
Training Transformers for KV Cache Compressibility

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Long-context language modeling is increasingly constrained by the Key–Value
2 (KV) cache, whose memory and decode-time access costs scale linearly with
3 the prefix length. This bottleneck has motivated a range of context-compression
4 methods, from token-level summarization to recent optimization-based KV cache
5 compression methods. These post-hoc methods operate on the KV cache of a
6 fixed pretrained model, so their effectiveness is fundamentally limited by how
7 well the model’s internal representations can be compressed. In this work, we
8 formalize the notion of KV compressibility and show that it is a property of
9 the *learned representations*, rather than of the context alone. We prove that
10 almost any sequence-to-vector function admits both highly compressible and
11 inherently non-compressible transformer implementations, highlighting the
12 need to guide transformers toward compressible representations during training.
13 Motivated by this, we propose **KV-Compression Aware Training (KV-CAT)**, a
14 continued pretraining procedure that incentivizes the emergence of compressible
15 representations. We introduce a train-time KV sparsification policy that masks KV
16 slots during training. This forces the model to use fewer KV slots and encourages it
17 to learn representations amenable to post-hoc compression. Empirically, we show
18 that KV-CAT improves the quality–budget tradeoff of downstream compression
19 methods across retrieval, long-context question answering, and perplexity-based
20 evaluation of compressed-prefix continuation.

21 1 Introduction

22 Language models (LMs) are increasingly deployed in long-horizon settings, from understanding
23 large codebases, long-form documents, and personal data repositories [41; 23; 1], to long-form
24 reasoning and continuously learning agents [8; 60; 35]. Using autoregressive transformer-based LMs
25 in these settings introduces a significant memory bottleneck: the Key–Value (KV) cache. During
26 inference, these models must store the key and value vectors for every token, at every layer and
27 (KV-)attention head. For long sequences, this cache can dominate both memory usage and decoding
28 cost, turning context length into a primary serving bottleneck [59; 34; 52].

29 A substantial body of recent work has sought to mitigate this bottleneck, broadly falling into two
30 categories. The first category focuses on designing more efficient alternatives to transformers.
31 Examples include linear attention mechanisms [29; 11; 20; 54], state space models [18; 17], sparse
32 attention variants [3; 57], and more. These methods significantly reduce the computational cost of
33 long-context language modeling. However, this efficiency typically comes at the cost of empirical
34 performance, and such models still lag behind transformers at scale.

35 The second category focuses on inference-time interventions applied to a fixed pretrained transformer.
36 These methods either operate on the input context or, more generally, directly on the KV cache.
37 Early approaches are largely heuristic, including textual summarization [46; 42], learned token
38 filtering [24], and policies based on attention patterns, recency, heavy-hitter behavior, or layer-wise

39 importance [59; 52; 34; 6]. More recently, optimization-based KV cache compression¹ methods have
 40 emerged as powerful techniques. For example, Eyuboglu et al. [16] use gradient-based optimization
 41 to match the distribution induced by the original cache, while Zweiger et al. [61] employ layer-wise
 42 objectives to reproduce its attention traces.

43 The success of KV cache compression suggests that the full KV cache often contains redundancies.
 44 However, as most approaches operate on a **fixed model**, they are inherently limited by how well that
 45 particular model’s representations can be compressed. Crucially, this compressibility is not deter-
 46 mined by the input sequence alone: two transformers can use very different internal representations
 47 when processing the same sequence, yet have identical next-token distributions. Consequently, some
 48 transformers’ KV caches may be more amenable to compression than others.

49 **This work.** In this paper, we study the following question: can transformer LMs be trained in a
 50 way that leads to KV caches that are more amenable to post-hoc compression? This shifts the target
 51 from the KV cache compression algorithm to the model itself. To address this question, we introduce
 52 the notion of **KV-compressibility**. Informally, a transformer is KV-compressible if there exists a
 53 compression policy that maps the KV cache of long input sequences to a shorter KV cache while
 54 preserving the model’s next token distribution.

55 Our theoretical results show that for almost any sequence-to-vector function, there exist transformer
 56 implementations whose prefix can be compressed to a single KV pair, as well as implementations for
 57 which any non-trivial compression incurs a constant error. To build intuition, we consider a motivating
 58 example: character histogram computation. We show that natural transformer implementations of this
 59 computation may produce incompressible token representations, whereas more structured alternatives
 60 are highly compressible. This motivates treating KV compressibility as an explicit training objective.

61 Guided by this perspective, we propose **KV-Compression Aware Training (KV-CAT)**, a continued
 62 pretraining (CPT) procedure that promotes the emergence of KV-compressible internal represen-
 63 tations. Starting from a pretrained transformer, KV-CAT introduces a train-time KV sparsification
 64 policy which masks out a constant fraction of the KV slots. The training objective combines a
 65 self-distillation loss, which matches the masked model’s distribution to the dense model’s distribution,
 66 and an NTP loss applied to the unmasked forward pass to preserve uncompressed model behavior.
 67 This exposes the model to the information bottleneck induced by KV cache compression during
 68 training, encouraging it to reorganize its representations into a more compressible form while
 69 maintaining performance in the uncompressed setting. Importantly, rather than replacing post-hoc
 70 compression methods, KV-CAT makes transformers more amenable to them.

71 To evaluate KV-CAT, we apply it to QWEN2.5 models [51] via continuous pretraining and test
 72 state-of-the-art optimization-based KV compression methods on the resulting checkpoints. We
 73 evaluate suffix completion under compressed prefixes, retrieval from compressed contexts, and
 74 compressed long-context QA on LongBench v2 [2]. Across compression budgets, model sizes, and
 75 compression methods, KV-CAT consistently improves the quality–budget tradeoff, yielding gains
 76 of up to **3.21**× in suffix perplexity retention, **5**× in optimization speed, **68%** in retrieval accuracy,
 77 and **39%** in long-context QA.

78 2 KV Cache Compression: Problem Formulation

79 In this section, we theoretically formalize the KV cache compression problem. A *KV cache compres-*
 80 *sion policy* for a transformer with L layers is a collection of functions $\mathcal{C} = (c_1, \dots, c_L)$, where each
 81 c_ℓ maps sequences of n key–value pairs to sequences of length $r(n) \leq n$. For KV pairs in layer ℓ

$$(K, V) = ([k_1, \dots, k_n]^\top, [v_1, \dots, v_n]^\top), \quad (1)$$

82 we write

$$c_\ell(K, V) = ([\tilde{k}_1, \dots, \tilde{k}_{r(n)}]^\top, [\tilde{v}_1, \dots, \tilde{v}_{r(n)}]^\top), \quad (2)$$

83 where $r(n)$ is referred to as the *compression budget*. Given a transformer M , a compression policy
 84 \mathcal{C} , and a sequence of tokens $\mathbf{a} = (a_1, \dots, a_n)$, we define a compressed model $M_{\mathcal{C}, \mathbf{a}}$ that shares the
 85 parameters of M but uses a modified forward pass. For an input sequence \mathbf{b} , the computation of
 86 $M_{\mathcal{C}, \mathbf{a}}(\mathbf{b})$ proceeds as in $M([\mathbf{a}, \mathbf{b}])$, except for the attention over prefix tokens, which is augmented

¹These methods are also commonly referred to as KV compaction methods in the literature.

87 the following way: at layer ℓ , let $(\mathbf{K}_a, \mathbf{V}_a)$ denote the KV cache of \mathbf{a} , based on the original model
 88 \mathbf{M} and let \mathbf{Y} be the representations of \mathbf{b} after layer $\ell - 1$ of $\mathbf{M}_{\mathbf{C},a}$. We compute

$$\mathbf{Q}_b = \mathbf{Y}\mathbf{W}_Q, \quad \mathbf{K}_b = \mathbf{Y}\mathbf{W}_K, \quad \mathbf{V}_b = \mathbf{Y}\mathbf{W}_V, \quad (3)$$

89 and replace the prefix KV cache $(\mathbf{K}_a, \mathbf{V}_a)$ with its compressed version $(\tilde{\mathbf{K}}_a, \tilde{\mathbf{V}}_a) = c_\ell(\mathbf{K}_a, \mathbf{V}_a)$.
 90 The attention computation then becomes

$$\text{AttnHead}(\mathbf{Y}) = \text{softmax} \left(\frac{1}{\sqrt{d_k}} \mathbf{Q}_b \begin{bmatrix} \tilde{\mathbf{K}}_a \\ \mathbf{K}_b \end{bmatrix}^\top \right) \begin{bmatrix} \tilde{\mathbf{V}}_a \\ \mathbf{V}_b \end{bmatrix}. \quad (4)$$

91 Finally, following standard convention in next-token prediction transformers, we define the output
 92 of $\mathbf{M}_{\mathbf{C},a}(\mathbf{b})$ to be the representation of the final token. Thus, compression modifies only the prefix
 93 KV pairs, leaving the remainder of the computation unchanged. We now formalize the notion of
 94 transformer *KV-compressibility*.

95 **Definition 2.1.** Fix $N \in \mathbb{N}$, $\varepsilon > 0$, and let $r : \mathbb{N} \rightarrow \mathbb{N}$ be a budget function. A transformer \mathbf{M} is said
 96 to be (N, ε, r) -compressible² if there exists a KV cache compression policy \mathbf{C} with budget r such
 97 that for every pair of sequences \mathbf{a} , \mathbf{b} of lengths n and k respectively, and with a combined length
 98 $n + k \leq N$, it holds that

$$\|\mathbf{M}([\mathbf{a}, \mathbf{b}]) - \mathbf{M}_{\mathbf{C},a}(\mathbf{b})\| < \varepsilon. \quad (5)$$

99 A more formal treatment of these terms and definitions is provided in Appendix B.1.

100 3 Motivation and Theoretical Results

101 Our analysis begins with the following theorem, which shows that almost any sequence-to-vector
 102 function can admit transformer implementations with different levels of KV-compressibility. Formal
 103 proofs of all results in this section are provided in Appendix B.

104 **Theorem 3.1.** Let A be a finite alphabet, and let $f : \bigcup_{n \leq N} A^n \rightarrow \mathbb{R}^{d_{\text{out}}}$ be a sequence-to-vector
 105 function. Suppose that there exists a sequence \mathbf{a} of length n and two sequences $\mathbf{b}^1, \mathbf{b}^2$ of length
 106 $k \leq N - n$ such that

$$f([\mathbf{a}, \mathbf{b}^1]) \neq f([\mathbf{a}, \mathbf{b}^2]). \quad (6)$$

107 Then for every $\varepsilon > 0$, there exists a transformer that approximates f and is $(N, \varepsilon, 1)$ -compressible.
 108 Additionally, for every $C > 0$, there exists a transformer of the same architecture, that approximates
 109 f but is not (N, C, r) -compressible for any budget function satisfying $r(n) < n$.

110 Before turning to training methods, we build intuition through a simple motivating example. We show
 111 that even for natural sequence-to-vector functions, simple transformer implementations can yield
 112 non-compressible representations, while more structured, highly compressible implementations exist.

113 **Motivating example: histogram computation.** Let $A = [m]$ be a finite alphabet. Given a sequence
 114 $\mathbf{a} = (a_1, \dots, a_n) \in A^n$, the histogram function computes the empirical distribution of symbols,

$$f_{\text{hist}}(\mathbf{a}) = \left(\frac{n_1}{n}, \dots, \frac{n_m}{n} \right), \quad (7)$$

115 where $n_i = |\{j : a_j = i\}|$. This is a natural sequence-to-vector function that depends only on
 116 aggregate statistics of the input. A straightforward way to implement this function with a 2-layer
 117 transformer is as follows. First, take the embedding map to be $\text{emb}(a_i, i) = \mathbf{e}_{a_i}$. We take the first
 118 transformer layer to be the identity and then apply a second attention layer with uniform attention
 119 weights (by setting $\mathbf{W}_Q = \mathbf{0}$). So that the representation of the final token becomes the average of all
 120 previous token embeddings, obtaining

$$\mathbf{M}(\mathbf{a}) = \frac{1}{n} \sum_{i=1}^n \mathbf{e}_{a_i} = f_{\text{hist}}(\mathbf{a}). \quad (8)$$

121 This gives a natural implementation of the histogram function with a very simple architecture (see
 122 Appendix B.2 for details). Interestingly, we find that this implementation is not compressible:

²We include N explicitly, as many natural sequence-to-vector functions require model dimension scaling with N (e.g., $O(N)$ or $O(\log N)$) in order to achieve arbitrarily accurate approximation; see e.g., Sanford et al. [48]; Yehudai et al. [55].

123 **Proposition 3.2.** For any $N \in \mathbb{N}$ and any budget function r with $r(N-1) < N-1$ there exists a con-
 124 stant $C > 0$ such that the above transformer implementation of f_{hist} is not (N, C, r) -compressible.

125 *Proof sketch.* Since $\mathbf{W}_O = \mathbf{0}$ in the first layer, compressing the prefix KV cache there has no effect
 126 on the suffix representations. In the second layer, $\mathbf{W}_Q = \mathbf{0}$, so attention weights are uniform and
 127 independent of the keys. Consequently, the compressed model computes the histogram by averaging
 128 suffix token embeddings with the compressed value vectors provided by the policy. As compression
 129 replaces the prefix by only $r(n)$ vectors, the normalization factor of this average changes. Thus,
 130 for any compression policy \mathcal{C} ,

$$M_{\mathcal{C}, \mathbf{a}}(\mathbf{b}) - M([\mathbf{a}, \mathbf{b}]) = \left(\frac{1}{r(n)+k} - \frac{1}{n+k} \right) \sum_{i=1}^k e_{b_i} + c(\mathbf{a}), \quad (9)$$

131 where $c(\mathbf{a})$ depends only on the prefix. The suffix-dependent term can vary over many values as \mathbf{b}
 132 changes, while $c(\mathbf{a})$ is fixed. Therefore, no non-trivial compression policy can uniformly approximate
 133 the full model for all suffixes, implying a constant lower bound on the compression error. \square

134 In contrast, a slight modification of the same architecture yields a highly compressible implementation.

135 **Proposition 3.3.** For every $N \in \mathbb{N}$ and $\varepsilon > 0$, there exists a f_{hist} transformer implementation with
 136 the same architecture as above that is $(N, \varepsilon, 1)$ -compressible.

137 The key idea is to leverage positional information, provided by the compression policy, to recover
 138 sequence lengths after compression and re-scale the averaged representation so that the correct his-
 139 togram can be reconstructed from a single KV pair. Importantly, since f_{hist} is permutation-invariant,
 140 standard training provides no incentive for the model to preserve positional information, making such
 141 compressible solutions unlikely to emerge without additional supervision. This further illustrates that
 142 compressibility depends on the learned representation and motivates compression-aware training.

143 4 Training for KV Cache Compressibility

144 As shown above, the same task may admit both compressible and non-compressible transformer
 145 solutions. This raises the question of whether transformers can be guided toward compressible repre-
 146 sentations during training. To this end, we introduce KV-Compression Aware Training (KV-CAT), a
 147 continued pretraining procedure that explicitly encourages compressible internal representations.

