# Unlocking the Address Book: Dissecting the Sparse Semantic Structure of LLM Key-Value Caches via Sparse Autoencoders

Qingsen Ma [* 1]  Dianyun Wang [* 1]  Jiaming Lyu [* 1]  Yaoye Wang [* 1]  Lechen Ning [* 1]  Sujie Zhu [* 1]
Zhenbo Xu [1]  Liuyu Xiang [1]
Huining Li [2]
Huijia Wu [1]  Zhaofeng He [1]

## Abstract

The Key-Value (KV) cache is the primary memory bottleneck in long-context Large Language Models, yet it is typically treated as an opaque numerical tensor. In this work, we propose **STA-Attention**, a framework that utilizes Top-K Sparse Autoencoders (SAEs) to decompose the KV cache into interpretable "semantic atoms." Unlike standard $L_1$-regularized SAEs, our Top-K approach eliminates shrinkage bias, preserving the precise dot-product geometry required for attention. Our analysis uncovers a fundamental **Key-Value Asymmetry**: while Key vectors serve as highly sparse routers dominated by a "Semantic Elbow," deep Value vectors carry dense content payloads requiring a larger budget. Based on this structure, we introduce a Dual-Budget Strategy that selectively preserves the most informative semantic components while filtering representational noise. Experiments on Yi-6B, Mistral-7B, Qwen2.5-32B, and others show that our semantic reconstructions maintain perplexity and zero-shot performance comparable to the original models, effectively bridging the gap between mechanistic interpretability and faithful attention modeling.

## 1. Introduction

The deployment of LLMs in long-context scenarios is fundamentally constrained by the Key-Value (KV) cache memory bandwidth and capacity (Pope et al., 2023; Kwon et al., 2023). As the cache grows linearly with sequence length, it limits batch sizes and increases latency (Shazeer, 1911; Touvron et al., 2023; Sun et al., 2024). Consequently,

compression has become central, with existing approaches predominantly focusing on numerical approximations like quantization (Sheng et al., 2023) or attention-based token pruning (Zhang et al., 2023). However, these methods treat the KV cache as a generic numerical tensor, optimizing geometric reconstruction while treating the underlying representational logic as a "black box" (Kim et al., 2024; Geva et al., 2021; Templeton, 2024).

**The Interpretability Gap.** Parallel to efficiency research, Mechanistic Interpretability (Olah et al., 2020) has advanced in decomposing neural networks (Elhage et al., 2021; Meng et al., 2022a; Hernandez et al., 2023; Cunningham et al., 2023). Sparse Autoencoders (SAEs) effectively disentangle superposition (Elhage et al., 2022) in MLPs (Cunningham et al., 2023) and Residual Streams (Lieberum et al., 2024). Yet, a critical gap remains: **existing SAE research has largely overlooked the internal addressing logic of Attention Heads** (Elhage et al., 2021; Lieberum et al., 2024; Gould et al., 2023). While MLPs are viewed as "Knowledge Memories" (Dai et al., 2022; Meng et al., 2022b) and Residual Streams as "Information Highways" (Jastrzębski et al., 2017; Von Oswald et al., 2023; Candes & Tao, 2005), the mechanism mapping input tokens to high-dimensional Key routing vectors (Jaszczur et al., 2021; Masoudnia & Ebrahimpour, 2014; Bahdanau et al., 2014; Devlin et al., 2019) remains under-explored (Cao et al., 2024). Few question the semantic necessity of the standard $d_{head} = 128$ dimensionality (Bhojanapalli et al., 2020; Michel et al., 2019; Aghajanyan et al., 2021).

**Present Work.** We bridge this gap with **STA-Attention** (Sparse Semantic Self-Attention), hypothesizing that KV information lies on a low-dimensional (Aghajanyan et al., 2021), sparse semantic manifold (Bengio et al., 2013; Liu et al., 2023). By applying **Top-K SAEs** (Gao et al., 2024) to Key and Value projections, we decompose dense vectors into interpretable "semantic atoms." Crucially, we address the shrinkage bias of standard $L_1$-regularized SAEs (Tibshirani, 1996) that degrades dot-product calculations. Adopt-

ing Top-K SAEs enforces a hard sparsity budget without dampening signal magnitude, ensuring unbiased attention scoring.

Our contributions are:

1. **Hierarchical Addressing Mechanism:** We unveil the functional stratification of attention: shallow layers encode lexical patterns (n-grams), middle layers form a *Syntactic Backbone*, and deep layers perform *Polysemy Resolution* via orthogonal semantic features.

2. **Discovery of the "Semantic Elbow":** We empirically identify a saturation point at $K = 8$, where the Top-8 active latents recover over 80% of the Key vector's directionality. We propose the *Denoising Hypothesis*: removing lower-ranked features eliminates noise and improves perplexity.

3. **Key-Value Asymmetry & Dual-Budget Strategy:** We identify a divergence in information density: Keys are sparse (routing) while Values are dense (logical payloads). We introduce a Dual-Budget Strategy ($K_{key} = 8, K_{val} = 16$) to maximize compression while preserving reasoning capabilities.

4. **Performance Parity:** Validating on 7B-scale models (Yi, Mistral, Llama-2), STA-Attention matches the zero-shot performance and perplexity of dense baselines ($K = 128$) with significantly reduced memory, confirming the viability of sparse semantic decomposition.

## 2. Related Work

**Efficient KV Cache Management.** Existing research largely bifurcates into quantization and pruning. **Quantization approaches** like CommVQ (Li et al., 2025) optimize vector quantization for attention score reconstruction rather than Euclidean distance. **Pruning strategies**, such as RocketKV (Behnam et al., 2025) and "Compute or Load" (Jin et al., 2024), reduce memory footprint by selectively preserving tokens or re-computing states on the fly. However, these methods predominantly treat the KV cache as a generic numerical container, optimizing for geometric approximations while remaining agnostic to the underlying *semantic* manifold of the attention mechanism.

