rui Men, Jianwei Zhang, Ivan Titov, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2026. [A unified view of attention and residual sinks: Outlier-driven rescaling is essential for transformer training](#). *Preprint*, arXiv:2601.22966.

Zihan Qiu, Zekun Wang, Bo Zheng, Zeyu Huang, Kaiyue Wen, Songlin Yang, Rui Men, Le Yu, Fei Huang, Suozhi Huang, Dayiheng Liu, Jingren Zhou, and Junyang Lin. 2025. [Gated attention for large language models: Non-linearity, sparsity, and attention-sink-free](#). *Preprint*, arXiv:2505.06708.

Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#). Technical report, OpenAI.

Yuval Ran-Milo. 2026. [Attention sinks are provably necessary in softmax transformers: Evidence from trigger-conditional tasks](#). *Preprint*, arXiv:2603.11487.

Valeria Ruscio, Umberto Nanni, and Fabrizio Silvestri. 2025. [What are you sinking? a geometric approach on attention sink](#). *Preprint*, arXiv:2508.02546.

Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. [Recursive deep models for semantic compositionality over a sentiment treebank](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA. Association for Computational Linguistics.

Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. 2023. [Roformer: Enhanced transformer with rotary position embedding](#). *Preprint*, arXiv:2104.09864.

Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. 2024. [Massive activations in large language models](#). *Preprint*, arXiv:2402.17762.

Shangwen Sun, Alfredo Canziani, Yann LeCun, and Jiachen Zhu. 2026. [The spike, the sparse and the sink: Anatomy of massive activations and attention sinks](#). *Preprint*, arXiv:2603.05498.

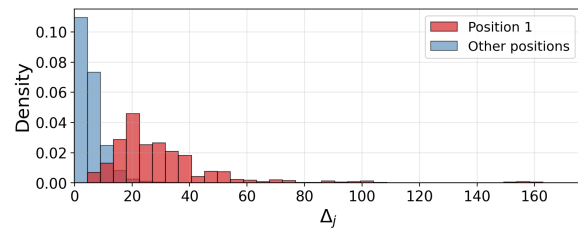


Figure 6: **Full (non-truncated) histogram of the source-agnostic shift $\Delta_{j,h}^{(l)}$** . Same data as fig. 1 with no axis restriction. Scores are normalized per (sentence, layer, head) slice so the minimum is zero. The first-position distribution (red) extends to a long right tail that lies well beyond the range of all other positions (blue).

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). *Preprint*, arXiv:1706.03762.

Yining Wang, Mi Zhang, Junjie Sun, Chenyue Wang, Min Yang, Hui Xue, Jialing Tao, Ranjie Duan, and Jiexi Liu. 2025. [Mirage in the eyes: Hallucination attack on multi-modal large language models with only attention sink](#). *Preprint*, arXiv:2501.15269.

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. [Efficient streaming language models with attention sinks](#). *Preprint*, arXiv:2309.17453.

Zhongzhi Yu, Zheng Wang, Yonggan Fu, Huihong Shi, Khalid Shaikh, and Yingyan Celine Lin. 2024a. [Unveiling and harnessing hidden attention sinks: Enhancing large language models without training through attention calibration](#). *Preprint*, arXiv:2406.15765.

Zhongzhi Yu, Zheng Wang, Yonggan Fu, Huihong Shi, Khalid Shaikh, and Yingyan Celine Lin. 2024b. [Unveiling and harnessing hidden attention sinks: Enhancing large language models without training through attention calibration](#). *Preprint*, arXiv:2406.15765.

A Further Experiments

A.1 Full Source-Agnostic Shift Histogram

Figure 6 shows the full (non-truncated) version of fig. 1.

A.2 Identifying Massive Activations in First-Position EPE

This section explains how we identify coordinates with unusually large absolute values in EPE_1 . We select coordinates whose absolute values exceed the mean absolute value by at least three standard deviations; in our model this criterion selects indices 138, 378 and 447. Elements at these indices are clear outliers, each more than 15 standard deviations away from the mean. Each such selected dimension exhibits the coordinate-level phenomenon

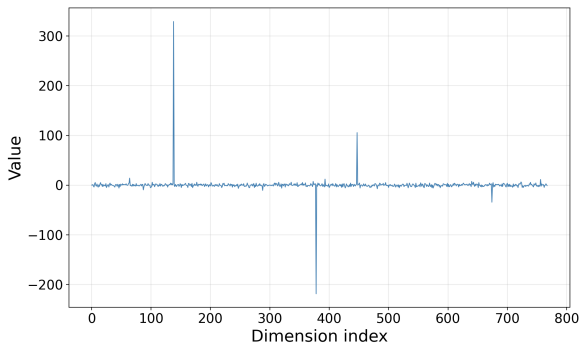


Figure 7: **A few coordinates of EPE_1 have massive activations.** Coordinate values of EPE_1 . Most coordinates are near zero; dims 138, 378, and 447 exhibit extremely large magnitudes.