148 **Train-time KV sparsification policy.** Our goal is to train transformers whose representations are
 149 more amenable to SOTA optimization-based KV cache compression methods. Ideally, one would
 150 incorporate such methods directly into the training loop; however, this is computationally prohibitive.
 151 A practical alternative is to introduce a simple and scalable training-time KV sparsification policy,
 152 such as random KV slot dropping, attention-based filtering, or a lightweight parameterized policy. In
 153 our ablations (Appendix F.1), we find that a parameterized policy provides the strongest performance,
 154 in terms of both downstream compression and retention of uncompressed accuracy. We hypothesize
 155 that this is because adaptive policies better align with the test-time behavior of optimization-based
 156 compressors. Accordingly, in our experiments, we implement KV-CAT using lightweight learned
 157 routers, described in Appendix C, while emphasizing that the framework is general and can
 158 accommodate a range of sparsification policies.

159 **Training objective.** We train the transformer using an augmented next-token prediction (NTP)
 160 objective. For a sequence $\mathbf{a} = (a_1, \dots, a_n)$, let $p_{\theta}^{\text{mask}}(\cdot | \mathbf{a}_{<i})$ denote the model distribution when
 161 masking is applied to the KV cache, and $p_{\theta}^{\text{dense}}(\cdot | \mathbf{a}_{<i})$ denote the distribution under the standard
 162 (unmasked) forward pass. During training, we optimize

$$\mathcal{L}(\theta) = \lambda_{\text{mask}} \mathcal{L}_{\text{mask}} + \lambda_{\text{anchor}} \mathcal{L}_{\text{anchor}} + \lambda_{\text{budget}} \mathcal{L}_{\text{budget}}. \quad (10)$$

163 The first term is given by

$$\mathcal{L}_{\text{mask}} = \frac{1}{n} \sum_{i=1}^n D_{\text{KL}}(\text{sg}[p_{\theta}^{\text{dense}}(\cdot | a_{<i})] || p_{\theta}^{\text{mask}}(\cdot | a_{<i})), \quad (11)$$

164 where sg denotes the stop-gradient operation. This term trains the masked forward pass to match
 165 the dense model distribution. The second term $\mathcal{L}_{\text{anchor}} = \frac{1}{n} \sum_{i=1}^n -\log p_{\theta}^{\text{dense}}(a_i | \mathbf{a}_{<i})$, preserves

Table 1: **KV-CAT preserves performance.** Values are normalized accuracy (%) over 1000 examples.

| Model | Variant | HellaSwag | WinoGrande | PIQA | OpenBookQA | ARC-E | ARC-C | Avg. |
|--------------|---------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| QWEN2.5-0.5B | Base | 54.6 | 52.9 | 70.4 | 38.2 | 60.8 | 32.2 | 51.5 |
| | KV-CAT | 53.6 | 54.1 | 72.1 | 37.4 | 63.6 | 32.4 | 52.2 |
| QWEN2.5-1.5B | Base | 68.6 | 60.5 | 76.2 | 40.0 | 73.3 | 45.2 | 60.6 |
| | KV-CAT | 66.4 | 60.5 | 75.5 | 40.6 | 73.9 | 43.5 | 60.1 |

166 the distribution of the dense model. When the KV-sparsification policy is learnable, we also include
 167 a budget loss designed to maintain a desirable retention rate, see Appendix C for details. At
 168 evaluation time, we use the unmasked forward pass, on top of which we can apply standard KV cache
 169 compression methods. The output of KV-CAT is a standard transformer whose representations are
 170 trained to be more amenable to compression.

171 5 Empirical Evaluation

172 Table 2: **Long-form QA after context compression.** We evaluate the
 173 models on seven tasks from LongBench v2. The KV caches of the
 174 questions’ contexts are compressed with a gradient-based method.

| Subdomain | 10% keep | | 20% keep | | 50% keep | |
|---------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| | Base | KV-CAT | Base | KV-CAT | Base | KV-CAT |
| Academic | 18.1 | 26.6 | 19.1 | 27.7 | 21.3 | 26.6 |
| Agent history QA | 20.0 | 30.0 | 25.0 | 30.0 | 30.0 | 30.0 |
| Knowledge graph reasoning | 33.3 | 33.3 | 26.7 | 26.7 | 20.0 | 26.7 |
| Legal | 27.3 | 48.5 | 27.3 | 45.5 | 27.3 | 42.4 |
| Many-shot learning | 28.6 | 19.0 | 28.6 | 28.6 | 38.1 | 33.3 |
| New language translation | 30.0 | 30.0 | 15.0 | 35.0 | 15.0 | 35.0 |
| Table QA | 11.1 | 27.8 | 22.2 | 22.2 | 38.9 | 33.3 |
| Total | 22.2 | 30.3 | 22.2 | 30.8 | 25.3 | 31.2 |

175
 176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186 with compressed context. Across experiments, we evaluate two state-of-the-art optimization-based
 187 KV compression methods, Attention Matching [61] and a gradient-based optimization method
 188 adapted from Eyuboglu et al. [16], on two QWEN2.5 model sizes (0.5B and 1.5B). Overall,
 189 KV-CAT preserves model performance in the uncompressed setting, while consistently improving
 190 the performance of downstream compression across all tasks and compression budgets, yielding gains
 191 of up to $3.21\times$ in suffix perplexity retention, **68%** in retrieval accuracy, and **39%** in long-context
 192 QA. Results for retrieval from a compressed context are in Table 3, and results for compressed
 193 long-context QA are in Table 2. For the full experimental setup and results, see Appendix D, for
 194 implementation details, see Appendix E, and for additional results and ablations, see Appendix F.

195 6 Conclusion

196 We argue that KV cache compressibility is determined
 197 not only by the input, but also by the transformer’s
 198 learned representation. We prove that almost any
 199 sequence-to-vector function admits both highly
 200 compressible and inherently non-compressible imple-
 201 mentations, motivating compression-aware training.
 202 Guided by this perspective, we introduce KV-CAT,
 203 a training procedure that encourages compressible
 204 representations by exposing the model to KV compres-
 205 sion during training. Empirically, models trained with
 206 KV-CAT consistently improve the quality–budget trade-
 207 off of state-of-the-art post-hoc compression methods
 208 across a range of tasks. Overall, our results suggest a
 209 complementary approach to post-hoc KV compression:
 210 training transformers to become compressible.

We evaluate the effect of KV-CAT training on transformer-based LLMs. Our evaluation considers four settings: (1) standard uncompressed evaluation, to verify that KV-CAT does not degrade base model performance, (2) suffix completion given a compressed prefix, (3) retrieval from a compressed context, and (4) long-context question answering

Table 3: **Needle retrieval from compressed haystack.** Exact-match accuracy (%) over 100 examples per keep ratio.

| Keep | QWEN2.5-0.5B | | QWEN2.5-1.5B | |
|------|--------------|-------------|--------------|-------------|
| | Base | KV-CAT | Base | KV-CAT |
| 5% | 18 | 17 | 20 | 20 |
| 10% | 12 | 15 | 20 | 20 |
| 15% | 15 | 16 | 21 | 21 |
| 20% | 15 | 13 | 22 | 23 |
| 25% | 20 | 22 | 24 | 28 |
| 30% | 23 | 34 | 41 | 44 |
| 35% | 21 | 32 | 46 | 54 |
| 40% | 24 | 38 | 42 | 55 |
| 50% | 28 | 47 | 49 | 67 |
| Mean | 19.6 | 26.0 | 31.7 | 36.9 |

211 **References**

- 212 [1] Simran Arora and Christopher Ré. Can foundation models help us achieve perfect secrecy?
213 *arXiv preprint arXiv:2205.13722*, 2022.
- 214 [2] Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng
215 Xu, Lei Hou, Yuxiao Dong, et al. Longbench v2: Towards deeper understanding and reason-
216 ing on realistic long-context multitasks. In *Proceedings of the 63rd Annual Meeting of the*
217 *Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3639–3664, 2025.
- 218 [3] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer.
219 *arXiv preprint arXiv:2004.05150*, 2020.
- 220 [4] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. PIQA: Reasoning
221 about physical commonsense in natural language. In *Proceedings of the AAAI Conference on*
222 *Artificial Intelligence*, volume 34, pages 7432–7439, 2020. doi: 10.1609/aaai.v34i05.6239.
223 URL <https://ojs.aaai.org/index.php/AAAI/article/view/6239>.
- 224 [5] Aydar Bulatov, Yuri Kuratov, and Mikhail S. Burtsev. Recurrent memory transformer, 2022.
- 225 [6] Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Yucheng Li, Tianyu Liu, Keming Lu, Wayne
226 Xiong, Yue Dong, Junjie Hu, and Wen Xiao. PyramidKV: Dynamic kv cache compression
227 based on pyramidal information funneling, 2025.
- 228 [7] Rujikorn Charakorn, Edoardo Cetin, Shinnosuke Uesaka, and Robert Tjarko Lange. Doc-to-lora:
229 Learning to instantly internalize contexts. *arXiv preprint arXiv:2602.15902*, 2026.
- 230 [8] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang
231 Hu, Yuhang Zhou, Te Gao, and Wanxiang Che. Towards reasoning era: A survey of long
232 chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.
- 233 [9] Tong Chen, Hao Fang, Patrick Xia, Xiaodong Liu, Benjamin Van Durme, Luke Zettlemoyer,
234 Jianfeng Gao, and Hao Cheng. Generative adapter: Contextualizing language models in
235 parameters with a single forward pass. *arXiv preprint arXiv:2411.05877*, 2024.
- 236 [10] Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models
237 to compress contexts, 2023.
- 238 [11] Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane,
239 Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking
240 attention with performers. *arXiv preprint arXiv:2009.14794*, 2020.
- 241 [12] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick,
242 and Oyvind Tafjord. Think you have solved question answering? try ARC, the AI2 reasoning
243 challenge. *arXiv preprint arXiv:1803.05457*, 2018. URL <https://arxiv.org/abs/1803.05457>.
244
- 245 [13] Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang,
246 and Nazli Goharian. A discourse-aware attention model for abstractive summarization of
247 long documents. In *Proceedings of the 2018 Conference of the North American Chapter of*
248 *the Association for Computational Linguistics: Human Language Technologies, Volume 2*
249 *(Short Papers)*, pages 615–621, New Orleans, Louisiana, 2018. Association for Computational
250 Linguistics. doi: 10.18653/v1/N18-2097. URL <https://aclanthology.org/N18-2097/>.
- 251 [14] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of*
252 *control, signals and systems*, 2(4):303–314, 1989.
- 253 [15] Yam Eitan. The centered convex body whose marginals have the heaviest tails. *arXiv preprint*
254 *arXiv:2110.14382*, 2021.
- 255 [16] Sabri Eyuboglu, Ryan Ehrlich, Simran Arora, Neel Guha, Dylan Zinsley, Emily Liu, Will Ten-
256 nien, Atri Rudra, James Zou, Azalia Mirhoseini, and Christopher Ré. Cartridges: Lightweight
257 and general-purpose long context representations via self-study, 2025.

- 258 [17] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces.
259 *arXiv preprint arXiv:2312.00752*, 2023.
- 260 [18] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
261 state spaces. *arXiv preprint arXiv:2111.00396*, 2021.
- 262 [19] Nathan Habib, Clémentine Fourier, Hynek Kydlíček, Thomas Wolf, and Lewis Tunstall. Lighte-
263 val: A lightweight framework for llm evaluation. [https://github.com/huggingface/
264 lighteval](https://github.com/huggingface/lighteval), 2023. GitHub repository.
- 265 [20] Alex Horn, Ali Kheradmand, and Mukul Prasad. Delta-net: Real-time network verification
266 using atoms. In *14th USENIX Symposium on Networked Systems Design and Implementation
267 (NSDI 17)*, pages 735–749, 2017.
- 268 [21] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*,
269 4(2):251–257, 1991.
- 270 [22] Sukjun Hwang, Brandon Wang, and Albert Gu. Dynamic chunking for end-to-end hierarchical
271 sequence modeling. *arXiv preprint arXiv:2507.07955*, 2025.
- 272 [23] Pranab Islam, Anand Kannappan, Douwe Kiela, Rebecca Qian, Nino Scherrer, and Bertie
273 Vidgen. Financebench: A new benchmark for financial question answering. *arXiv preprint
274 arXiv:2311.11944*, 2023.
- 275 [24] Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. LLMingua: Com-
276 pressing prompts for accelerated inference of large language models, 2023.
- 277 [25] Huiqiang Jiang, Qianhui Wu, Xufang Luo, Dongsheng Li, Chin-Yew Lin, Yuqing Yang, and Lili
278 Qiu. LongLLMingua: Accelerating and enhancing llms in long context scenarios via prompt
279 compression, 2024.
- 280 [26] Samuel Karlin and William J Studden. Optimal experimental designs. *The Annals of Mathemat-
281 ical Statistics*, 37(4):783–815, 1966.
- 282 [27] Samuel Karlin and William J Studden. Tchebycheff systems: With applications in analysis and
283 statistics. (*No Title*), 1966.
- 284 [28] Samuel Karlin and Zvi Ziegler. Chebyshevian spline functions. *Siam Journal on Numerical
285 Analysis*, 3(3):514–543, 1966.
- 286 [29] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers
287 are rnns: Fast autoregressive transformers with linear attention. In *International conference on
288 machine learning*, pages 5156–5165. PMLR, 2020.
- 289 [30] Jang-Hyun Kim, Jinuk Kim, Sangwoo Kwon, Jae W Lee, Sangdoo Yun, and Hyun Oh Song.
290 Kvzip: Query-agnostic kv cache compression with context reconstruction. *arXiv preprint
291 arXiv:2505.23416*, 2025.
- 292 [31] Junhyuck Kim, Jongho Park, Jaewoong Cho, and Dimitris Papailiopoulos. Lexico: Extreme
293 KV cache compression via sparse coding over universal dictionaries. In *Proceedings of the
294 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine
295 Learning Research*, pages 30672–30687, 2025.
- 296 [32] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman
297 Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and
298 Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances
299 in Neural Information Processing Systems*, volume 33, pages 9459–9474, 2020.
- 300 [33] Yucheng Li, Bo Dong, Frank Guerin, and Chenghua Lin. Compressing context to enhance
301 inference efficiency of large language models. In *Proceedings of the 2023 Conference on
302 Empirical Methods in Natural Language Processing*, pages 6342–6353, 2023. doi: 10.18653/
303 v1/2023.emnlp-main.391.