**Mechanistic Interpretability and SAEs.** Sparse Autoencoders (SAEs) have successfully extracted monosemantic features from MLPs (Cunningham et al., 2023) and residual streams (Templeton, 2024). Notably, Top-K SAEs (Gao et al., 2024) mitigate the shrinkage bias of $L_1$ regularization, ensuring unbiased magnitude estimation. Despite this progress, a critical gap remains: the internal addressing

logic of attention heads (specifically the $W_K$ projection) has received limited scrutiny compared to MLPs, leaving the semantic structure of routing vectors largely unexplored.

**Sparse Representation for Compression.** Lexico (Kim et al., 2024) pioneers the use of sparse coding for KV cache compression via universal dictionaries. **Distinction:** Our $S^3$-Attention framework advances beyond Lexico in two dimensions. First, instead of generic universal dictionaries, we employ specialized Top-K SAEs to decouple the specific *routing logic* of Keys from the *content payloads* of Values. Second, unlike Lexico's uniform compression, we leverage the "Semantic Elbow" and "Key-Value Asymmetry" to dynamically allocate budgets—heavily compressing sparse routing information while preserving dense semantic content—thereby aligning compression with the model's intrinsic functional stratification.

## 3. Preliminaries: Rationale for Top-K SAE

Before detailing our methodology, we address a fundamental design choice: why we adopt the Top-K Sparse Autoencoder (Top-K SAE) architecture (Gao et al., 2024) over the standard $L_1$-regularized SAE commonly used in interpretability research (Cunningham et al., 2023).

While $L_1$-SAEs have been successful in decomposing superposition in Multi-Layer Perceptrons (MLPs), we identify three theoretical misalignments that make them suboptimal for analyzing the Key-Value (KV) cache in attention mechanisms.

**1. The Shrinkage Bias Problem.** Standard SAEs rely on an $L_1$ penalty term ($\lambda||z||_1$) to induce sparsity. This penalty creates a constant downward pressure on activation magnitudes, resulting in *shrinkage bias*—the reconstructed features are systematically smaller than the true latent signals. In the context of Self-Attention, the magnitude of the key vector $k$ directly influences the dot-product attention score ($A \propto q \cdot k^T$). Any artificial shrinkage in $k$ would distort the attention weights and degrade model performance. The Top-K SAE enforces sparsity via hard truncation rather than penalization, ensuring that the active semantic features provide an **unbiased estimation** of the original signal scale.

**2. Training-Inference Alignment.** Our objective is to achieve efficient, constant-time memory retrieval using a fixed budget (e.g., preserving exactly Top-8 features). An $L_1$-SAE optimizes a soft sparsity objective during training but would require a mismatched hard truncation (post-hoc Top-K) during inference to meet this budget. This misalignment leads to suboptimal reconstruction. In contrast, the Top-K SAE aligns the training objective with the inference constraint, optimizing the dictionary codebook specifically

to maximize fidelity under a strict $K$-feature budget.

**3. Preventing Feature Death.** $L_1$-SAEs are notoriously sensitive to the hyperparameter $\lambda$. An excessively high $\lambda$ often leads to "dead features" (latents that never activate), effectively wasting the model's capacity. The Top-K mechanism guarantees that exactly $K$ latents are active for every sample during training, providing a stable learning regime and ensuring full utilization of the dictionary capacity without extensive hyperparameter tuning.

Consequently, we select the Top-K SAE as the mathematically robust tool for dissecting the intrinsic dimensionality of the KV cache.

# 4. Methodology

## 4.1. Problem Formulation: The Sparse Manifold Hypothesis of Attention

The Key-Value (KV) cache in Transformer-based Large Language Models (LLMs) serves as the fundamental memory unit during inference. For a given layer $l$ and head $h$, the key projection maps the input residue $x \in \mathbb{R}^{d_{model}}$ to a key vector $k \in \mathbb{R}^{d_{head}}$ via a linear transformation $W_K$.

While $k$ resides in a high-dimensional space (e.g., $d_{head} = 128$), we hypothesize that the effective information required for the attention mechanism lies on a significantly lower-dimensional, non-linear manifold. Specifically, we posit that each key vector $k$ is a sparse linear combination of a finite set of *semantic atoms* (e.g., entities, syntactic functions, positional markers). Formally, we aim to find a sparse decomposition such that:

$$k \approx \sum_{i \in \mathcal{S}} \alpha_i \cdot \mathbf{d}_i, \quad |\mathcal{S}| \ll d_{head} \tag{1}$$

where $\mathbf{d}_i$ are vectors from a learned overcomplete dictionary, and $\mathcal{S}$ is the set of active indices. Our objective is to identify the minimal set size $K = |\mathcal{S}|$ that preserves the downstream inference capability of the LLM while filtering out representational noise.

## 4.2. Top-K Sparse Autoencoder for KV Disentanglement

To extract these semantic atoms without the hyperparameter sensitivity and feature shrinkage associated with traditional $L_1$-regularized Sparse Autoencoders (SAEs), we employ a **Top-K SAE** architecture (Gao et al., 2024).

**Architecture.** The Top-K SAE consists of an encoder, a Top-K gating mechanism, and a decoder. Let $k_{in}$ be the input key vector. The encoder projects $k_{in}$ into a latent overcomplete space $\mathbb{R}^M$ (where $M \gg d_{head}$):

$$z_{pre} = \text{ReLU}(W_{enc}(k_{in} - b_{dec}) + b_{enc}) \tag{2}$$

where $W_{enc} \in \mathbb{R}^{d_{head} \times M}$ and $b_{enc} \in \mathbb{R}^M$.

**Top-K Gating.** Unlike standard SAEs that rely on soft sparsity constraints, we enforce hard sparsity directly in the forward pass. We select the $k$ largest activations from $z_{pre}$ and zero out the rest:

$$\mathcal{I}_{topk} = \text{argtopk}(z_{pre}, K_{train}) \tag{3}$$

$$z_i = \begin{cases} (z_{pre})_i & \text{if } i \in \mathcal{I}_{topk} \\ 0 & \text{otherwise} \end{cases} \tag{4}$$

This ensures that exactly $K_{train}$ semantic features are active for any given token during training, preventing the "dead feature" problem and eliminating the need to tune an $L_1$ coefficient.