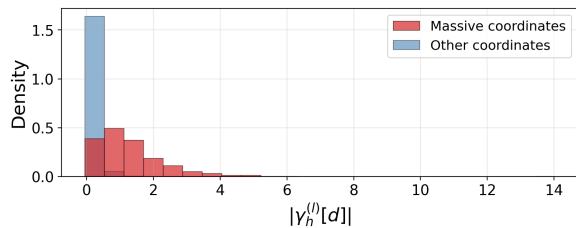


Figure 8: **Full (non-truncated) histogram of $|\gamma_h^{(l)}[d]|$.** Same data as fig. 4 with no axis restriction. The massive-activation coordinates of EPE_1 (red) extend to a long right tail that lies well beyond the range of all other coordinates (blue).

described in section 3.2.5 (i.e., large $|\gamma_h^{(l)}[d]|$ and a strong contribution to the source-agnostic shift). See fig. 7 for a visualization of EPE_1 . It is clear visually that coordinates 138, 378, and 447 have conspicuously larger norms than other indices.

A.3 Full Coordinate-Level Alignment Histogram

Figure 8 shows the full (non-truncated) version of fig. 4.

B Related Work

Our work continues a growing line of efforts to understand attention sinks and related extreme-token phenomena, including empirical, mechanistic, and functional analyses by Xiao et al. (2024), Yu et al. (2024a), Gu et al. (2025), Guo et al. (2024), Sun et al. (2024), Cancedda (2024), Barbero et al. (2025), Ruscio et al. (2025), Qiu et al. (2026), Sun et al. (2026), Hong and Lee (2025), de Llano et al. (2026), Ran-Milo (2026), and Dial (2025). Of particular relevance are the following works, which resonate most closely with our findings; below we clarify how our analysis relates to and extends each of them.

Gu et al. (2025) show that attention sinks emerge during pretraining, that their location depends on factors such as the loss and the data distribution,

and that sinks behave more like key-bias-like extra scores than semantically meaningful content. Our analysis is complementary but substantially more fine-grained. Rather than stopping at the conclusion that BOS semantics are not the right explanation, we show what specifically replaces that explanation in GPT-2: the sink is driven by the positional signal at position 1 after it is transformed by the first-layer MLP. Concretely, we separate BOS token signal from positional information with a direct intervention: nullifying the BOS token leaves the sink essentially unchanged, whereas removing the first positional embedding largely destroys it; moreover, swapping the EPE strongly affects the sink, while swapping the raw positional embedding does not.

Sun et al. (2024) identify massive activations as outlier coordinates whose values are largely input-invariant, function as indispensable bias terms, and concentrate attention on the tokens that carry them. We build directly on this observation, but add the missing upstream explanation for the GPT-2 sink: in our setting, the relevant massive activations are not merely correlated with the sink token, but are generated from the first positional encoding through the first-layer MLP. Put differently, we do not only show that massive activations matter for sinks; we trace the sink-relevant massive coordinates back to a concrete source, namely the MLP-transformed positional signal at position 1.

Dial (2025) is especially close in spirit to our work. He shows in GPT-2-small that a prominent sink-associated coordinate grows in an early MLP layer, tracing its progression through the MLP computation. Our analysis sharpens this in two ways. First, we formalize the relevant object as the effective positional encoding (EPE), and verify quantitatively that it closely tracks the net positional signal written by the first layer. Second, we show causally that the crucial object is not merely an early-MLP feature, but specifically the MLP-processed positional signal: swapping the EPE sharply disrupts the sink, whereas swapping the raw positional embedding does not.

Dial (2025) also shows that downstream queries tend to align with the key-weight direction associated with the sink-related massive coordinate. We extend this by identifying the precise parameter-level mechanism through which this alignment affects attention scores. The key object is the alignment between $EPE_1 W_{k,h}^{(l)}$ and the learned query bias $b_{Q,h}^{(l)}$, which induces a large source-agnostic score shift $\Delta_{1,h}^{(l)}$ toward position 1, connecting the observed alignment to a concrete sink-forming circuit rather than treating it as an isolated geometric regularity.

B.1 EPE and PE Swap Details

Interventions (d) and (e) test whether the sink is controlled by the direction of the EPE (or raw PE) at position 1. Let m_j denote the first-layer MLP output at position j , and let $\hat{u}_j = u_j/\|u_j\|$ be the unit-norm direction to swap ($u_j = \text{EPE}_j$ for intervention d; $u_j = p_j$ for intervention e). We compute $\alpha = m_1 \cdot \hat{u}_1$, the projection of m_1 onto the position-1 direction, project this component out of m_1 , and replace it with the position-2 direction. We apply the mirror operation at position 2 using the same coefficient:

$$m_1 \leftarrow m_1 - \alpha \hat{u}_1 + \alpha \hat{u}_2, \quad m_2 \leftarrow m_2 + \alpha \hat{u}_1 - \alpha \hat{u}_2.$$

Using the same scalar α for the removed and added components makes the swap magnitude-preserving: only the positional direction changes, not its norm.