- 304 [34] Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye,
305 Tianle Cai, Patrick Lewis, and Deming Chen. SnapKV: Llm knows what you are looking for
306 before generation, 2024.
- 307 [35] Zhuoling Li, Xiaogang Xu, Zhenhua Xu, SerNam Lim, and Hengshuang Zhao. Larm: Large
308 auto-regressive model for long-horizon embodied intelligence. *arXiv preprint arXiv:2405.17424*,
309 2024.
- 310 [36] Yewei Liu, Xiyuan Wang, Yansheng Mao, Yoav Gelbery, Haggai Maron, and Muhan Zhang.
311 Shine: A scalable in-context hypernetwork for mapping context to lora in a single pass. *arXiv*
312 *preprint arXiv:2602.06358*, 2026.
- 313 [37] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture
314 models, 2016.
- 315 [38] CA Micchelli and Allan Pinkus. Moment theory for weak chebyshev systems with applications
316 to monosplines, quadrature formulae and best one-sided L^1 -approximation by spline functions
317 with fixed knots. *SIAM Journal on Mathematical Analysis*, 8(2):206–230, 1977.
- 318 [39] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct
319 electricity? a new dataset for open book question answering. In *Proceedings of the 2018*
320 *Conference on Empirical Methods in Natural Language Processing*, pages 2381–2391, Brussels,
321 Belgium, 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1260. URL
322 <https://aclanthology.org/D18-1260/>.
- 323 [40] Jesse Mu, Xiang Lisa Li, and Noah Goodman. Learning to compress prompts with gist tokens,
324 2023.
- 325 [41] Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. Using
326 an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International*
327 *Conference on Software Engineering*, pages 1–13, 2024.
- 328 [42] Emre Okular. Context Engineering - Short-Term Memory Management with Sessions from
329 OpenAI Agents SDK, September 2025.
- 330 [43] Matanel Oren, Michael Hassid, Nir Yarden, Yossi Adi, and Roy Schwartz. Transformers are
331 multi-state RNNs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural*
332 *Language Processing*, pages 18724–18741, 2024. doi: 10.18653/v1/2024.emnlp-main.1043.
- 333 [44] Guilherme Penedo, Hynek Kydlíček, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro
334 Von Werra, Thomas Wolf, et al. The fineweb datasets: Decanting the web for the finest text data
335 at scale. *Advances in Neural Information Processing Systems*, 37:30811–30849, 2024.
- 336 [45] Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap.
337 Compressive transformers for long-range sequence modelling. *arXiv preprint*, 2019. URL
338 <https://arxiv.org/abs/1911.05507>.
- 339 [46] Prithvi Rajasekaran, Ethan Dixon, Carly Ryan, and Jeremy Hadfield. Effective context engi-
340 neering for ai agents, September 2025.
- 341 [47] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. WinoGrande: An
342 adversarial winograd schema challenge at scale. In *Proceedings of the AAAI Conference on*
343 *Artificial Intelligence*, volume 34, pages 8732–8740, 2020. doi: 10.1609/aaai.v34i05.6399.
344 URL <https://ojs.aaai.org/index.php/AAAI/article/view/6399>.
- 345 [48] Clayton Sanford, Daniel J Hsu, and Matus Telgarsky. Representational strengths and limitations
346 of transformers. *Advances in Neural Information Processing Systems*, 36:36677–36707, 2023.
- 347 [49] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. Social IQa:
348 Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on*
349 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*
350 *on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China,
351 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1454. URL <https://aclanthology.org/D19-1454/>.
- 352

- 353 [50] Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. QUEST:
354 Query-aware sparsity for efficient long-context LLM inference. In *Proceedings of the 41st In-*
355 *ternational Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning*
356 *Research*, pages 47901–47911, 2024.
- 357 [51] Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL [https://](https://qwenlm.github.io/blog/qwen2.5/)
358 qwenlm.github.io/blog/qwen2.5/.
- 359 [52] Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
360 language models with attention sinks, 2024.
- 361 [53] Guangxuan Xiao, Jiaming Tang, Jingwei Zuo, Junxian Guo, Shang Yang, Haotian Tang, Yao
362 Fu, and Song Han. Duoattention: Efficient long-context LLM inference with retrieval and
363 streaming heads. In *International Conference on Learning Representations*, 2025. URL
364 <https://arxiv.org/abs/2410.10819>.
- 365 [54] Songlin Yang, Jan Kautz, and Ali Hatamizadeh. Gated delta networks: Improving mamba2
366 with delta rule. *arXiv preprint arXiv:2412.06464*, 2024.
- 367 [55] Gilad Yehudai, Haim Kaplan, Guy Dar, Royi Rassin, Asma Ghandeharioun, Mor Geva, and
368 Amir Globerson. When can transformers count to n? *arXiv preprint arXiv:2407.15160*, 2024.
- 369 [56] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov,
370 and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30,
371 2017.
- 372 [57] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti,
373 Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird:
374 Transformers for longer sequences. *Advances in neural information processing systems*, 33:
375 17283–17297, 2020.
- 376 [58] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can
377 a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the*
378 *Association for Computational Linguistics*, pages 4791–4800, Florence, Italy, 2019. Association
379 for Computational Linguistics. doi: 10.18653/v1/P19-1472. URL [https://aclanthology.](https://aclanthology.org/P19-1472/)
380 [org/P19-1472/](https://aclanthology.org/P19-1472/).
- 381 [59] Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao
382 Song, Yuandong Tian, Christopher Ré, Clark Barrett, Zhangyang Wang, and Beidi Chen. H_2O :
383 Heavy-hitter oracle for efficient generative inference of large language models, 2023.
- 384 [60] Junhao Zheng, Chengming Shi, Xidi Cai, Qiuke Li, Duzhen Zhang, Chenxing Li, Dong Yu,
385 and Qianli Ma. Lifelong learning of large language model based agents: A roadmap. *IEEE*
386 *Transactions on Pattern Analysis and Machine Intelligence*, 2026.
- 387 [61] Adam Zweiger, Xinghong Fu, Han Guo, and Yoon Kim. Fast kv compaction via attention
388 matching, 2026.

389 A Related Work

390 **Token-level context compression and retrieval.** One approach to long-context inference reduces
391 the number of input tokens before prompting, rather than modifying the model’s hidden state.
392 Retrieval-augmented generation keeps the prompt short by selecting only relevant external
393 text [32], while prompt-compression methods prune or rewrite the provided context in token
394 space [24; 25; 33; 42; 46]. These methods are complementary to latent KV cache compression.

395 **Learned memory and soft-token compression.** Another line of work uses learned memory or soft
396 tokens as compact substitutes for long contexts. Recurrent Memory Transformers add memory tokens
397 that carry information between segments [5]; AutoCompressors adapt language models to map earlier
398 segments into summary vectors that act as soft prompts [10]; and gist-tokens train models to compress
399 prompts into reusable learned tokens [40]. DuoAttention instead learns which attention heads require
400 full retrieval over the KV cache, and which can run with a streaming-style state [53]. These methods

401 modify the model’s runtime interface, architecture, or cache policy. Another direction [9; 7; 36] uses
 402 hypernetworks to encode long contexts into model parameters, typically via LoRA adapters, rather
 403 than storing them in the prompt or KV cache.

404 **Post-hoc KV cache compression.** Another line of work focuses on reducing the KV cache of
 405 pretrained transformers at inference time. Token-selection and eviction methods keep a subset of
 406 original KV slots using attention, recency, or query-dependent importance signals [59; 43; 34; 6; 50].
 407 KVzip selects keys using reconstruction-style objectives over the context [30]. Other post-hoc
 408 methods move beyond choosing a subset of the original tokens. Cartridges optimize a compact latent
 409 KV cache for each context [16], Attention Matching constructs compact keys and values to preserve
 410 attention behavior [61], and Lexico represents KV vectors with sparse codes over learned universal
 411 dictionaries [31]. Our work is orthogonal to these compressors: we ask whether the model can be
 412 trained so that the same post-hoc methods work better.

413 B Theory

414 B.1 Transformers and KV Cache Compression

415 **Notation.** For a set A , let A^* denote the set of all finite sequences over A . We denote by e_i the i -th
 416 standard (one-hot encoded) basis vector. For $n \in \mathbb{N}$, we write $[n] = \{1, \dots, n\}$.

417 For matrices $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{Y} \in \mathbb{R}^{n \times d'}$, we denote their row-wise concatenation (along the feature
 418 dimension) by

$$[\mathbf{X}, \mathbf{Y}] \in \mathbb{R}^{n \times (d+d')}. \quad (12)$$

419 For matrices $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{Z} \in \mathbb{R}^{n' \times d}$, we denote their column-wise concatenation (along the
 420 sequence dimension) by

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Z} \end{bmatrix} \in \mathbb{R}^{(n+n') \times d}. \quad (13)$$

421 When convenient, we sometimes write $[\mathbf{X}^\top, \mathbf{Z}^\top]^\top$ instead of $\begin{bmatrix} \mathbf{X} \\ \mathbf{Z} \end{bmatrix}$.

422 For vectors $\mathbf{v}, \mathbf{u} \in \mathbb{R}^n$, we treat them as column vectors in $\mathbb{R}^{n \times 1}$ and write

$$[\mathbf{v}, \mathbf{u}] \in \mathbb{R}^{n \times 2}. \quad (14)$$

423 **Definition B.1** (Attention Head). Let $d_{\text{in}}, d_{\text{out}} \in \mathbb{N}$. An *attention head* is a function

$$\text{AttnHead} : (\mathbb{R}^{d_{\text{in}}})^* \rightarrow (\mathbb{R}^{d_{\text{out}}})^* \quad (15)$$

424 defined as follows. For a sequence $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d_{\text{in}}}$,

$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q, \quad \mathbf{K} = \mathbf{X}\mathbf{W}_K, \quad \mathbf{V} = \mathbf{X}\mathbf{W}_V, \quad (16)$$

425 where $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d_{\text{in}} \times d_{\text{out}}}$. The output is

$$\text{AttnHead}(\mathbf{X}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \mathbf{V}, \quad (17)$$

426 where the softmax is applied row-wise.

427 **Definition B.2** (Transformer Block). Fix $h, d_{\text{in}}, d_{\text{out}}, d_{\text{ff}} \in \mathbb{N}$. A *transformer block* is a function

$$\text{Block} : (\mathbb{R}^{d_{\text{in}}})^* \rightarrow (\mathbb{R}^{d_{\text{out}}})^* \quad (18)$$

428 defined as follows. For a sequence

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d_{\text{in}}}, \quad (19)$$

429 we define

$$\text{Block}(\mathbf{X}) = \text{FFN}([\text{head}_1(\mathbf{X}), \dots, \text{head}_h(\mathbf{X})]\mathbf{W}_O + \mathbf{X}), \quad (20)$$

430 where each head_i is an attention head with input and output dimensions $d_{\text{in}}, d_{\text{out}}$ respectively, and
 431 $\mathbf{W}_O \in \mathbb{R}^{h \cdot d_{\text{out}} \times d_{\text{out}}}$, and FFN is a 2-layer feed forward network applied row-wise, i.e.,

$$\text{FFN}(\mathbf{X})_i = \text{FFN}(\mathbf{x}_i), \quad i = 1, \dots, n, \quad (21)$$

432 with

$$\text{FFN}(\mathbf{x}) = \sigma(\mathbf{x}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2, \quad (22)$$

433 where $\mathbf{W}_1 \in \mathbb{R}^{d_{\text{out}} \times d_{\text{ff}}}$, $\mathbf{W}_2 \in \mathbb{R}^{d_{\text{ff}} \times d_{\text{out}}}$, $\mathbf{b}_1 \in \mathbb{R}^{d_{\text{ff}}}$, and $\mathbf{b}_2 \in \mathbb{R}^{d_{\text{out}}}$.

434 **Note.** In the theoretical analysis, we omit layer normalization, causal masking, and the residual
 435 connection following the FFN for simplicity, as is standard in prior theoretical treatments (See e.g.
 436 [48; 55]). Our results extend to the full architecture with only minor modifications. All empirical
 437 evaluations are conducted with these components included.

438 **Definition B.3** (Token Embedding). Let A be a finite alphabet. A *token embedding* function is a
 439 function

$$\text{emb} : A \times \mathbb{N} \rightarrow \mathbb{R}^{d_{\text{model}}}. \quad (23)$$

440 Given a sequence $\mathbf{a} = (a_1, \dots, a_n) \in A^*$, its embedding is the sequence

$$\text{emb}(\mathbf{a}) = [\text{emb}(a_1, 1), \dots, \text{emb}(a_n, n)]^\top \in \mathbb{R}^{n \times d_{\text{model}}}. \quad (24)$$

441 **Definition B.4** (Transformer). Let A be a finite alphabet. A *transformer* is a tuple

$$\mathbf{M} = (\text{Block}_1, \dots, \text{Block}_L, \text{emb}), \quad (25)$$

442 where each $\text{Block}_\ell : (\mathbb{R}^{d_{\ell-1}})^* \rightarrow (\mathbb{R}^{d_\ell})^*$ is a transformer block, and $\text{emb} : A^* \rightarrow (\mathbb{R}^{d_0})^*$ is a token
 443 embedding function.

444 The transformer \mathbf{M} defines a function

$$\mathbf{M} : A^* \rightarrow \mathbb{R}^{d_L} \quad (26)$$

445 as follows. For $\mathbf{a} = (a_1, \dots, a_n) \in A^*$, let

$$\mathbf{X}^{(0)} = \text{emb}(\mathbf{a}), \quad \mathbf{X}^{(\ell)} = \text{Block}_\ell(\mathbf{X}^{(\ell-1)}), \quad \ell = 1, \dots, L. \quad (27)$$

446 Then

$$\mathbf{M}(\mathbf{a}) = \mathbf{X}_n^{(L)}, \quad (28)$$

447 i.e., the output is the representation of the final token.

448 **Definition B.5** (KV Cache Compression). A *KV cache compression policy* is a tuple $\mathbf{C} =$
 449 (c_1, \dots, c_L) , where each

$$c_\ell : (\mathbb{R}^{d_\ell} \times \mathbb{R}^{d_\ell})^* \rightarrow (\mathbb{R}^{d_\ell} \times \mathbb{R}^{d_\ell})^* \quad (29)$$

450 maps a sequence of key–value pairs to a shorter sequence.

451 Concretely, for

$$(\mathbf{K}, \mathbf{V}) = ([\mathbf{k}_1, \dots, \mathbf{k}_n]^\top, [\mathbf{v}_1, \dots, \mathbf{v}_n]^\top), \quad (30)$$

452 we write

$$c_\ell(\mathbf{K}, \mathbf{V}) = (\tilde{\mathbf{K}}, \tilde{\mathbf{V}}) = ([\tilde{\mathbf{k}}_1, \dots, \tilde{\mathbf{k}}_{r(n)}]^\top, [\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_{r(n)}]^\top), \quad (31)$$

453 where $r(n) \leq n$. The function $r(\cdot)$ is called the *compression budget*.

454 Given a transformer $\mathbf{M} = (\text{emb}, \text{Block}_1, \dots, \text{Block}_L)$, a compression policy \mathbf{C} , and a *context*
 455 sequence $\mathbf{a} = (a_1, \dots, a_n)$, we define the associated *compressed transformer* $\mathbf{M}_{\mathbf{C}, \mathbf{a}}$.

456 For an input sequence $\mathbf{b} = (b_1, \dots, b_k)$, the output $\mathbf{M}_{\mathbf{C}, \mathbf{a}}(\mathbf{b})$ is obtained by running the forward pass of
 457 \mathbf{M} on the concatenated sequence $[\mathbf{a}, \mathbf{b}] = (a_1, \dots, a_n, b_1, \dots, b_k)$, with the following modifications.

458 For each ℓ and each attention head at the ℓ -th block, after computing keys and values, the compression
 459 function c_ℓ is applied to the KV pairs corresponding to the context tokens. The compressed KV pairs
 460 replace the original ones in the attention computation.