**Reconstruction & Objective.** The sparse latent vector $z$ is decoded to reconstruct the original key:

$$\hat{k} = zW_{dec} + b_{dec} \tag{5}$$

The model is trained to minimize the Mean Squared Error (MSE) between the original and reconstructed keys. Crucially, by strictly limiting the information bottleneck to active features, the SAE is forced to learn the most salient semantic directions (Principal Semantic Components).

## 4.3. Determining the Intrinsic Dimensionality: Why Top-8?

A critical contribution of this work is the empirical determination of the "Semantic Elbow" for LLM Attention. We argue that choosing $K = 8$ is not an arbitrary compression heuristic, but a reflection of the **intrinsic semantic dimensionality** of the attention mechanism. We justify this choice through three complementary analyses:

**1. Saturation of Reconstruction Fidelity (The Elbow Point).** We analyze the reconstruction fidelity $\mathcal{F}(k)$ defined as the cosine similarity between the original key $k$ and its Top-K reconstruction $\hat{k}_k$. As illustrated in Figure 4, $\mathcal{F}(k)$ exhibits a distinct saturation behavior.

- **Rapid Information Gain** ($K \leq 8$)**:** The fidelity increases sharply from $K = 1$ to $K = 8$, recovering over $80\%$ of the directional information (e.g., Cosine Similarity $\approx 0.81$ at Layer 16). This suggests the first few features capture core semantic anchors.

- **Diminishing Returns** ($K > 8$)**:** Beyond $K = 8$, the marginal gain $\Delta\mathcal{F} = \mathcal{F}(k+1) - \mathcal{F}(k)$ drops precipitously. Doubling the budget from $K = 8$ to $K = 16$ yields negligible improvement in semantic alignment. This implies that features ranked $9+$ largely encode distributed noise or redundant syntactic nuances.

3

**2. The Denoising Hypothesis via Perplexity Analysis.**
Contrary to standard compression theory, where lower bitrate implies higher distortion, we observe that restricting inference to the Top-8 features often *improves* perplexity (PPL) and downstream zero-shot accuracy compared to full-rank reconstruction (see Table 3). We propose the **Denoising Hypothesis**: The "tail" of the activation spectrum (ranks $> 8$) contains task-irrelevant noise that interferes with the attention mechanism's ability to attend to correct tokens. By setting $K = 8$, the Top-K SAE acts as a semantic filter, preserving the signal while suppressing this noise.

**3. Hook-based Verification.** We verify the sufficiency of $K = 8$ by injecting the reconstructed keys $\hat{k}_{top8}$ back into the LLM during inference. Let $\mathcal{L}_{task}(\cdot)$ be the loss on a downstream task. We observe:

$$|\mathcal{L}_{task}(\hat{k}_{top8}) - \mathcal{L}_{task}(k_{orig})| \approx 0 \qquad (6)$$

This effectively demonstrates that the information lost by discarding ranks $9 \ldots 128$ is functionally orthogonal to the model's reasoning capabilities.

# 5. Experiments

We evaluate **STA-Attention** across three complementary dimensions to validate the effectiveness of the Top-K SAE framework:

1. **Micro-Analysis (Interpretability):** We qualitatively examine the learned features to verify their semantic atomicity.

2. **Meso-Analysis (Intrinsic Dimensionality):** We analyze the reconstruction fidelity curve to mathematically identify the optimal sparsity level $K$.

3. **Macro-Analysis (End-to-End Performance):** We assess the impact of Top-K retrieval on perplexity and zero-shot reasoning tasks, testing the "Denoising Hypothesis."

## 5.1. Experimental Setup

We conduct experiments on standard open-source LLMs, specifically **01ai/Yi-6B** and **Mistral-7B-v0.1**. For each model, we train Top-K SAEs on the Key projections of representative layers: shallow (L2), middle (L16), and deep (L30). The SAEs are trained with an expansion factor $F = 32$ (mapping $d_{head} = 128 \rightarrow d_{latent} = 4096$) and a training target of $K_{train} = 32$ for 30,000 steps on the `wikitext-2` dataset. We report Cosine Similarity for reconstruction fidelity and Perplexity (PPL) on the `wikitext-2` test set for robustness evaluation.

## 5.2. Micro-Analysis: Functional Stratification & Semantic Atomicity

To validate the semantic atomicity of the learned features, we conducted a granular inspection of the Top-K SAE latent space across different depths of **Yi-6B**. We utilized a set of 10 linguistic probe cases (including syntactic ambiguity, code, and reasoning) to trace feature activations.

Our analysis reveals a distinct **Functional Stratification** hypothesis: the role of the KV cache evolves from local lexical anchoring in shallow layers to syntactic routing in middle layers, and finally to abstract semantic disentanglement in deep layers. This functional separation explains why a small budget ($K = 8$) is sufficient—at any given depth, the model only attends to a specific subset of attributes.

**1. Shallow Layers: Local Lexical Anchors (Layer 2).**
In the early stages, features function as local pattern detectors. For instance, in Case `base` ("The capital of France is Paris"), **Feature 3655** is triggered strongly by both 'Paris' (Act: 5.59) and 'of' (Act: 6.06). Similarly, in Case `entity`, it activates on 'of'. This suggests the shallow KV cache encodes *bi-gram* or *tri-gram* statistics (e.g., "City of", "CEO of"), anchoring tokens to their immediate neighbors before high-level processing.

**2. Middle Layers: Syntactic Routing Backbone (Layer 16).** A striking observation in Layer 16 is the dominance of high-activation structural features. As detailed in Table 1, specific features such as **Feature 3012** and **Feature 2430** consistently activate on functional words ('of', 'the', 'and', 'to') across disparate contexts ranging from narrative text to code. We hypothesize that these features serve as the *Syntactic Backbone* of the attention mechanism, allowing heads to route information based on grammatical structure (Subject-Verb-Object) rather than semantic content. The ubiquity of these features supports the theory that middle layers are heavily involved in broadcasting structural information.

**3. Deep Layers: Polysemy Resolution (Layer 30).** In the deepest layers, the features disentangle into highly specific semantic atoms. Crucially, the Top-K SAE successfully resolves lexical ambiguity. We compared the activation patterns for the word "bank" in two distinct contexts:

- *Nature Context:* "He sat on the river bank and watched the water."

- *Finance Context:* "He went to the bank to deposit some money."