461 More precisely, let $(\mathbf{K}_\mathbf{a}, \mathbf{V}_\mathbf{a})$ denote the KV cache produced by the prefix \mathbf{a} at some attention head of
 462 Block_ℓ of the original model \mathbf{M} , and let $\mathbf{Y} \in \mathbb{R}^{k \times d_{\ell-1}}$ be the input corresponding to \mathbf{b} in the forward
 463 pass computation of $\mathbf{M}_{\mathbf{C}, \mathbf{a}}$ after block $\ell - 1$. We first compute

$$\mathbf{Q}_\mathbf{b} = \mathbf{Y} \mathbf{W}_Q, \quad \mathbf{K}_\mathbf{b} = \mathbf{Y} \mathbf{W}_K, \quad \mathbf{V}_\mathbf{b} = \mathbf{Y} \mathbf{W}_V. \quad (32)$$

464 We then compress the context KV pairs:

$$(\tilde{\mathbf{K}}_\mathbf{a}, \tilde{\mathbf{V}}_\mathbf{a}) = c_\ell(\mathbf{K}_\mathbf{a}, \mathbf{V}_\mathbf{a}). \quad (33)$$

465 The resulting attention computation is

$$\text{AttnHead}(\mathbf{Y}) = \text{softmax} \left(\frac{1}{\sqrt{d_k}} \mathbf{Q}_\mathbf{b} \begin{bmatrix} \tilde{\mathbf{K}}_\mathbf{a} \\ \mathbf{K}_\mathbf{b} \end{bmatrix}^\top \right) \begin{bmatrix} \tilde{\mathbf{V}}_\mathbf{a} \\ \mathbf{V}_\mathbf{b} \end{bmatrix}. \quad (34)$$

466 **Definition B.6.** [KV compressibility] Let $N \in \mathbb{N}$, $\varepsilon > 0$, and let $r : \mathbb{N} \rightarrow \mathbb{N}$ be a budget function. A
 467 transformer M is said to be (N, ε, r) -compressible if there exists a KV cache compression policy C
 468 with budget r such that for every pair of sequences \mathbf{a}, \mathbf{b} of lengths n and k with combined length
 469 satisfying $n + k \leq N$, it holds that

$$\|M([\mathbf{a}, \mathbf{b}]) - M_{C, \mathbf{a}}(\mathbf{b})\| < \varepsilon. \quad (35)$$

470 We include N explicitly in the definition above, as many natural sequence-to-vector functions require
 471 model dimension scaling with N (e.g., $O(N)$ or $O(\log N)$) in order to achieve arbitrarily accurate
 472 approximation; see Sanford et al. [48]; Yehudai et al. [55] for examples.

473 B.2 Motivating example: histogram computation

474 In this section we provide formal proofs of Propositions 3.2 and 3.3.

475 *Proof of Proposition 3.2.* The construction of $M = (\text{emb}, \text{Block}_1, \text{Block}_2)$ is straightforward. First,
 476 for $j \in [m]$ and $i \in [N]$, define

$$\text{emb}(j, i) = \mathbf{u}_j + \mathbf{p}_i, \quad (36)$$

477 where $\mathbf{u}_j, \mathbf{p}_i \in \mathbb{R}^m$ encode token value and token position respectively, and assume that emb is
 478 injective over $[m] \times [N]$.

479 Let $\mathbf{X} \in \mathbb{R}^{n \times m}$. We define $\text{Block}_1(\mathbb{R}^m)^* \rightarrow (\mathbb{R}^{2m})^*$ by

$$\text{Block}_1(\mathbf{X})_i \approx \rho_1(\mathbf{X}_i) \quad (37)$$

480 where $\rho_1 : \mathbb{R}^m \rightarrow \mathbb{R}^{2m}$ is applied row wise \mathbf{X} and satisfies for each $j, i \in [m] \times [N]$:

$$\rho_1(\mathbf{u}_j + \mathbf{p}_i) = \begin{bmatrix} \mathbf{e}_j \\ \mathbf{0} \end{bmatrix}. \quad (38)$$

481 Block_1 can be implemented by zeroing out all attention projections and relying on the residual
 482 connection together with a feedforward network that approximates ρ_1 to arbitrary precision (this can
 483 be obtained as a two-layer feedforward network is a universal approximator of continuous functions
 484 on compact sets [14; 21]).

485 Let $\mathbf{X} \in \mathbb{R}^{n \times 2m}$, and write each row as \mathbf{x}_i . We define

$$\text{Block}_2(\mathbf{X})_i = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \quad (39)$$

486 i.e., each output token is replaced by the average of all input tokens. This operation can be imple-
 487 mented using a single attention head by setting

$$\mathbf{W}_Q = \mathbf{W}_K = \mathbf{0}, \quad (40)$$

488 so that all attention weights are uniform, and choosing

$$\mathbf{W}_V = \mathbf{W}_O = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{I} & \mathbf{0} \end{pmatrix}. \quad (41)$$

489 This choice extracts the first half of each vector and places it in the second half before averaging.
 490 After applying a residual connection, the intermediate representation of the i -th token is then

$$\left[\frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \right] \quad (42)$$

491 Finally, the feedforward network is taken to be

$$\text{FFN}(\mathbf{x}, \mathbf{y}) = \mathbf{y}, \quad (43)$$

492 resulting in a final representation of $\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. For input sequence \mathbf{a} of length n we thus have

$$M(\mathbf{a}) \approx \frac{1}{n} \sum_{i=1}^n \mathbf{e}_{a_i} = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{a_i=1}, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{a_i=2}, \dots, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{a_i=m} \right) = f_{\text{hist}}(\mathbf{a}), \quad (44)$$

493 Thus M approximates f_{hist} to arbitrary precision. Now let $C = (c_1, c_2)$ be any compression policy
 494 with budget function $r(\cdot)$ satisfying $r(N-1) < N-1$. Fix a prefix sequence \mathbf{a} of length n and a
 495 suffix sequence \mathbf{b} of length k , and consider the computation of $M_{C,\mathbf{a}}(\mathbf{b})$.

496 First, observe that in the first block all attention projection matrices (and in particular \mathbf{W}_O) are
 497 zero. Hence, c_1 has no effect on the representations of the tokens of \mathbf{b} , and after the first block the
 498 representations are exactly as in the uncompressed model.

499 Next, consider the second block. Let $(\mathbf{K}_\mathbf{a}, \mathbf{V}_\mathbf{a})$ denote the KV cache corresponding to the prefix \mathbf{a}
 500 computed during the computation of $M(\mathbf{a})$, and let $(\tilde{\mathbf{K}}_\mathbf{a}, \tilde{\mathbf{V}}_\mathbf{a}) = c_2(\mathbf{K}_\mathbf{a}, \mathbf{V}_\mathbf{a})$ be the compressed
 501 cache. Since $\mathbf{W}_Q = 0$, all attention scores are uniform for both M and $M_{C,\mathbf{a}}$ and so

$$M_{C,\mathbf{a}}(\mathbf{b}) = \frac{1}{r(n)+k} \left(\sum_{i=1}^{r(n)} \tilde{\mathbf{v}}_i + \sum_{i=1}^k \mathbf{e}_{b_i} \right). \quad (45)$$

$$M([\mathbf{a}, \mathbf{b}]) = \frac{1}{n+k} \left(\sum_{i=1}^n \mathbf{e}_{a_i} + \sum_{i=1}^k \mathbf{e}_{b_i} \right). \quad (46)$$

502 For notational convenience, define

$$\tilde{\mathbf{v}} = \frac{1}{r(n)+k} \sum_{i=1}^{r(n)} \tilde{\mathbf{v}}_i, \quad \tilde{\mathbf{a}} = \frac{1}{n+k} \sum_{i=1}^n \mathbf{e}_{a_i}, \quad \tilde{\mathbf{b}} = \sum_{i=1}^k \mathbf{e}_{b_i}. \quad (47)$$

503 Then the difference between the compressed and full outputs can be written as

$$\|M_{C,\mathbf{a}}(\mathbf{b}) - M([\mathbf{a}, \mathbf{b}])\| = \left\| \tilde{\mathbf{v}} - \tilde{\mathbf{a}} + \left(\frac{1}{r(n)+k} - \frac{1}{n+k} \right) \tilde{\mathbf{b}} \right\|. \quad (48)$$

504 We now derive a lower bound. First, choose $\mathbf{b} = (1, 1, \dots, 1)$ and examine the first coordinate. This
 505 yields

$$\|M_{C,\mathbf{a}}(\mathbf{b}) - M([\mathbf{a}, \mathbf{b}])\| \geq \left| \tilde{v}_1 - \tilde{a}_1 + \frac{k(r(n)-n)}{(n+k)(r(n)+k)} \right|. \quad (49)$$

506 Next, choosing $\mathbf{b} = (2, 2, \dots, 2)$ removes the contribution of the last term, giving

$$\|M_{C,\mathbf{a}}(\mathbf{b}) - M([\mathbf{a}, \mathbf{b}])\| \geq |\tilde{v}_1 - \tilde{a}_1|. \quad (50)$$

507 Combining the two inequalities, we conclude that for any choice of $\tilde{\mathbf{v}}$, there exists a \mathbf{b} such that

$$\|M_{C,\mathbf{a}}(\mathbf{b}) - M([\mathbf{a}, \mathbf{b}])\| \geq \frac{k(n-r(n))}{2(n+k)(r(n)+k)}. \quad (51)$$

508 Finally, taking $n = N-1$ and $k = 1$ yields

$$\|M_{C,\mathbf{a}}(\mathbf{b}) - M([\mathbf{a}, \mathbf{b}])\| \geq \frac{N-1-r(N-1)}{2(N)(r(N-1)+1)} = C > 0, \quad (52)$$

509 which completes the proof. \square

510 **Note:** for the weakest compression case $r(n) = n-1$, the constant above is $C = O(N^{-2})$, and so as
 511 N grows the approximation bound becomes weaker, but for the extreme compression case $r(n) = c$
 512 for some $c \in \mathbb{N}$, C can be chosen independently of N .

513 *proof of Proposition 3.3.* We construct $M = (\text{emb}, \text{Block}_1, \text{Block}_2)$ similarly to the above construc-
 514 tion, with a few modifications. Like before, for $j \in [m]$ and $i \in [N]$, define

$$\text{emb}(j, i) = \mathbf{u}_j + \mathbf{p}_i, \quad (53)$$

515 where $\mathbf{u}_j, \mathbf{p}_i \in \mathbb{R}^m$ encode token value and token position respectively, and assume that emb is
 516 injective over $[m] \times [N]$.

517 Let $\mathbf{X} \in \mathbb{R}^{n \times m}$. We define $\text{Block}_1(\mathbb{R}^m)^* \rightarrow (\mathbb{R}^{4m})^*$ by

$$\text{Block}_1(\mathbf{X})_i \approx \rho_1(\mathbf{X}_i) \quad (54)$$

518 where $\rho_1 : \mathbb{R}^m \rightarrow \mathbb{R}^{4m}$ is applied row wise \mathbf{X} and satisfies for each $j, i \in [m] \times [N]$:

$$\rho_1(\mathbf{u}_j + \mathbf{p}_i) = \begin{bmatrix} \mathbf{e}_j \\ \mathbf{p}_i \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}. \quad (55)$$

519 Where each of the vector blocks above is of size m . That is, Block_1 is identical to the construction in
 520 the proof of Proposition 3.2, except that it preserves the positional component. As before, Block_1
 521 can be implemented by zeroing out all attention projections and relying on the residual connection
 522 together with a feedforward network that approximates ρ_1 to arbitrary precision. This follows from
 523 the universal approximation property of two-layer feedforward networks for continuous functions on
 524 compact sets [14; 21].

525 We construct Block_2 to be composed of a single attention head with

$$\mathbf{W}_Q = \mathbf{W}_K = \mathbf{0}, \mathbf{W}_O = \mathbf{I} \quad (56)$$

526 so that all attention weights are uniform, and choosing

$$\mathbf{W}_V = \begin{pmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}. \quad (57)$$

527 This choice extracts the first m -length block of each vector and places it in the third block before
 528 averaging, keeping the first and second vector blocks the same, and zeroing out the last vector block.
 529 After applying a residual connection, the intermediate representation of the i -th token a_i is then

$$\begin{bmatrix} \mathbf{e}_{a_i} \\ \mathbf{p}_i \\ \frac{1}{n} \sum_{j=1}^n \mathbf{e}_j \\ \mathbf{0} \end{bmatrix}. \quad (58)$$

530 finally, we choose the last feed forward network to approximate a map ρ_2 satisfying for all $i, j, k \in$
 531 $[N]$:

$$\rho_2(\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}) = \begin{cases} \mathbf{y} & \text{if } \mathbf{z} = \mathbf{0} \\ \frac{j-i+1}{j} \mathbf{y} & \text{if } \mathbf{z} = \frac{\mathbf{p}_i}{k+1} \text{ and } \mathbf{x} = \mathbf{p}_j. \end{cases} \quad (59)$$

532 Such a function can be realized to arbitrary precision by a feedforward network: the above specifica-
 533 tion is defined on a union of disjoint compact sets, and therefore admits a continuous extension to a
 534 compact domain, which can in turn be approximated by a standard feedforward network via universal
 535 approximation.

536 Equations 58 and 59 imply that, in the uncompressed setting, for any sequence \mathbf{a} of length $n \leq N$,
 537 we recover as before

$$\mathbf{M}(\mathbf{a}) \approx \frac{1}{n} \sum_{i=1}^n \mathbf{e}_{a_i} = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{1}_{a_i=1}, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{a_i=2}, \dots, \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{a_i=m} \right) = f_{\text{hist}}(\mathbf{a}), \quad (60)$$

538 and so \mathbf{M} approximates f_{hist} . We now define a compression policy $\mathbf{C} = (\mathbf{c}_1, \mathbf{c}_2)$. We first note that,
 539 like before, SINCE $\mathbf{W}_O = \mathbf{0}$ in the first block, \mathbf{c}_1 does not effect the forward pass of $\mathbf{M}_{\mathbf{C}, \mathbf{a}}$.

540 Given a prefix \mathbf{a} of length n with KV cache $(\mathbf{K}_a, \mathbf{V}_a)$, we compress it into a single KV pair:

$$c_2(\mathbf{K}_a, \mathbf{V}_a) = \left(\mathbf{0}, \begin{bmatrix} \sum_{i=1}^n \mathbf{e}_{a_i} \\ \mathbf{0} \\ \mathbf{0} \\ \mathbf{p}_n \end{bmatrix} \right). \quad (61)$$

541 That is, we store a single value vector containing (i) the unnormalized histogram of the prefix and (ii)
 542 its length encoded via \mathbf{p}_n . Now consider processing a suffix $\mathbf{b} = (b_1, \dots, b_k)$. Since the compressed
 543 prefix consists of a single KV entry with zero key, attention again produces a uniform average over
 544 the compressed prefix and suffix tokens. The representation of the final token b_k becomes

$$\begin{bmatrix} \mathbf{e}_{b_i} \\ \mathbf{p}_{n+k} \\ \frac{1}{k+1} \left(\sum_{i=1}^n \mathbf{e}_{a_i} + \sum_{i=1}^k \mathbf{e}_{b_i} \right) \\ \frac{1}{k+1} \mathbf{p}_n \end{bmatrix}. \quad (62)$$

545 Applying ρ_2 and using equation 59, we obtain

$$M_{C,a}(\mathbf{b}) \approx \frac{1 + (n+k) - n}{n+k} \cdot \frac{1}{k+1} \left(\sum_{i=1}^n \mathbf{e}_{a_i} + \sum_{i=1}^k \mathbf{e}_{b_i} \right) = \frac{1}{n+k} \sum_{i=1}^{n+k} \mathbf{e}_{c_i}, \quad (63)$$

546 where c_i ranges over $[\mathbf{a}, \mathbf{b}]$. This matches the histogram of the concatenated sequence (up to ε),
 547 completing the proof. □

548

549 B.3 Compressible Transformers for General Functions

550 In this section, we present a formal proof of Theorem 3.1. The proof is divided into two lemmas:
 551 the first (Lemma B.7) establishes the existence of compressible transformer approximations, and the
 552 second (Lemma B.8) demonstrates the existence of non-compressible ones. The theorem then follows
 553 immediately.

554 **Lemma B.7.** *Let A be a finite alphabet, and let $f : \bigcup_{n \leq N} A^n \rightarrow \mathbb{R}^{d_{out}}$ be any sequence-to-vector*
 555 *function. Then for every $\varepsilon > 0$ There exists a transformer such that*

556 1. (**Approximation**) *For every sequence $\mathbf{a} = (a_1, \dots, a_n) \in A^n$ with $n \leq N$,*

$$\|f(\mathbf{a}) - M(\mathbf{a})\| < \varepsilon. \quad (64)$$

557 2. (**Maximal compressibility**) *There exists a KV cache compression policy C with compression*
 558 *budget*

$$r(n) \equiv 1 \quad (65)$$

559 *such that for every prefix $\mathbf{a} \in A^n$ and suffix $\mathbf{b} \in A^k$ with $k+n \leq N$,*

$$\|M([\mathbf{a}, \mathbf{b}]) - M_{C,a}(\mathbf{b})\| < \varepsilon. \quad (66)$$

560 *Proof.* We assume that the token embedding map $\text{emb} : A \times \mathbb{N} \rightarrow \mathbb{R}^{d_0}$ is defined by

$$\text{emb}(a, i) = \mathbf{u}_a + \mathbf{p}_i \quad (67)$$

561 where $\|\mathbf{p}_i\| = \|\mathbf{p}_j\|$ for all $i, j \in [N]$, and $\mathbf{u}_a, \mathbf{p}_i \in \mathbb{R}^{d_0}$ for all $a \in A$. We further assume injectivity
 562 of the embedding function, that is

$$(a, i) \neq (b, j) \Rightarrow \text{emb}(a, i) \neq \text{emb}(b, j). \quad (68)$$

563 Such embeddings can be obtained, for example, using a learned lookup table for tokens combined
 564 with a positional encoding (e.g., sinusoidal) with sufficiently large frequencies.

565 Let $\mathbf{u}_\emptyset \in \mathbb{R}^{d_0}$ be a padding vector such that $\mathbf{v}_\emptyset \neq \text{emb}(a, i)$ for all $a \in A$ and $i \in \mathbb{N}$. For a sequence
 566 $\mathbf{a} = (a_1, \dots, a_n) \in A^*$ with $n \leq N$, define

$$U_{\mathbf{a}} = \{ \mathbf{P}[\text{emb}(a_1, 1), \dots, \text{emb}(a_n, n), \mathbf{u}_\emptyset, \dots, \mathbf{u}_\emptyset]^\top \mid \mathbf{P} \in \mathbb{R}^{N \times N} \text{ is a permutation matrix} \} \subset \mathbb{R}^{N \times d_0}. \quad (69)$$

567 That is, we embed the sequence, pad it to length N using \mathbf{u}_\emptyset , and then apply an arbitrary permutation
 568 to the order of the elements.

569 By injectivity of the embeddings, the sets $\{U_{\mathbf{a}} : \mathbf{a} \in A^*, |\mathbf{a}| \leq N\}$ are pairwise disjoint. Conse-
 570 quently, there exists a continuous function

$$\bar{f} : \mathbb{R}^{N \times d_0} \rightarrow \mathbb{R}^{d_{\text{out}}} \quad (70)$$

571 such that for every $\mathbf{a} \in A^*$ with $|\mathbf{a}| \leq N$ and every $[\mathbf{v}_1, \dots, \mathbf{v}_N]^\top \in U_{\mathbf{a}}$,

$$\bar{f}(\mathbf{v}_1, \dots, \mathbf{v}_N) = f(\mathbf{a}). \quad (71)$$

572 Without loss of generality, we may assume that \bar{f} is permutation-invariant. That is, for every
 573 permutation $\sigma \in S_N$ and every $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_N]^\top \in \mathbb{R}^{N \times d_0}$,

$$\bar{f}(\mathbf{v}_1, \dots, \mathbf{v}_N) = \bar{f}(\mathbf{v}_{\sigma(1)}, \dots, \mathbf{v}_{\sigma(N)}). \quad (72)$$

574 Indeed, if \bar{f} is not permutation-invariant, we can define

$$\tilde{f}(\mathbf{v}_1, \dots, \mathbf{v}_N) = \frac{1}{|S_N|} \sum_{\sigma \in S_N} \bar{f}(\mathbf{v}_{\sigma(1)}, \dots, \mathbf{v}_{\sigma(N)}). \quad (73)$$

575 Then \tilde{f} is continuous and permutation-invariant. Moreover, for every sequence $\mathbf{a} \in A^*$ and every
 576 $\mathbf{V} \in U_{\mathbf{a}}$,

$$\tilde{f}(\mathbf{V}) = \bar{f}(\mathbf{V}), \quad (74)$$

577 since $U_{\mathbf{a}}$ is closed under permutations.

578 Because \bar{f} is continuous and permutation-invariant, it follows from Zaheer et al. [56] that there exists
 579 a pair of functions $\phi : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_1}$ and $\rho : \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_{\text{out}}}$ such that

$$\bar{f}(\mathbf{v}_1, \dots, \mathbf{v}_N) = \rho \left(\sum_{i=1}^N \phi(\mathbf{v}_i) \right) \quad (75)$$

580 (Note that some functions require $d_1 = O(N)$). We further assume without loss of generality that
 581 $\phi(\mathbf{v}_\emptyset) = \mathbf{0}$. We now construct a compressible transformer $M = (\text{emb}, \text{Block}_1, \text{Block}_2)$ approxim-
 582 ating f using this decomposition. The construction of this transformer closely resembles the one used
 583 in the proof of Proposition 3.3.

584 For the first block Block_1 , we set all attention projection matrices to zero,

$$\mathbf{W}_Q = \mathbf{W}_K = \mathbf{W}_V = \mathbf{W}_O = \mathbf{0}, \quad (76)$$

585 so that the block reduces to its feedforward component.

586 Let $\tilde{\phi} : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{4d_0}$ be a continuous function satisfying

$$\tilde{\phi}(\text{emb}(a, i)) = \begin{bmatrix} \phi(\text{emb}(a, i)) \\ \mathbf{p}_i \\ \mathbf{0} \\ \mathbf{0} \end{bmatrix}. \quad (77)$$

587 for all $a \in A$ and $i \in [N]$ where of the four vector blocks above are of size d_0 . That is, $\tilde{\phi}$ computes ϕ
 588 on the embedding while preserving positional information and padding with additional zeroes. We
 589 choose the feedforward map FFN_1 in this block to approximate $\tilde{\phi}$ to accuracy $\varepsilon' > 0$, to be specified
 590 later (This is possible as 2-layer feed forward networks are universal approximators of continuous
 591 functions on compact sets [14; 21]). Thus, for any sequence $\mathbf{a} = (a_1, \dots, a_n)$,

$$\text{Block}_1(\text{emb}(\mathbf{a})) \approx \begin{bmatrix} \tilde{\phi}(\text{emb}(a_1, 1)) \\ \vdots \\ \tilde{\phi}(\text{emb}(a_n, n)) \end{bmatrix}. \quad (78)$$

592 For the second block Block_2 , we again use a single attention head, and set

$$\mathbf{W}_Q = \mathbf{W}_K = \mathbf{0}, \mathbf{W}_O = \mathbf{I} \quad (79)$$

593 The value projection \mathbf{W}_V is chosen to extract the ϕ -component from the output of $\tilde{\phi}$, and copy it to
 594 the third d_1 -sized vector block, keep the first $2d_1$ coordinates the same, and zero out the remaining
 595 coordinates. That is:

$$\mathbf{W}_V = \begin{pmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}. \quad (80)$$

596 Since $\mathbf{W}_Q = \mathbf{W}_K = \mathbf{0}$, all attention scores are identical, and hence the attention weights are uniform.
 597 Applying the attention update in Block_2 to the outputs of Block_1 and using a residual connection,
 598 the representation at position i becomes approximately

$$\begin{bmatrix} \tilde{\phi}(\text{emb}(a_i, i)) \\ \mathbf{p}_i \\ \frac{1}{n} \sum_{j=1}^n \phi(\text{emb}(a_j, j)) \\ \mathbf{0} \end{bmatrix}. \quad (81)$$

599 We now choose the feedforward network of Block_2 to approximate up to ε'' (to be chosen later) a
 600 continuous function $\tilde{\rho}$ satisfying

$$\tilde{\rho}(\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}) = \begin{cases} \rho(j \cdot \mathbf{y}) & \text{if } \mathbf{z} = \mathbf{0} \text{ and } \mathbf{x} = \mathbf{p}_j. \\ \rho((j - i + 1) \cdot \mathbf{y}) & \text{if } \mathbf{z} = \frac{\mathbf{p}_i}{k+1} \text{ and } \mathbf{x} = \mathbf{p}_j. \end{cases} \quad (82)$$

601 Intuitively, this map recovers the position index i from the positional embedding \mathbf{p}_i , rescales the
 602 averaged sum as necessary to recover the unweighted sum, and applies ρ . In particular, at the final
 603 position n , the output of M is approximately

$$\rho \left(n \cdot \frac{1}{n} \sum_{j=1}^n \phi(\text{emb}(a_j, j)) \right) = \rho \left(\sum_{j=1}^n \phi(\text{emb}(a_j, j)) \right) = f(\mathbf{a}). \quad (83)$$

604 By choosing $\varepsilon', \varepsilon''$ sufficiently small, it follows that M approximates f to arbitrary precision.

605 We now show that there exists a compression policy with $r(n) = 1$ that satisfies Equation 66.

606 Define $\mathbf{C} = (\mathbf{c}_1, \mathbf{c}_2)$ as follows. For any sequence of key-value pairs (\mathbf{K}, \mathbf{V}) with $\mathbf{V} =$
 607 $[\mathbf{v}_1, \dots, \mathbf{v}_n]^\top$,

$$\mathbf{c}_1(\mathbf{K}, \mathbf{V}) = (\mathbf{0}, \mathbf{0}), \quad (84)$$

608 and

$$\mathbf{c}_2(\mathbf{K}, \mathbf{V}) = (\tilde{\mathbf{K}}, \tilde{\mathbf{V}}), \quad (85)$$

609 where $\tilde{\mathbf{K}} = \mathbf{0}$ and

$$\tilde{\mathbf{V}} = \begin{pmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \sum_{i=1}^n \mathbf{v}_i + \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{I} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix} \mathbf{v}_n. \quad (86)$$

610 clearly $r(n) = 1$.

611 For a prefix \mathbf{a} of length n and a suffix \mathbf{b} of length k , we begin by analyzing the forward pass of
 612 $\text{M}_{\mathbf{C}, \mathbf{a}}(\mathbf{b})$. After the first block, the representations of the tokens b_i are given by

$$\begin{bmatrix} \tilde{\phi}(\text{emb}(b_1, 1+n)) \\ \vdots \\ \tilde{\phi}(\text{emb}(b_k, k+n)) \end{bmatrix}. \quad (87)$$

613 This follows because the embedding layer of $M_{C,\mathbf{a}}$ is identical to that of M , applied to the concatenated
614 sequence $[\mathbf{a}, \mathbf{b}]$. Moreover, the computation of Block_1 is applied independently to each token
615 representation, and is therefore unaffected by compression. Next, consider the attention computation
616 in the second block. Since $\mathbf{W}_Q = 0$, all attention scores are zero, and hence the attention weights are
617 uniform, equal to $\frac{1}{1+k}$. By the definition of c_2 , it follows that after the attention update in Block_2 , the
618 representation of each token b_i is

$$\begin{bmatrix} \phi(\text{emb}(b_i, n+i)) \\ \mathbf{e}_{n+i} \\ \frac{1}{1+k} \left(\sum_{j=1}^n \phi(\text{emb}(a_j, j)) + \sum_{j=1}^k \phi(\text{emb}(b_j, n+j)) \right) \\ \frac{1}{k+1} \mathbf{p}_n \end{bmatrix}. \quad (88)$$

619 Finally, after applying the feed forward network FFN, by choosing $\varepsilon', \varepsilon''$ small enough, Equation 82
620 implies that the final output of $M_{C,\mathbf{a}}(\mathbf{b})$ is approximately

$$\rho \left((k+n-n+1) \cdot \frac{1}{k+1} \left(\sum_{j=1}^n \phi(\text{emb}(a_j, j)) + \sum_{j=1}^k \phi(\text{emb}(b_j, n+j)) \right) \right) \approx f([\mathbf{a}, \mathbf{b}]) \quad (89)$$

621 and is in an ε approximator of $M_{C,\mathbf{a}}(\mathbf{b})$, completing the proof. Finally, after applying the feedforward
622 network FFN, and choosing $\varepsilon', \varepsilon''$ sufficiently small, Equation 82 implies that the final output of
623 $M_{C,\mathbf{a}}(\mathbf{b})$ is approximately

624 and thus constitutes an ε -approximation of $M_{C,\mathbf{a}}(\mathbf{b})$, completing the proof. □

625

626 **Lemma B.8.** *Let A be a finite alphabet, and let $f : \bigcup_{n \leq N} A^n \rightarrow \mathbb{R}^{\text{d}_{\text{out}}}$ be any sequence-to-vector*
627 *function. Assume additionally that for some prefix sequence $\bar{\mathbf{a}}$ of length $n < N$ there exists two suffix*
628 *sequences $\mathbf{b}^1, \mathbf{b}^2$ both of length $k \leq N - n$ such that*

$$f([\mathbf{a}, \mathbf{b}^1]) \neq f([\bar{\mathbf{a}}, \mathbf{b}^2]) \quad (90)$$

629 *Then for every $\varepsilon, C > 0$ There exists a transformer such that*

630 1. (**Approximation**) *For every sequence $\mathbf{a} = (a_1, \dots, a_n) \in A^n$ with $n \leq N$,*

$$\|f(\mathbf{a}) - M(\mathbf{a})\| < \varepsilon. \quad (91)$$

631 2. (**Non-compressibility**) *For every KV cache compression policy C with compression budget*
632 *satisfying $r(n) < n$*

633 *there exists a suffix $\mathbf{b} \in A^k$ with $k+n \leq N$ such that,*

$$\|M([\mathbf{a}, \mathbf{b}]) - M_{C,\mathbf{a}}(\mathbf{b})\| > C. \quad (92)$$

634 *Proof.* By the same argument as in the proof of Lemma B.7, there exist functions $\phi : \mathbb{R}^{\text{d}_0} \rightarrow \mathbb{R}^{\text{d}_1}$
635 and $\rho : \mathbb{R}^{\text{d}_1} \rightarrow \mathbb{R}^{\text{d}_{\text{out}}}$ such that for every $\mathbf{a} = (a_1, \dots, a_n)$,

$$f(\mathbf{a}) = \rho \left(\sum_{i=1}^n \phi(\text{emb}(a_i, i)) \right). \quad (93)$$

636 Let $\bar{\mathbf{a}}$ be the prefix from the theorem statement, and let $\mathbf{b}^1, \mathbf{b}^2 \in A^k$ be suffixes such that

$$f([\bar{\mathbf{a}}, \mathbf{b}^1]) \neq f([\bar{\mathbf{a}}, \mathbf{b}^2]). \quad (94)$$

637 It follows from equations 90 and 93 that

$$\sum_{i=1}^k \phi(\text{emb}(b_i^1, i)) \neq \sum_{i=1}^k \phi(\text{emb}(b_i^2, i)), \quad (95)$$

638 where $\mathbf{b}_1 = (b_1^1, \dots, b_k^1)$ and $\mathbf{b}_2 = (b_1^2, \dots, b_k^2)$.