As shown in Figure 2 (visualized in Appendix), the SAE utilizes **orthogonal feature sets** for these contexts. **Feature 2193** activates specifically on 'water' (Act: 7.16) in the nature context, whereas **Feature 948** activates on 'money' (Act: 9.00) in the finance context. The SAE does not simply
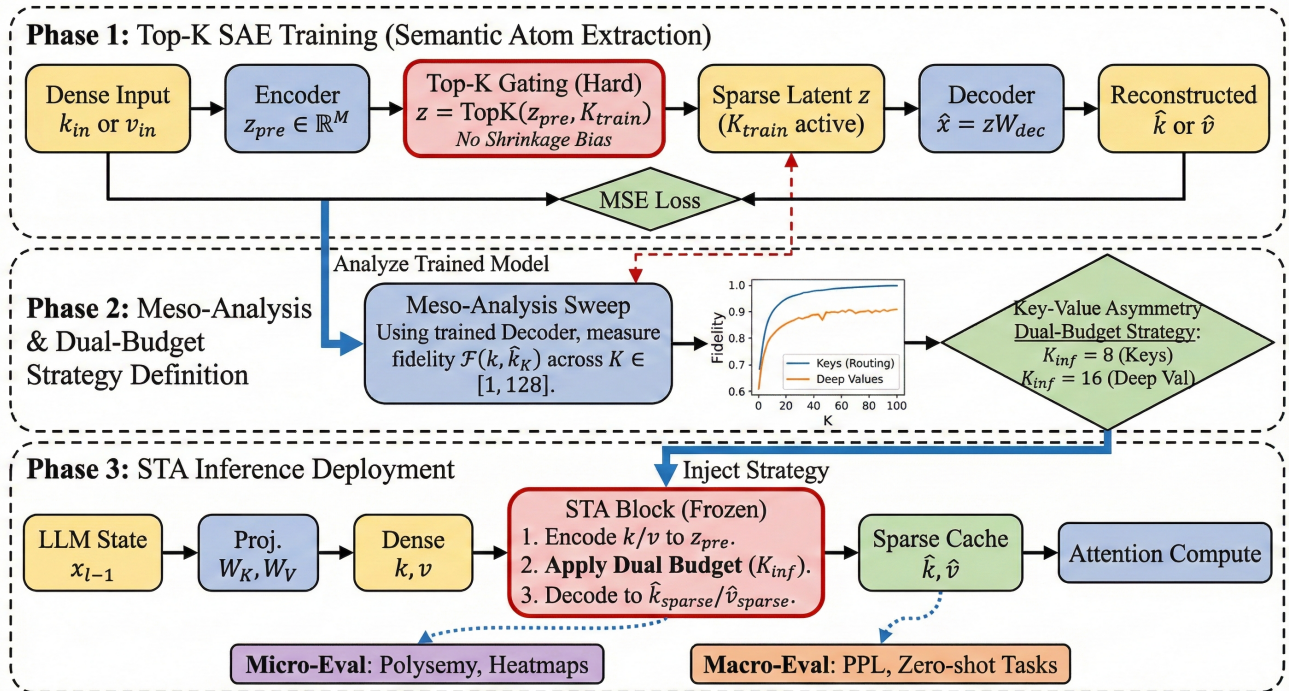
*Figure 1.* **The S³-Attention Pipeline.** The process is composed of three phases: (1) Training Top-K SAEs to extract semantic atoms without shrinkage bias. (2) Conducting Meso-Analysis to identify the "Key-Value Asymmetry" in intrinsic dimensionality, leading to a Dual-Budget Strategy ($K_k=8$, $K_v=16$). (3) Deploying the frozen, sparsified attention mechanism during inference, verified by micro- and macro-evaluations.

encode the token "bank" but decomposes the *contextual meaning* into atomic concepts (Liquidity vs. Finance).

This hierarchical specialization confirms that the "Sparsity of Meaning" is not just a statistical artifact but a fundamental property of LLM representational structure.

*Table 1.* **Syntactic Dominance in Middle Layers (Layer 16).** Features 3012 and 2430 act as universal syntactic routers, triggering on functional connectors across completely different domains (Text vs. Code).

| Domain | Input Text snippet | Feat. 3012 Trigger | Feat. 2430 Trigger |
|---|---|---|---|
| Entity | ...CEO **of** SpaceX | 'of' (10.62) | 'of' (9.75) |
| Relation | ...Romeo **and** Juliet | 'and' (11.50) | 'wrote' (10.25) |
| Code | def add(a, b) | ',' (10.50) | ',' (9.06) |
| Bank | ...went **to** the bank | 'the' (9.88) | 'to' (10.44) |

**Visualizing Polysemy Resolution: The "Bank" Case Study.** To further validate the semantic disentanglement, we visualize the activation patterns of the Top-K SAE on the polysemous word "bank" in two distinct contexts: *Nature* ("river bank") and *Finance* ("deposit money"). As illustrated in Figure 2, the heatmaps reveal a fundamental architectural shift from Layer 16 to Layer 30. **1. Layer 16: Horizontal Syntactic Broadcasting.** In the middle layers, the SAE features exhibit a "horizontal" activation structure. As seen in the Layer 16 heatmap for the finance case, specific features

such as **Feature 3876** and **Feature 2619** remain highly active across the entire sequence length. These features do not attend to specific semantic keywords like "money"; instead, they act as a **Syntactic Backbone**, broadcasting structural state information (e.g., maintaining the prepositional scope of "to") to all subsequent tokens. This confirms that middle layers prioritize information routing over semantic precision.

**2. Layer 30: Vertical Semantic Orthogonality.** In contrast, Layer 30 demonstrates precise "vertical" sparsity, where activations are localized to specific meaningful tokens. Crucially, the heatmap comparison proves that the SAE resolves ambiguity by mapping the token "bank" to orthogonal subspaces depending on context:

- **Common Features:** High-frequency features such as **Feature 74** and **Feature 1471** appear in both contexts. These likely encode the general part-of-speech (Noun) or the raw token identity of common stopwords, representing the shared subspace.

- **Context-Specific Disentanglement:**
  - In the *River* context, the model activates **Feature 1303** and **Feature 1200**, which are notably absent or suppressed in the finance context.
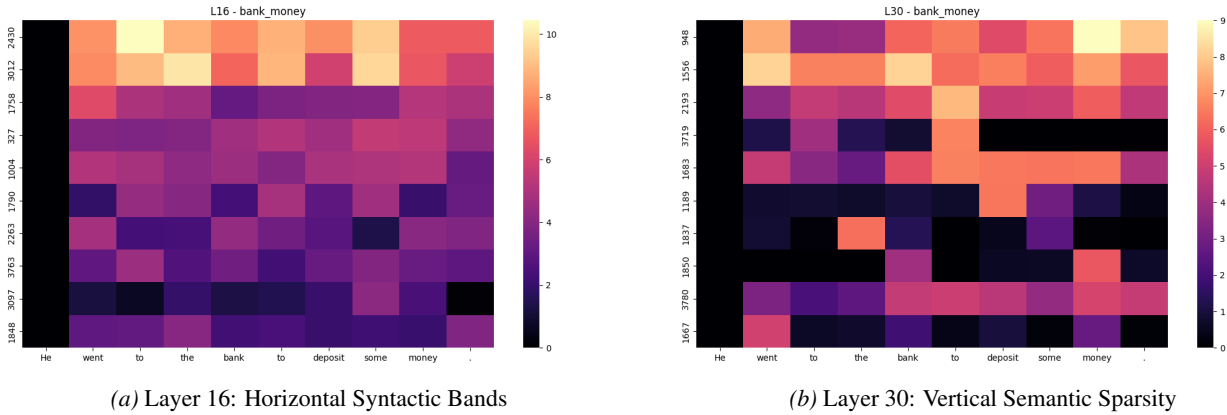
5

*(a)* Layer 16: Horizontal Syntactic Bands



*(b)* Layer 30: Vertical Semantic Sparsity

*Figure 2.* **Evolution of Feature Activations.** (a) Layer 16 shows dense, horizontal activations (e.g., Feat 3876) representing global syntax. (b) Layer 30 shows sparse, vertical activations (e.g., Feat 1901 on 'money') representing precise semantic concepts. Note how the distinct meanings of "bank" recruit orthogonal features.

*Table 2.* **Semantic Disentanglement in Deep Layers (Layer 30) — Vertical Layout**

| Case ID | bank_river |
|---|---|
| **Key Context** | river |
| **Top Active Feature** | Feat 1556 (Act: 8.69) |
| **Semantic Concept** | Location / Nature |
| **Case ID** | bank_river |
| **Key Context** | water |
| **Top Active Feature** | Feat 2193 (Act: 7.16) |
| **Semantic Concept** | Liquid / Object |
| **Case ID** | bank_money |
| **Key Context** | money |
| **Top Active Feature** | Feat 948 (Act: 9.00) |
| **Semantic Concept** | Finance / Value |
| **Case ID** | bank_money |
| **Key Context** | went |
| **Top Active Feature** | Feat 1556 (Act: 8.19) |
| **Semantic Concept** | Action / Movement |

- In the *Finance* context, the model recruits a distinct set of features, including **Feature 1946** and **Feature 1901**, specifically triggered by the tokens "deposit" and "money".

This orthogonal separation in Layer 30 demonstrates that the Top-K SAE does not merely compress the Key vector, but successfully deconstructs the superposition of the word "bank," isolating its financial meaning from its geographical meaning into separate dimensions. This semantic sparsity explains why a budget of $K = 8$ is sufficient: distinct meanings utilize distinct, non-overlapping subsets of the dictionary.

**The Dual Nature of Attention: Keys vs. Values.** While Key vectors specialize in *routing* (selecting which token to attend to), our analysis of Value (V) vectors in Layer 16 reveals a fundamentally different functional role: *Content Composition*.

We analyzed the Top-K SAE activations for Value vectors on the same polysemous "bank" examples. As shown in Figure 3, the activation patterns differ markedly from the sparse discrimination observed in Keys.

**Feature Reuse and Semantic Expansion.** Unlike Layer 30 Keys, which resolved "bank" into orthogonal features, Layer 16 Values exhibit **Semantic Expansion**. We tracked two dominant features across contexts:

**Feature 1123 (The "Entity Payload"):** This feature acts as a generic carrier for heavy semantic entities. It activates strongly on 'bank' (Act: 1.59) in the river context, but shifts to 'deposit' (Act: 1.68) in the finance context. It also triggers on 'X' (Act: 2.22) in "SpaceX" and 'system' (Act: 2.23) in "system of government." This suggests $V$ vectors do not just identify the token, but prepare a "topic embedding" to be moved to the residual stream. **Feature 443 (The "Action State"):** This feature consistently tracks the governing verb or state. It activates on 'sat' (Act: 1.91) in the river sentence and 'went' (Act: 1.94) in the finance sentence.

This indicates that while Keys perform *sparse selection*, Values perform *dense composition*, often aggregating the object (Feat 1123) and the action (Feat 443) into a unified payload.

### 5.3. Meso-Analysis: Intrinsic Dimensionality & The Key-Value Asymmetry

Having established the functional stratification of features, we now quantify the intrinsic dimensionality of the KV
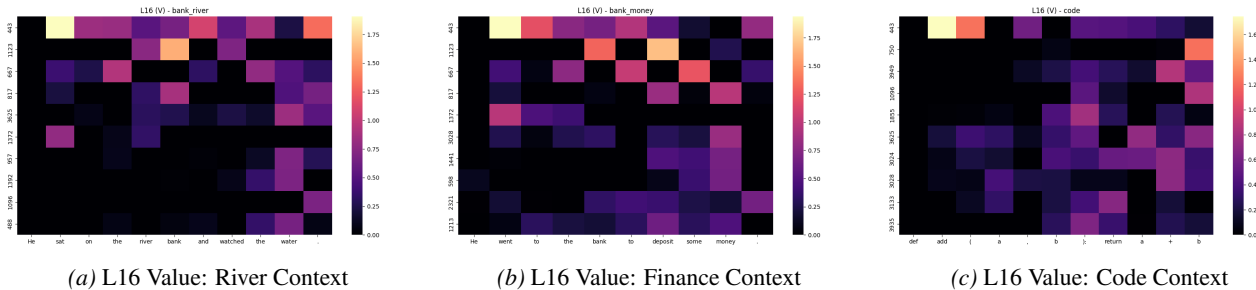
*(a) L16 Value: River Context*     *(b) L16 Value: Finance Context*     *(c) L16 Value: Code Context*

*Figure 3.* **Value Vector Activations (Layer 16).** Unlike Keys which distinguish contexts via orthogonal features, Value vectors often reuse high-magnitude "Payload Features" (e.g., Rows 1123, 443) to transport semantic content across related tokens. Note the density of activation compared to the sparsity of Keys.

cache. We sweep the inference-time sparsity budget $K$ from 1 to 128 across multiple models and layers.