639 Define

$$H = \text{conv} \left(\left\{ \sum_{i=1}^k \phi(\text{emb}(b_i, i)) \mid (b_1, \dots, b_k) \in A^k \right\} \right) + \sum_{i=1}^n \phi(\text{emb}(a_i, i)), \quad (96)$$

640 and let D denote the diameter of H . Since ρ is continuous and H is compact, the image $\rho(H)$ is also
641 compact. Consequently, there exists $\mathbf{u}_0 \in \mathbb{R}^{d_{\text{out}}}$ such that

$$\text{dist}(\mathbf{u}_0, \rho(H)) > C. \quad (97)$$

642 We construct the transformer M similarly to the one constructed in the proof of Lemma B.7, with two
643 modifications to the second block.

644 First, we replace $\mathbf{W}_O = \mathbf{I}$ with

$$\mathbf{W}_O = \begin{pmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{pmatrix}. \quad (98)$$

645 Second, we modify the second feedforward function to approximate a continuous function $\bar{\rho}$ satisfying

$$\bar{\rho}(\mathbf{w}, \mathbf{x}, \mathbf{y}, \mathbf{z}) = \begin{cases} \rho(i \cdot \mathbf{y}), & \text{if } \|\mathbf{x} - \mathbf{p}_i\| < \varepsilon' \text{ and } i \cdot \mathbf{y} \in H_{\frac{\delta}{2}}, \\ \mathbf{u}_0, & \text{if } \|\mathbf{x} - \mathbf{p}_i\| < \varepsilon' \text{ and } i \cdot \mathbf{y} \notin H_{\delta}, \end{cases} \quad (99)$$

646 where

$$\delta = \left(\frac{n+k}{n-1+k} - 1 \right) \frac{D}{2}, \quad (100)$$

647 and

$$H_{\varepsilon} := \{\mathbf{z} \in \mathbb{R}^d : \text{dist}(\mathbf{z}, H) < \varepsilon\}. \quad (101)$$

648 By the same reasoning as in Lemma B.7, Equation 91 holds.

649 Let $\mathbf{C} = (\mathbf{c}_1, \mathbf{c}_2)$ be any compression policy with compression budget $r(n) < n$. Since the first block
650 of M zeroes out all attention heads and depends only on the token embeddings, \mathbf{c}_1 has no effect on
651 the forward pass of $M_{\mathbf{C}, \mathbf{a}}(\mathbf{b})$. Thus, after the first block, the representations of the tokens b_i are

$$\begin{bmatrix} \tilde{\phi}(\text{emb}(b_1, 1+n)) \\ \vdots \\ \tilde{\phi}(\text{emb}(b_k, k+n)) \end{bmatrix}. \quad (102)$$

652 In the second block, since $\mathbf{W}_Q = \mathbf{0}$, all attention scores are identical, and hence the attention weights
653 are uniform, equal to $\frac{1}{k+r(n)}$. It follows that the representation of each token b_i after the second block
654 is

$$\bar{\rho} \left(\phi(\text{emb}(\mathbf{b}_i, i+n)), \mathbf{p}_{i+n}, \frac{1}{r(n)+k} \left(\tilde{\mathbf{v}} + \sum_{j=1}^k \phi(\text{emb}(b_j, j+n)) \right), \mathbf{0} \right), \quad (103)$$

655 where the compressed prefix cache is given by

$$\mathbf{c}_2(\mathbf{K}_{\mathbf{a}}, \mathbf{V}_{\mathbf{a}}) = (\tilde{\mathbf{K}}_{\mathbf{a}}, \tilde{\mathbf{V}}_{\mathbf{a}}) = ([\tilde{\mathbf{k}}_1, \dots, \tilde{\mathbf{k}}_{r(n)}]^\top, [\tilde{\mathbf{v}}_1, \dots, \tilde{\mathbf{v}}_{r(n)}]^\top) \quad (104)$$

656 and

$$\tilde{\mathbf{v}} = \left(\sum_{i=1}^{r(n)} \tilde{\mathbf{v}}_i \right)_{2d_1:3d_1}. \quad (105)$$

657 Since $\frac{n+k}{r(n)+k} > 1$ by Lemma B.9, there exists a sequence $\bar{\mathbf{b}}$ of length k such that

$$\text{dist}\left(\frac{n+k}{r(n)+k}\left(\mathbf{v} + \sum_{j=1}^k \phi(\text{emb}(\bar{\mathbf{b}}_j, j+n))\right), H\right) > \left(\frac{n+k}{r(n)+k} - 1\right) \frac{D}{2} \geq \delta. \quad (106)$$

658 Therefore, by Equation 99, the final representation of the token $\bar{\mathbf{b}}_k$, which is the output of $M_{C,\bar{\mathbf{a}}}(\bar{\mathbf{b}})$, is
659 equal to \mathbf{u}_0 .

660 On the other hand, by construction, $M([\bar{\mathbf{a}}, \bar{\mathbf{b}}]) \in \rho(H)$. Hence,

$$\|M([\bar{\mathbf{a}}, \bar{\mathbf{b}}]) - M_{C,\bar{\mathbf{a}}}(\bar{\mathbf{b}})\| > C, \quad (107)$$

661 completing the proof.

662

□

663 **Lemma B.9.** Let $\mathbf{x}_1, \dots, \mathbf{x}_k \in \mathbb{R}^d$ be distinct points with $k \geq 2$, and let

$$H := \text{conv}\{\mathbf{x}_1, \dots, \mathbf{x}_k\}.$$

664 Let $\alpha > 1$. Then for every $\mathbf{v} \in \mathbb{R}^d$ there exists $i \in [k]$ satisfying

$$\text{dist}(\alpha \mathbf{x}_i + \mathbf{v}, H) \geq \frac{\alpha - 1}{2} \text{diam}(H).$$

665 *Proof.* Let

$$D := \text{diam}(H) = \sup_{\mathbf{p}, \mathbf{q} \in H} \|\mathbf{p} - \mathbf{q}\|.$$

666 Since the points $\mathbf{x}_1, \dots, \mathbf{x}_k$ are distinct and $k \geq 2$, we have $D > 0$.

667 Choose $\mathbf{p}, \mathbf{q} \in H$ such that $\|\mathbf{p} - \mathbf{q}\| = D$ (such a pair exists since H is compact), and define

$$\mathbf{u} := \frac{\mathbf{p} - \mathbf{q}}{\|\mathbf{p} - \mathbf{q}\|}.$$

668 The set $\{\langle \mathbf{u}, \mathbf{x} \rangle \mid \mathbf{x} \in H\}$ forms a closed interval, and the width of H in direction \mathbf{u} is

$$\max_{\mathbf{x} \in H} \langle \mathbf{u}, \mathbf{x} \rangle - \min_{\mathbf{x} \in H} \langle \mathbf{u}, \mathbf{x} \rangle = D.$$

669 For any $\mathbf{v} \in \mathbb{R}^d$, the width of $\alpha H + \mathbf{v}$ in direction \mathbf{u} is therefore

$$\max_{\mathbf{x} \in \alpha H + \mathbf{v}} \langle \mathbf{u}, \mathbf{x} \rangle - \min_{\mathbf{x} \in \alpha H + \mathbf{v}} \langle \mathbf{u}, \mathbf{x} \rangle = \alpha D.$$

670 Let

$$\varepsilon := \frac{\alpha - 1}{2} D.$$

671 Suppose, for contradiction, that

$$\text{dist}(\alpha \mathbf{x}_i + \mathbf{v}, H) < \varepsilon \quad \text{for all } i = 1, \dots, k.$$

672 Since

$$\alpha H + \mathbf{v} = \text{conv}\{\alpha \mathbf{x}_1 + \mathbf{v}, \dots, \alpha \mathbf{x}_k + \mathbf{v}\},$$

673 and since the ε -neighborhood of H ,

$$H_\varepsilon := \{\mathbf{z} \in \mathbb{R}^d : \text{dist}(\mathbf{z}, H) < \varepsilon\},$$

674 is convex, it follows that

$$\alpha H + \mathbf{v} \subseteq H_\varepsilon.$$

675 Consequently, the width of $\alpha H + \mathbf{v}$ in direction \mathbf{u} is strictly smaller than the width of H_ε in direction
676 \mathbf{u} . This width is at most

$$D + 2\varepsilon.$$

677 Therefore,

$$\alpha D < D + 2\varepsilon.$$

678 But by the definition of ε ,

$$D + 2\varepsilon = D + (\alpha - 1)D = \alpha D,$$

679 which is a contradiction. Hence there exists some $i \in \{1, \dots, k\}$ such that

$$\text{dist}(\alpha \mathbf{x}_i + \mathbf{v}, H) \geq \varepsilon.$$

680

□

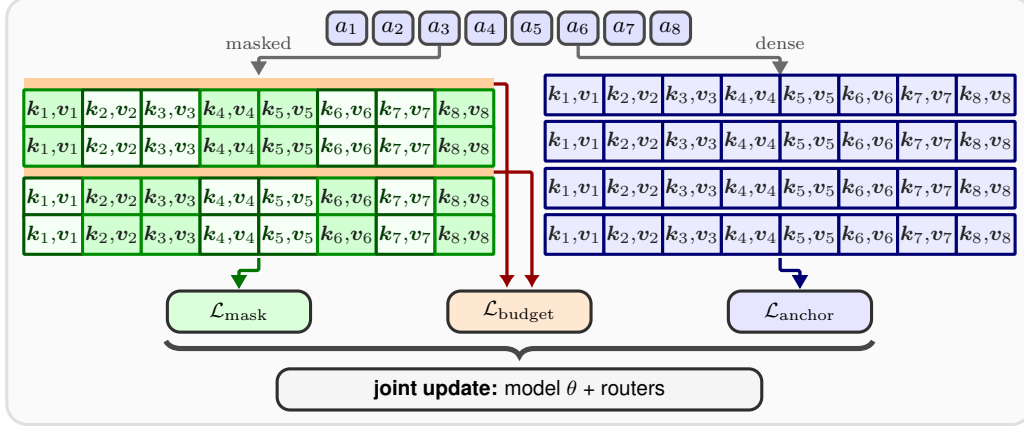


Figure 1: **KV-Compression Aware Training (KV-CAT)**. A context a goes through both masked (left) and dense (right) forward passes. In the masked forward pass, routers (orange) compute masks for groups of consecutive layers, marking KV slots as active (green) or inactive (muted green). In the dense forward pass (blue), all KV slots are kept. The output of the masked forward pass is used to compute $\mathcal{L}_{\text{mask}}$, router distributions are used to compute $\mathcal{L}_{\text{budget}}$, and the outputs of the dense forward pass are used to compute $\mathcal{L}_{\text{anchor}}$. These are jointly used to update the parameters.

681 **Note.** We believe that a finer-grained analysis of transformer compressibility may be possible through
 682 the theory of Chebyshev systems [27], a classical framework with broad applications across both
 683 pure and applied mathematics [26; 38; 15; 28]. Chebyshev systems provide tools for controlling the
 684 zeros, sign changes, and interpolation structure of linear combinations of functions, which makes
 685 them a natural candidate for studying the degrees of freedom available in attention-like mixtures of
 686 value functions. We leave a systematic development of this connection to future work.

687 C KV-CAT implementation details

688 **Masked attention forward pass.** At each layer, the train-time KV sparsification policy defines
 689 which previous KV slots are visible to attention. For layer ℓ , let $m_j^{(\ell)}$ be the most recent mask
 690 produced by the sparsification policy. The attention for query token t is evaluated over

$$\mathcal{A}_t^\ell = \{j \leq t : m_j^{(\ell)} = 1\} \quad (108)$$

691 Thus all valid queries are still processed by the transformer, but their access to past KV slots is
 692 restricted to the active, unmasked KV slots. The masks are independent across compression points: a
 693 token dropped in one layer may be selected again in a later layer. No masking is used at evaluation
 694 time in our comparisons.

695 **Routing mechanism.** Starting from a pretrained decoder-only transformer, we insert learned *routers*
 696 between consecutive layers and train them jointly with the model. At layer ℓ , the router takes as
 697 input the token representations from layer $\ell - 1$ and outputs a scalar score in $[0, 1]$ for each token,
 698 indicating its importance. These scores are thresholded using a (non-trainable) hyperparameter τ to
 699 produce a binary mask: tokens with scores above τ are active in attention, while the rest are masked
 700 out. For efficiency, routers are implemented as linear attention modules. Routers are initialized with
 701 all tokens active (i.e., all scores are set to 1), so training begins from the standard dense transformer
 702 and gradually learns to mask tokens, using a budget loss (defined below). To further reduce overhead,
 703 routers are shared across groups of consecutive layers.

704 **Router parameterization.** For the main experiments in the paper, we use lightweight learnable
 705 router modules, parameterized as linear attention, as our KV sparsification policy. For a sequence
 706 $x_{1:T}$, let $h_t^\ell \in \mathbb{R}^d$ denote the hidden state of token t before decoder layer ℓ . We insert routers at a
 707 small set of compression layers \mathcal{C} (e.g., for QWEN2.5-0.5B we use $\mathcal{C} = \{1, 7, 13, 19\}$). Each router
 708 independently predicts which KV slots are masked for the following group of layers. Given hidden

709 states $H^\ell = (h_1^\ell, \dots, h_T^\ell)$ and, the router first computes

$$\tilde{h}_t = \text{LN}(h_t^\ell), \quad q_t = \phi(W_Q \tilde{h}_t), \quad k_t = \phi(W_K \tilde{h}_t), \quad r_t = W_V \tilde{h}_t, \quad \phi(z) = \text{ELU}(z) + 1. \quad (109)$$

710 It then forms a causal linear-attention summary without storing a router KV cache,

$$S_t = \sum_{j \leq t} k_j r_j^\top, \quad z_t = \sum_{j \leq t} k_j, \quad a_t = W_O \frac{q_t^\top S_t}{q_t^\top z_t + \epsilon}. \quad (110)$$

711 The keep probability is computed from a cosine score between a pointwise projection of the current
712 state and the state after the router summary:

$$u_t = \frac{W_P h_t^\ell}{\|W_P h_t^\ell\|_2}, \quad w_t = \frac{h_t^\ell + \alpha a_t}{\|h_t^\ell + \alpha a_t\|_2}, \quad p_t = \frac{1 - \langle u_t, w_t \rangle}{2}. \quad (111)$$

713 The binary routing decision is

$$m_t^\ell = \mathbf{1}\{p_t > \tau\}, \quad (112)$$

714 with $\tau = 0.5$ in all reported runs. During training, this hard threshold is optimized with a straight-
715 through estimator: the forward pass uses the binary mask above, while the backward pass passes
716 gradients through the pre-threshold keep probability p_t to the router parameters. We initialize
717 $W_P = -I$ and $\alpha = 0$, which gives $p_t = 1$ for every valid token, so training starts from the dense
718 model.