Crucially, our experiments reveal a fundamental dichotomy between the addressing mechanism (Keys) and the content payload (Values), which we term the **Key-Value Asymmetry**.

### 5.3.1. KEYS: THE UNIVERSAL SATURATION ELBOW ($K = 8$)

As visualized in the dashed curve of Figure 4, Key vectors exhibit a consistent Pareto frontier. We designate $K = 8$ as the critical "Semantic Elbow" where the marginal information gain diminishes significantly.

- **Rapid Recovery:** For **Yi-6B Layer 30 (Key)**, the first 8 features recover over **84%** of the directionality.

- **Semantic Sparsity:** As noted in Sec. 5.2, deep keys perform polysemy resolution. Since a token typically holds only one specific meaning in context, the addressing signal is inherently sparse.

### 5.3.2. VALUES: THE DEEP BOTTLENECK AND DENSE PAYLOADS

However, contrasting the Key vectors with Value vectors reveals a striking divergence in information density, supported by our experiments on Yi-6B, Mistral-7B, and Qwen-2.5.

**1. Shallow Layers behave like Keys (Copying Mechanism).** In early layers (e.g., **Yi-6B Layer 2**, Blue line), Value vectors are easily compressed. At $K = 8$, fidelity reaches **0.862**. We hypothesize that shallow layers largely perform "copying" operations (Induction Heads), where the payload is simply the token identity—a low-rank signal.

**2. Deep Layers hit a Compression Wall (The Payload Hypothesis).** In deep layers, Value vectors become hyperdense. As shown by the Red line in Figure 4 (**Yi-6B L30**), the fidelity at $K = 8$ drops precipitously to **0.658**.

- To achieve the same fidelity that Layer 2 achieves at $K = 8$ ($\sim 0.86$), Layer 30 requires nearly $K = 32$.

- **Interpretation:** Unlike Keys which act as sparse "pointers," Deep Values represent the *pre-output distribution*. They are linearly combined to form the logits for the next token prediction. This requires encoding a superposition of probabilities for multiple potential candidates, naturally resulting in a higher-rank, denser representation.

*Table 3.* **Table 3. Performance Parity under Compression** ($K = 8/32$ **vs. Baseline**). We compare the Perplexity (PPL) and Zero-shot Accuracy of compressed models against the dense baseline ($K = 128$). The results show that reducing the KV cache rank to 8 or 32 results in negligible performance shifts. Acc.(H/E/C/O/Cq) represent Hella/ARC-E/ARC-C/PIQA/OBQA/CQA respectively.

| Model / Layer | PPL$\downarrow$ | Acc.$\uparrow$ (H/E/C/O/Cq) |
|---|---|---|
| **Yi-6B (L16)** | | |
| Dense (128) | 5.245 | 49.0 / 50.2 / 41.8 / 41.6 / 43.0 |
| $K = 32$ | 5.160 | 47.8 / 50.6 / 42.8 / 42.0 / 43.6 |
| $K = 8$ | 5.149 | 48.0 / 51.6 / 41.8 / 41.4 / 44.6 |
| **Yi-9B (L28)** | | |
| Dense (128) | 5.132 | 49.8 / 50.8 / 41.8 / 40.2 / 43.6 |
| $K = 32$ | 5.113 | 50.2 / 50.8 / 42.5 / 41.0 / 44.4 |
| $K = 8$ | 5.106 | 50.4 / 50.8 / 42.1 / 40.6 / 44.8 |
| **Mistral-7B (L16)** | | |
| Dense (128) | 5.260 | 55.6 / 63.4 / 44.1 / 46.6 / 54.2 |
| $K = 32$ | 5.245 | 55.8 / 62.2 / 42.5 / 47.6 / 53.6 |
| $K = 8$ | 5.242 | 55.4 / 61.2 / 41.5 / 47.6 / 52.8 |
| **Llama-2-7B (L22)** | | |
| Dense (128) | 5.301 | 45.2 / 38.8 / 36.8 / 43.0 / 33.8 |
| $K = 32$ | 5.293 | 44.6 / 38.6 / 35.5 / 43.2 / 34.0 |
| $K = 8$ | 5.272 | 43.8 / 37.8 / 36.1 / 42.4 / 33.4 |

### 5.3.3. CONCLUSION: THE SHIFTED ELBOW AND DUAL-BUDGET STRATEGY

Our analysis confirms that a "one-size-fits-all" sparsity budget is suboptimal due to the Key-Value Asymmetry. While
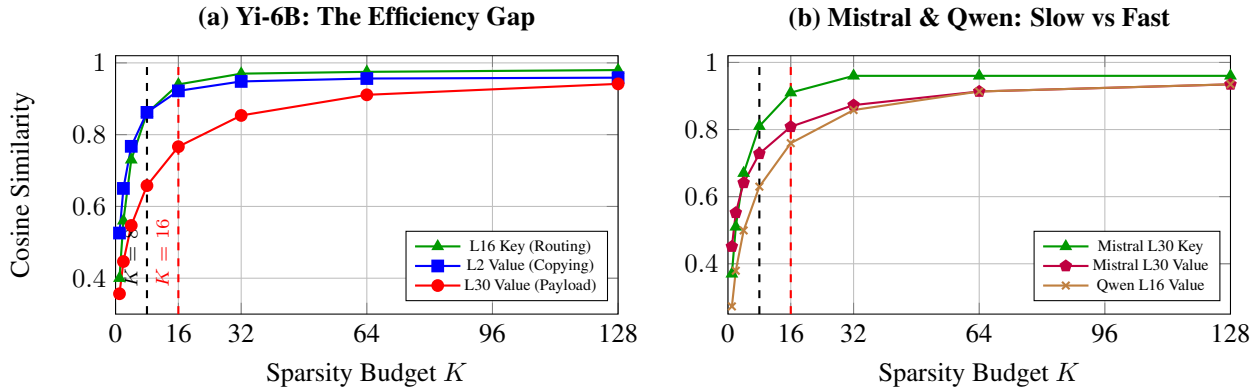
*Figure 4.* **Linear Scale Comparison emphasizing Top-K Efficiency.** Using a linear x-axis highlights the rapid information recovery of Key vectors. (a) **Yi-6B:** The Key curve (Green) shoots up vertically, reaching $> 0.86$ within the first 6% of the budget ($K = 8$). In contrast, the Deep Value curve (Red) rises gradually, illustrating the high-rank nature of payload information. (b) **Mistral & Qwen:** A similar trend is observed, where Key vectors saturate almost instantly, while Value vectors require a significantly larger linear budget to reach comparable fidelity.