719 **Training objective.** Let $p_\theta^{\text{mask}}(\cdot | x_{<t})$ denote the masked forward pass and $p_\theta^{\text{dense}}(\cdot | x_{<t})$ the
720 same model with masking disabled. We optimize

$$\begin{aligned} \mathcal{L} = & \lambda_{\text{mask}} \frac{1}{T-1} \sum_{t=2}^T D_{\text{KL}}(\text{sg}[p_\theta^{\text{dense}}(\cdot | x_{<t})] \| p_\theta^{\text{mask}}(\cdot | x_{<t})) \\ & + \lambda_{\text{anchor}} \frac{1}{T-1} \sum_{t=2}^T -\log p_\theta^{\text{dense}}(x_t | x_{<t}) + \lambda_{\text{budget}} \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} \mathcal{B}_c. \end{aligned} \quad (113)$$

721 where sg stops gradients through the dense teacher. When using the router based KV sparsification
722 policy, we use $\lambda_{\text{budget}} \neq 0$ and the budget term follows the load-balancing objective from Hwang
723 et al. [22]. For target keep rate ρ , define

$$F_c = \frac{1}{T} \sum_{t=1}^T m_t^c, \quad G_c = \frac{1}{T} \sum_{t=1}^T p_t^c, \quad \mathcal{B}_c = \frac{F_c G_c}{\rho} + \frac{(1-F_c)(1-G_c)}{1-\rho}. \quad (114)$$

724 All transformer and router parameters are updated during continued pretraining; in the reported runs
725 $\rho = 0.5$, $\lambda_{\text{mask}} = 1$, $\lambda_{\text{anchor}} = 1$, and $\lambda_{\text{budget}} = 0.1$.

726 D Main Experiment Overview

727 We evaluate the effect of KV-CAT on the performance of downstream optimization-based KV
728 cache compression techniques. Our empirical study is guided by four central questions: **(Q1)**
729 Does KV-CAT training preserve base model performance in uncompressed settings? **(Q2)**
730 Does KV-CAT improve compression quality under a fixed KV-optimization budget? **(Q3)** Does
731 KV-CAT improve retrieval from a KV-compressed context? **(Q4)** Do these gains transfer to
732 KV-compressed long-context question answering? Across all experiments, we compare the original
733 pretrained model to the model obtained via KV-CAT continued pretraining, applying the same
734 downstream compression method at matched KV retention and optimization budgets. We provide
735 full experimental details in Appendix E, and additional results in Appendix F.

736 **KV-CAT training setup.** We apply KV-CAT to QWEN2.5-0.5B and QWEN2.5-1.5B checkpoints
737 via continued pretraining on FineWeb-Edu [44]. In both cases, we train all model parameters, and
738 introduce four learned routers shared across four layer groups. Each model is trained on a total of
739 5.24×10^9 tokens with a max learning rate of 10^{-4} . See further details in Appendix E.1.

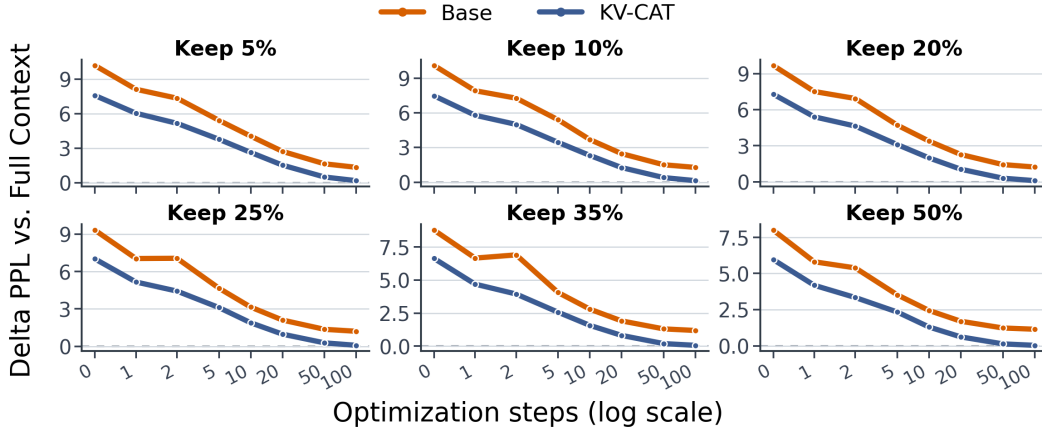


Figure 2: **KV-CAT speeds up gradient-based KV cache compression.** We plot the gap in suffix perplexity under full/compressed-prefix inference throughout gradient-based KV cache optimization. Each panel fixes a different KV keep ratio. Across ratios, the KV-CAT checkpoint achieves a comparable Δ PPL in fewer optimization steps than the base model, yielding up to a $5\times$ speedup.

740 **KV cache compression methods.** Across experiments, we use two post-hoc KV cache compression
 741 procedures: Attention Matching [61] and a gradient-based KV cache optimization method adapted
 742 from Eyuboglu et al. [16]. Attention Matching constructs a compact prefix cache layer-by-layer,
 743 approximating the dense model’s attention traces. The gradient-based method instead directly
 744 optimizes the compact prefix KV cache to match the dense model’s logits on suffix tokens. Both
 745 methods compress the KV cache of a given prefix $\langle \text{text} \rangle$ by leveraging supervision from a suffix,
 746 which together form a *reconstruction sequence*. For the retrieval and QA evaluations, we follow
 747 Zweiger et al. [61], and use a reconstruction sequence of the form $\langle \text{text} \rangle \langle \text{instruction} \rangle \langle \text{text} \rangle$,
 748 where $\langle \text{instruction} \rangle$ is set to “Reproduce the preceding passage verbatim.”, and the loss is applied only on the second repeated $\langle \text{text} \rangle$ span. Additional implementation details are
 749 given in Appendix E.
 750

751 **Preserving uncompressed performance (Q1).** We evaluate the effect of KV-CAT on the per-
 752 formance of the model *with no KV cache compression applied*. We compare the base model and
 753 the KV-CAT-trained checkpoint on six standard multiple-choice benchmarks using the LightEval
 754 harness [19]. As shown in Table 1, KV-CAT closely matches the base model, with an average gain of
 755 0.7 accuracy points for QWEN2.5-0.5B and a drop of 0.5 points for QWEN2.5-1.5B. These results
 756 indicate that KV-CAT preserves standard dense behavior rather than trading it off for compressibility.

757 **Compression under fixed KV-optimization budget (Q2).** We evaluate compression quality
 758 at different keep ratios, under a fixed KV-optimization budget. To do this, we sample held-out
 759 prefix-suffix pairs from FineWeb and directly optimize a compacted prefix KV cache to produce
 760 the corresponding suffix. We optimize the KV slots via Attention Matching using 256 query states
 761 and 8 NNLS iterations across all keep ratios. We report three different similarity metrics between
 762 the suffix next token distribution under the full/compressed prefix, averaged over 128 pairs per keep
 763 ratio. As seen in Table 4, under a fixed optimization budget, the KV-CAT-trained model consistently
 764 outperforms the base model across compression ratios, model sizes, and similarity measures, and
 765 achieves up to $3.21\times$ better retention of suffix perplexity. We additionally evaluate gradient-based
 766 KV cache optimization on the same prefix-suffix pairs. As seen in Figure 2, the KV-CAT-trained
 767 model achieves better suffix perplexity retention throughout optimization, requiring up to $5\times$
 768 fewer optimization steps to match the base model’s performance. For additional details, see Appendix E.3
 769 and for cross-domain evaluation see Appendix F.2.

770 **Retrieval from a compressed context (Q3).**
 771 We conduct a needle-in-a-haystack experi-
 772 ment in which we first optimize a compact
 773 KV cache for the haystack prefix, and then
 774 evaluate exact-match passkey retrieval from
 775 the resulting compressed cache. We compress
 776 the haystack using the gradient-based method

Table 4: **Suffix completion under prefix compression.** Each example uses a 768-token prefix and a 256-token suffix. Prefix KVs are compressed using Attention Matching. We report: (1) Δ PPL, the difference in suffix perplexity under the full prefix vs. compressed prefix; (2) KL, the divergence between the two token distributions; and (3) Top-1, percentage of tokens on which the two distributions’ top-1 agree.

| Model | Keep | Variant | Δ PPL (\downarrow) | KL (\downarrow) | Top-1 (\uparrow) |
|--------------|------|---------|-------------------------------|---------------------|----------------------|
| QWEN2.5-0.5B | 5% | Base | 0.780 | 0.0562 | 88.5% |
| | | KV-CAT | 0.461 | 0.0385 | 91.1% |
| | 10% | Base | 0.662 | 0.0483 | 89.3% |
| | | KV-CAT | 0.422 | 0.0321 | 91.4% |
| | 20% | Base | 0.777 | 0.0549 | 88.7% |
| | | KV-CAT | 0.422 | 0.0321 | 91.4% |

777 described earlier, training the compacted KVs
 778 for 100 steps on the reconstruction sequence.
 779 We find that KV-CAT improves retrieval
 780 accuracy over the base model, particularly
 781 at moderate compression budgets. KV-CAT
 782 increases mean retrieval accuracy by 6.4
 783 points for QWEN2.5-0.5B and 5.2 points for
 784 QWEN2.5-1.5B (Table 3), achieving an 11–
 785 19 point improvement at 30–50 percent keep
 786 ratios. See additional details in Appendix E.4.

787 **Compressed long-context question answer-**
 788 **ing (Q4).** We evaluate on seven LongBench
 789 v2 [2] tasks, where the context of each ques-
 790 tion is compressed using the reconstruction
 791 sequence before model prediction. As shown
 792 in Table 2, compared to the base model, the KV-CAT-trained model achieves better average accuracy
 793 across all KV retention ratios, with an average improvement of up to **39%**. These results demonstrate
 794 that KV-CAT yields consistent downstream gains, even when applied to LLMs using only an NTP
 795 objective and no post-training. See additional details in Appendix E.5.

796 E Extended Experimental Details

797 We now describe additional details for the empirical evaluation described in the previous section.
 798 We describe the continued-pretraining setup, including the data, architecture hyperparameters, and
 799 compression objective, and then give additional details on the post-hoc compression evaluations and
 800 experimental setups used throughout the paper. Unless otherwise stated, evaluations compare the
 801 base model to the KV-CAT checkpoint with masking disabled. All experiments were run on up to 8
 802 H100 GPUs.

Table 5: KV-CAT continued-pretraining hyperparameters for QWEN2.5 checkpoints.

| Setting | QWEN2.5-0.5B | QWEN2.5-1.5B |
|---|------------------------------|------------------------------|
| Base checkpoint | Qwen/Qwen2.5-0.5B | Qwen/Qwen2.5-1.5B |
| Training data | FineWeb-Edu train, streaming | FineWeb-Edu train, streaming |
| Sequence length | 1024 tokens | 1024 tokens |
| Compression layers | 0, 6, 12, 18 | 0, 7, 14, 21 |
| Router feature dimensions | 64 | 64 |
| Train threshold τ_{train} | 0.5 | 0.5 |
| Target keep rate | 50% | 50% |
| Compression loss weight λ_{mask} | 1 | 1 |
| Dense anchor weight λ_{anchor} | 1 | 1 |
| Budget weight λ_{budget} | 0.1 | 0.1 |
| Optimizer | AdamW | AdamW |
| Weight decay | 0.01 | 0.01 |
| Gradient clipping | 1.0 max norm | 1.0 max norm |
| Peak learning rate | 10^{-4} | 10^{-4} |
| Minimum learning rate | 5×10^{-6} | 5×10^{-6} |
| Warmup steps | 600 | 600 |
| Tokens per optimizer step | 131,072 | 131,072 |
| Training hardware | 8×H100 GPUs | 8×H100 GPUs |
| Training steps | 40k | 40k |
| Training tokens | 5.24×10^9 | 5.24×10^9 |

803 E.1 Continued pretraining runs

804 We train two model sizes, initialized from the public QWEN2.5-0.5B and QWEN2.5-1.5B [51]
 805 checkpoints, using continued pretraining on FineWeb-Edu [44] at a context length of 1024. The
 806 compression-aware objective is the one described in the method section: a compressed self-distillation

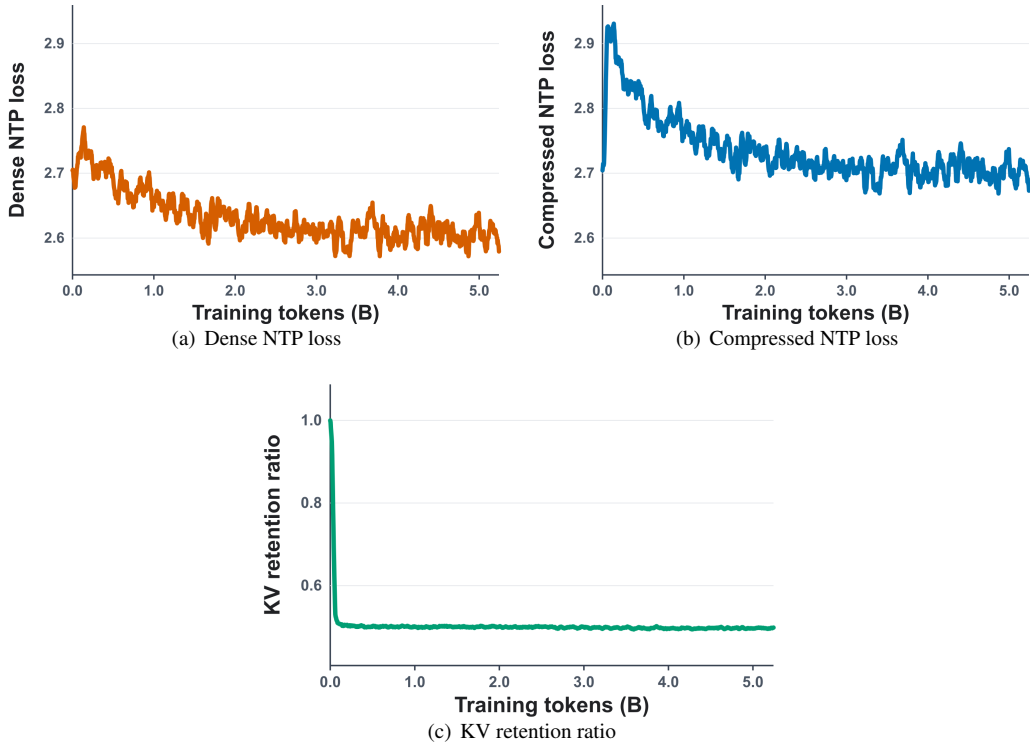


Figure 3: **QWEN2.5-0.5B KV-CAT continued pretraining run.** We plot dense validation next-token-prediction loss, compressed-path validation next-token-prediction loss, and the realized KV retention ratio over the 5.24B-token KV-CAT continued-pretraining run. The retention trace reports the fraction of KV slots kept by the learned routers, whose budget target is 50%.

807 loss matched to the model’s dense distribution, an uncompressed next-token prediction anchor, and a
 808 budget penalty. All transformer and compression module parameters are updated.