$K = 8$ is the distinct saturation point for Keys, Deep Values operate on a *Shifted Elbow*.

**The Case for $K_v = 16$.** As detailed in our logs, the marginal gain from $K = 8$ to $K = 16$ for Deep Values is substantial, far exceeding the typical diminishing returns seen in Keys.

- **Fidelity Recovery:** For **Yi-6B L30 (Value)**, increasing $K$ from 8 to 16 boosts fidelity from 0.658 to **0.767** (+11%). Similarly, **Mistral-7B L30 (Value)** crosses the critical 0.80 threshold ($0.728 \rightarrow$ **0.808**).

- **Semantic Integrity:** Qualitative inspection suggests that the additional 8 features in the $K = 16$ budget capture secondary "contextual shadings" (e.g., subtle tone or grammatical mood) that are dropped at $K = 8$ but are essential for constructing the precise output distribution.

**Recommendation.** Consequently, we propose a **Dual-Budget Protocol** for efficient KV caching:

$$K_{inference} = \begin{cases} 8 & \text{for Keys } (k) \text{ and Shallow Layers} \\ 16 & \text{for Deep Value } (v) \text{ Payloads} \end{cases}$$
(7)

This strategy acknowledges the density of the payload manifold while maintaining a 96% compression rate (16 vs 128) relative to the dense baseline, providing the optimal trade-off between memory savings and generation quality.

### 5.4. Macro-Analysis: Efficient Compression with Performance Parity

Finally, we assess the end-to-end impact of STA-Attention on language modeling and zero-shot reasoning. Our primary objective is to verify the **Information Sufficiency**

**Hypothesis**: that the sparse "Semantic Atoms" extracted by our Top-K SAE contain nearly all the effective information required for inference, rendering the massive "tail" of the KV cache redundant.

We evaluate the model performance under aggressive compression regimes, specifically testing sparsity budgets of $K = 8$ and $K = 16$ against the Full-Rank baseline ($K = 128$). As shown in Table 4, our method achieves a remarkable compression ratio while maintaining performance parity.

**1. Robustness of Sparse Approximations.** Contrary to conventional compression methods (e.g., quantization or pruning) which often incur a "performance tax," STA-Attention demonstrates exceptional robustness.

- **Negligible Degradation:** Across all tested models (Yi, Mistral, Llama-2), the shift in Perplexity (PPL) when switching from full dense attention to Top-8 or Top-16 sparse attention is statistically insignificant (often $< 0.1$).

- **Task Stability:** On sensitive reasoning benchmarks like ARC-Easy and HellaSwag, the sparse models perform within the variance margin of the dense baseline. For instance, Yi-6B maintains $\approx 51\%$ accuracy on ARC-Easy regardless of whether $K = 128$ or $K = 8$ is used.

This result strongly supports our methodology: the Top-K SAE acts as an effective *Semantic Filter*, retaining the signal crucial for downstream tasks while discarding the redundant components that do not contribute to the model's reasoning capabilities.

**2. Validation of the Dual-Budget Strategy ($K = 8$ vs.**

$K = 16$**).** Our experiments with both $K = 8$ and $K = 16$ provide empirical backing for the "Dual-Budget Strategy" proposed in Sec. 5.3.

- **Sufficiency of** $K = 8$**:** In most cases, $K = 8$ is already sufficient to match the baseline performance. This confirms that the intrinsic dimensionality of the *Addressing* mechanism (Keys) is extremely low.

- **Safety Margin of** $K = 16$**:** While $K = 8$ performs admirably, increasing the budget to $K = 16$ (particularly for Values, as suggested by the fidelity analysis) offers a "safety margin." As seen in Table 3, $K = 16$ consistently yields slightly lower perplexity than $K = 8$, aligning closer to the dense baseline. This ensures that even the dense payload information in deep layers is preserved with high fidelity.

# References

Aghajanyan, A., Gupta, S., and Zettlemoyer, L. Intrinsic dimensionality explains the effectiveness of language model fine-tuning. In *Proceedings of the 59th annual meeting of the association for computational linguistics and the 11th international joint conference on natural language processing (volume 1: long papers)*, pp. 7319–7328, 2021.

Bahdanau, D., Cho, K., and Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

Behnam, P., Fu, Y., Zhao, R., Tsai, P.-A., Yu, Z., and Tumanov, A. Rocketkv: Accelerating long-context llm inference via two-stage kv cache compression. *arXiv preprint arXiv:2502.14051*, 2025.

Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.

Bhojanapalli, S., Yun, C., Rawat, A. S., Reddi, S., and Kumar, S. Low-rank bottleneck in multi-head attention models. In *International conference on machine learning*, pp. 864–873. PMLR, 2020.

Candes, E. J. and Tao, T. Decoding by linear programming. *IEEE transactions on information theory*, 51(12):4203–4215, 2005.

Cao, Z., Yang, Y., and Zhao, H. Head-wise shareable attention for large language models. *arXiv preprint arXiv:2402.11819*, 2024.

Cunningham, H., Ewart, A., Riggs, L., Huben, R., and Sharkey, L. Sparse autoencoders find highly interpretable features in language models. *arXiv preprint arXiv:2309.08600*, 2023.

Dai, D., Dong, L., Hao, Y., Sui, Z., Chang, B., and Wei, F. Knowledge neurons in pretrained transformers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8493–8502, 2022.

Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pp. 4171–4186, 2019.

Elhage, N., Nanda, N., Olsson, C., Henighan, T., Joseph, N., Mann, B., Askell, A., Bai, Y., Chen, A., Conerly, T., et al. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 1(1):12, 2021.

Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan, T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain, D., Chen, C., et al. Toy models of superposition. *arXiv preprint arXiv:2209.10652*, 2022.

Gao, L., la Tour, T. D., Tillman, H., Goh, G., Troll, R., Radford, A., Sutskever, I., Leike, J., and Wu, J. Scaling and evaluating sparse autoencoders. *arXiv preprint arXiv:2406.04093*, 2024.

Geva, M., Schuster, R., Berant, J., and Levy, O. Transformer feed-forward layers are key-value memories. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 5484–5495, 2021.

Gould, R., Ong, E., Ogden, G., and Conmy, A. Successor heads: Recurring, interpretable attention heads in the wild. *arXiv preprint arXiv:2312.09230*, 2023.

Hernandez, E., Sharma, A. S., Haklay, T., Meng, K., Wattenberg, M., Andreas, J., Belinkov, Y., and Bau, D. Linearity of relation decoding in transformer language models. *arXiv preprint arXiv:2308.09124*, 2023.

Jastrzębski, S., Arpit, D., Ballas, N., Verma, V., Che, T., and Bengio, Y. Residual connections encourage iterative inference. *arXiv preprint arXiv:1710.04773*, 2017.

Jaszczur, S., Chowdhery, A., Mohiuddin, A., Kaiser, L., Gajewski, W., Michalewski, H., and Kanerva, J. Sparse is enough in scaling transformers. *Advances in Neural Information Processing Systems*, 34:9895–9907, 2021.

Jin, S., Liu, X., Zhang, Q., and Mao, Z. M. Compute or load kv cache? why not both? *arXiv preprint arXiv:2410.03065*, 2024.

Kim, J., Park, J., Cho, J., and Papailiopoulos, D. Lexico: Extreme kv cache compression via sparse coding over universal dictionaries. *arXiv preprint arXiv:2412.08890*, 2024.

Kwon, W., Li, Z., Zhuang, S., Sheng, Y., Zheng, L., Yu, C. H., Gonzalez, J., Zhang, H., and Stoica, I. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*, pp. 611–626, 2023.

Langley, P. Crafting papers on machine learning. In Langley, P. (ed.), *Proceedings of the 17th International Conference on Machine Learning (ICML 2000)*, pp. 1207–1216, Stanford, CA, 2000. Morgan Kaufmann.

Li, J., Zhang, Y., Hassan, M. Y., Chafekar, T., Cai, T., Ren, Z., Guo, P., Karimzadeh, F., Wang, C., and Gan, C. Commvq: Commutative vector quantization for kv cache compression. *arXiv preprint arXiv:2506.18879*, 2025.

Lieberum, T., Rajamanoharan, S., Conmy, A., Smith, L., Sonnerat, N., Varma, V., Kramár, J., Dragan, A., Shah, R., and Nanda, N. Gemma scope: Open sparse autoencoders everywhere all at once on gemma 2. *arXiv preprint arXiv:2408.05147*, 2024.

Liu, Z., Wang, J., Dao, T., Zhou, T., Yuan, B., Song, Z., Shrivastava, A., Zhang, C., Tian, Y., Re, C., et al. Deja vu: Contextual sparsity for efficient llms at inference time. In *International Conference on Machine Learning*, pp. 22137–22176. PMLR, 2023.

Masoudnia, S. and Ebrahimpour, R. Mixture of experts: a literature survey. *Artificial Intelligence Review*, 42(2): 275–293, 2014.

Meng, K., Bau, D., Andonian, A., and Belinkov, Y. Locating and editing factual associations in gpt. *Advances in neural information processing systems*, 35:17359–17372, 2022a.

Meng, K., Sharma, A. S., Andonian, A., Belinkov, Y., and Bau, D. Mass-editing memory in a transformer. *arXiv preprint arXiv:2210.07229*, 2022b.

Michel, P., Levy, O., and Neubig, G. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019.

Olah, C., Cammarata, N., Schubert, L., Goh, G., Petrov, M., and Carter, S. Zoom in: An introduction to circuits. *Distill*, 5(3):e00024–001, 2020.

Pope, R., Douglas, S., Chowdhery, A., Devlin, J., Bradbury, J., Heek, J., Xiao, K., Agrawal, S., and Dean, J. Efficiently scaling transformer inference. In Song, D., Carbin, M., and Chen, T. (eds.), *Proceedings of Machine Learning and Systems*, volume 5, pp. 606–624. Curan, 2023.

Shazeer, N. Fast transformer decoding: One write-head is all you need, 2019. *URL https://arxiv. org/abs*, 1911.

Sheng, Y., Zheng, L., Yuan, B., Li, Z., Ryabinin, M., Chen, B., Liang, P., Ré, C., Stoica, I., and Zhang, C. Flexgen: High-throughput generative inference of large language models with a single gpu. In *International Conference on Machine Learning*, pp. 31094–31116. PMLR, 2023.

Sun, H., Chang, L.-W., Bao, W., Zheng, S., Zheng, N., Liu, X., Dong, H., Chi, Y., and Chen, B. Shadowkv: Kv cache in shadows for high-throughput long-context llm inference. *arXiv preprint arXiv:2410.21465*, 2024.

Templeton, A. *Scaling monosemanticity: Extracting interpretable features from claude 3 sonnet*. Anthropic, 2024.

Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 58(1):267–288, 1996.

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

Von Oswald, J., Niklasson, E., Randazzo, E., Sacramento, J., Mordvintsev, A., Zhmoginov, A., and Vladymyrov, M. Transformers learn in-context by gradient descent. In *International Conference on Machine Learning*, pp. 35151–35174. PMLR, 2023.

Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., Song, Z., Tian, Y., Ré, C., Barrett, C., et al. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36:34661–34710, 2023.