809 The router layers are inserted at layers 0, 6, 12, and 18 for QWEN2.5-0.5B and layers 0, 7, 14, and
 810 21 for QWEN2.5-1.5B. Each module is a causal linear-attention boundary predictor (as described
 811 in Appendix C) with 64 feature dimensions, a threshold $\tau = 0.5$, and a target keep rate of 50%.
 812 The objective weights are $\lambda_{\text{mask}} = 1$, $\lambda_{\text{anchor}} = 1$, and $\lambda_{\text{budget}} = 0.1$. We use AdamW with peak
 813 learning rate 10^{-4} , 600 warmup steps, minimum learning rate 5×10^{-6} , weight decay 0.01, and
 814 gradient clipping at norm 1.0. Both model sizes use a batch of 131,072 tokens per optimizer step,
 815 corresponding to 128 sequences of length 1024, and are trained on 8 GPUs. Runs are configured for
 816 40k optimizer steps, which see 5.24×10^9 tokens.

817 At evaluation time, the boundary predictors are turned off, and the post-hoc compressor named in
 818 each experiment is the only cache-reduction procedure whose quality is measured. Figure 3 reports
 819 the compressed and uncompressed validation NTP loss as well as the KV retention ratio throughout
 820 the KV-CAT continued pretraining for QWEN2.5-0.5B. Table 5 reports all hyperparameters and
 821 training setup for both runs.

822 E.2 No compression QA evaluation

823 To check that our training procedure does not substantially degrade ordinary model behavior, we
 824 evaluate normalized multiple-choice accuracy on a variety of question answering benchmarks us-
 825 ing the LightEval harness [19]. The evaluation suite contains HellaSwag [58], WinoGrande [47],
 826 PIQA [4], Social IQa [49], OpenBookQA [39], ARC-Easy and ARC-Challenge [12]. We use 1000
 827 validation examples per task, except for OpenBookQA, whose validation split contains 500 examples
 828 in this setup. The reported average is the unweighted mean over these eight tasks. The purpose of this
 829 evaluation is to measure whether the trained checkpoints maintain language modeling performance

830 on standard short-context multiple-choice tasks before we test their behavior under hoc KV cache
 831 compression.

832 E.3 Suffix perplexity under prefix KV cache compression

833 The suffix-prediction experiments evaluate whether the continued-pretrained checkpoints are easier
 834 targets for post-hoc prefix KV cache compression. The in-domain evaluation data consists of held-out
 835 blocks from FineWeb. Each evaluation example is a 1024-token excerpt split into a 768-token prefix
 836 and a 256-token suffix. We first run the model on the full 1024-token sequence and record the suffix
 837 logits and per-layer attention traces. The first 768 tokens define the prefix cache to be compressed.
 838 The following 256 suffix tokens are never removed or compressed: after a compact prefix cache
 839 has been constructed, the same suffix tokens are appended normally and scored under the original
 840 transformer layers. This protocol is shared by the Attention Matching table and the gradient based
 841 optimization step-curve figures, but the two experiments fit the compact cache differently.

842 **Attention-Matching.** We use the AM-HighestAttnKeys Attention matching variant from Zweiger
 843 et al. [61]. For a keep ratio $\rho \in \{0.05, 0.1, 0.2, 0.4\}$, Attention Matching constructs a compact prefix
 844 cache independently for each layer and key-value head. Let $K, V \in \mathbb{R}^{L \times d}$ denote the dense prefix
 845 keys and values for one key-value head, where $L = 768$, and let $Q \in \mathbb{R}^{M \times d}$ denote suffix query
 846 states from the corresponding grouped query heads. We subsample at most 256 suffix queries per
 847 key-value head. Attention Matching first selects $\lceil \rho L \rceil$ source keys with the largest root-mean-square
 848 attention weight over the sampled suffix queries. The selected keys form the compact keys C_1 . A
 849 scalar log-bias vector β is then fit by two iterations of box-constrained nonnegative least squares
 850 so that the compact cache approximately matches the dense attention normalizer. Finally, compact
 851 values C_2 are fit by ridge least squares to reconstruct the dense attention outputs,

$$\text{softmax}\left(\frac{QC_1^\top}{\sqrt{d_k}} + \beta\right) C_2 \approx \text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right) V.$$

852 We use ridge coefficient 10^{-4} with spectral scaling and solve the value regression with a least-squares
 853 solver. After constructing the compact prefix cache, we run the model on the same prefix-suffix
 854 sequence while replacing the dense prefix cache with the compact cache for all subsequent suffix
 855 positions. We report the increase in suffix perplexity relative to the same model’s dense-prefix forward
 856 pass, KL divergence to the native dense-prefix logits, and top-1 agreement with the native logits. All
 857 metrics are averaged over 128 evaluation examples.

858 **Gradient-based KV cache compression.** In this setting, we directly optimize a small set of
 859 continuous key and value vectors for each example while keeping all model parameters fixed. The
 860 resulting optimization curves measure how much information about the prefix can be preserved in a
 861 compact KV cache when the cache itself is allowed to adapt to the example. The dense model is first
 862 run on the full prefix-suffix sequence to produce teacher logits on the suffix. For a target keep ratio

$$\rho \in \{0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.5\},$$

863 we allocate a compact prefix cache with

$$m = \lceil \rho L \rceil$$

864 KV slots per layer and key-value head, where $L = 768$ is the dense prefix length. The compact cache
 865 consists of learnable key and value tensors

$$\tilde{K}_{\ell,h}, \tilde{V}_{\ell,h} \in \mathbb{R}^{m \times d}$$

866 for each layer ℓ and key-value head h . We initialize $\tilde{K}_{\ell,h}$ and $\tilde{V}_{\ell,h}$ from the first m dense prefix KV
 867 states for that layer and head. The compact cache is optimized per example using the suffix teacher
 868 distribution from the dense full-prefix forward pass. Let z_t^* denote the dense teacher logits at suffix
 869 position t , and let $z_t(\tilde{K}, \tilde{V})$ denote the logits obtained when the dense prefix cache is replaced by the
 870 compact cache. We minimize

$$\mathcal{L}_{\text{KV}} = \frac{1}{|\mathcal{S}|} \sum_{t \in \mathcal{S}} D_{\text{KL}}\left(\text{softmax}(z_t^*) \parallel \text{softmax}(z_t(\tilde{K}, \tilde{V}))\right),$$

871 where \mathcal{S} is the set of suffix prediction positions. The optimization updates only the compact KV
872 tensors; model weights are never updated. We use AdamW for 100 gradient steps with learning rate
873 10^{-2} , weight decay 0, and gradient clipping at norm 1.0. We record intermediate compact caches
874 after steps

$$\{0, 1, 2, 5, 10, 20, 50, 100\}$$

875 to produce the optimization-step curves. After each recorded optimization step, we evaluate next-
876 token prediction on the same suffix using the compact cache in place of the dense prefix cache. The
877 reported perplexity is computed against the ground-truth suffix tokens.

878 This evaluation should be interpreted as a controlled probe of the cache-construction problem in
879 post-hoc compression. For each example, the compact prefix cache is fit using supervision derived
880 from the same prefix–suffix sequence on which it is evaluated: Attention Matching uses suffix
881 attention/query traces to construct the compact cache, while the gradient-based method uses suffix
882 teacher logits. Thus, the comparison asks whether the trained checkpoints make the per-example
883 compact-cache optimization problem easier under matched cache budgets and matched supervision,
884 rather than testing transfer from a cache fit on one suffix to unseen suffixes (which is tested in the
885 next experiments).

886 E.4 Needle-in-a-haystack retrieval under KV cache compression

887 We evaluate retrieval under prefix KV cache compression using a deterministic needle-in-a-haystack
888 task. Each example contains neutral filler text, passkey-like distractor numbers, a single six-digit
889 passkey, and a final query asking for the passkey. The prompt length before the final query is 1024
890 tokens. The model is evaluated using the compressed prompt and the query.

891 **Example construction.** To generate each example, we sample a six-digit string $a \in$
892 $\{000000, \dots, 999999\}$, which defines the passkey stored in the prefix and the answer expected
893 at evaluation. For example, if $a = 483920$, the needle inserted into the haystack is

Memory record: special_passkey=483920.

894 The final query is

Memory record: special_passkey=

895 and the answer is exactly

483920.

896 The passkey is inserted at relative depths

$$d \in \{0, 0.25, 0.5, 0.75, 1\},$$

897 where $d = 0.0$ places the needle at the beginning of the haystack and $d = 1$ places it at the end.
898 In addition to the answer passkey, examples contain distractors: passkey-like six-digit strings that
899 are explicitly not equal to the answer. These distractors are inserted into the sequence so that the
900 model cannot solve the task by copying an arbitrary number. For example, if the answer is 483920, a
901 possible distractor sentences is or

Ignore the misleading special passkey candidate 672041.

902 During example construction, the haystack is assembled from filler text before and after the needle.
903 At each filler-sentence draw, with probability 0.35 we insert a distractor sentence; otherwise we insert
904 a neutral sentence unrelated to the retrieval key. Distractor codes are sampled uniformly from the
905 six-digit code space and rejected if they match the answer. The needle is then inserted at the requested
906 relative depth in the haystack, and the query is appended after the haystack. For each compression
907 ratio, we evaluate 20 instances per depth totaling at a 100 examples.

908 **Training the compact KV cache.** For each example, we train a new compact prefix KV cache
909 while keeping all model weights frozen. Let $\langle \text{haystack} \rangle$ denote the haystack prefix, including the
910 filler text, needle and any distractors, but excluding the final query and answer. We construct an
911 reconstruction sequence

$\langle \text{haystack} \rangle \langle \text{instruction} \rangle \langle \text{haystack} \rangle,$

912 where the instruction is

Reproduce the preceding passage verbatim.

913

Preserve every word, digit, and punctuation mark.

914 The compact cache replaces the KV states for the first `<haystack>`. The loss is applied only on
915 the second `<haystack>`, so the optimizer must store enough information in the compact cache to
916 reconstruct the original prefix, including the passkey and distractors.

917 For a keep ratio ρ , the compact cache contains

$$m = \lceil \rho L \rceil$$

918 KV slots per layer and KV head, where L is the haystack length. The compact keys and values are
919 initialized by selecting m random haystack KV states and then optimized directly with AdamW. As
920 in Appendix E.3, we use a KL objective between the suffix logits (computed on the second copy of
921 `<haystack>`) produced by the dense model and the ones produced when attending to the compact
922 cache. For each example, this objective is optimized by running AdamW for 100 steps with learning
923 rate 10^{-2} , weight decay 0, and gradient clipping at norm 1.0.

924 **Evaluation.** At evaluation time, the haystack prefix is represented by the optimized compact KV
925 cache, while the final query is appended normally and remains uncompressed. We then perform
926 greedy answer generation and compare the generated completion against the six-digit passkey. In
927 the table, we report native-success-conditioned exact-match retrieval accuracy: for each model, keep
928 ratio, and compact-cache optimization step, we report the percentage of those examples for which the
929 compressed-prefix model also generates exactly the passkey.

930 E.5 Long-form question answering under KV cache compression

931 We evaluate long-context question answering under post-hoc KV cache compression on LongBench
932 v2. Because the compression procedure optimizes a separate compact KV cache for each prompt,
933 this evaluation is substantially more expensive than ordinary decoding. We therefore report results on
934 a fixed subset of seven LongBench v2 subdomains: Academic, Agent history QA, Knowledge graph
935 reasoning, Legal, Many-shot learning, New language translation, and Table QA. This subset contains
936 221 examples in total.

937 For each LongBench example, we construct a multiple-choice prompt from the context, question, and
938 four answer choices. The prefix is

Please read the following text and answer the question below.

939 followed by the context from the benchmark. The suffix is

What is the correct answer to this question: `<question>`

940 followed by the four choices (A), (B), (C), and (D), and the answer-format prefix

The correct answer is (.

941 We score the four continuations A), B), C), and D) by total log-likelihood and predict the answer
942 with the highest score. No free-form generation is used.

943 For each example, model, and keep ratio $\in \{0.1, 0.2, 0.5\}$, we optimize a compact KV cache for
944 example’s context while keeping all model weights frozen. The optimization uses a reconstruction
945 sequence as in Appendix E.4. As common in LongBench evaluations, if the reconstruction sequence
946 length is larger than the context length of the model (32,768 tokens) we truncate the context using
947 middle truncation (keeping a prefix and suffix of equal length) [2]. The KC objective is a KL
948 divergence between the dense model’s next-token distribution and the compact-cache model’s next-
949 token distribution on all tokens of the repeated prefix. The compressed caches are initialized from
950 the first prefix KV slots. We optimize each compact cache for 100 AdamW steps with learning rate
951 10^{-2} , weight decay 0, gradient clipping norm 1.0, and gradient checkpointing enabled. At evaluation
952 time, the optimized compact KV cache replaces the full KV cache of the context prefix. The question,
953 answer choices, and answer prefix are then processed normally, without compression. We report
954 accuracy, grouped by subdomain and averaged over the 221 examples.

Table 6: KV-CAT implemented with a ROUTER policy retains base model performance better than the fixed RAND and ATTN policies.

| Model | Variant | HellaSwag | WinoGrande | PIQA | OpenBookQA | ARC-E | ARC-C | Avg. |
|--------------|---------|-------------|-------------|-------------|-------------|-------------|-------------|-------------|
| QWEN2.5-0.5B | Base | 54.6 | 52.9 | 70.4 | 38.2 | 60.8 | 32.2 | 51.5 |
| | ROUTER | 53.6 | 54.1 | 72.1 | 37.4 | 63.6 | 32.4 | 52.2 |
| | RAND | 52.6 | 52.9 | 70.2 | 35.4 | 61.8 | 31.1 | 50.7 |
| | ATTN | 50.7 | 53.6 | 70.0 | 35.6 | 60.2 | 30.9 | 50.2 |

955 F Additional Experimental Results

956 This section collects supplementary results that are complementary to Section 5. In Section F.1, we
 957 ablate the choice of the train-time sparsification policy used in KV-CAT. In Section F.2, we test
 958 whether the Attention Matching performance gains we observe in Section 5 generalize to prefix-suffix
 959 from additional text corpora.

960 F.1 Train-time KV sparsification policy ablation

961 We use this ablation to choose the train-time KV sparsification policy used in the main experiments.
 962 Starting from QWEN2.5-0.5B, we train three checkpoints with the same self-distillation and dense-
 963 anchor losses and the same 50% target keep rate, varying only the policy that selects active KV slot
 964 during the masked forward pass. RAND uniformly samples an exact-count 50% subset of valid KV
 965 slots at each sparsification point. ATTN is an H2O-inspired [59] attention-mass baseline: it computes
 966 dense causal attention at the routed layer, sums the attention mass received by each source token
 967 over heads and valid query positions, and keeps the top 50% source KV slots. ROUTER uses the
 968 lightweight learned linear-attention routers described in Appendix C. Because RAND and ATTN
 969 satisfy the target keep rate by construction, we set $\lambda_{\text{budget}} = 0$ for those runs; for ROUTER, we use
 970 the budget regularizer described in Appendix C. In all comparisons below, the trained checkpoints
 971 are evaluated in the unmasked forward pass mode, with the training-time sparsification mechanism
 972 disabled, and post-hoc compression is applied afterward.

973 **Uncompressed performance.** ROUTER is the only sparsification policy in this ablation that pre-
 974 serves dense downstream accuracy at or above the base model average. Its average score is 52.2,
 975 compared to 51.5 for the base model, 50.7 for RAND, and 50.2 for ATTN. The fixed policies still
 976 retain much of the base model’s zero-shot performance after continued pretraining.

977 **Attention Matching Evaluation.** All checkpoints are evaluated with the same Attention Matching
 978 procedure described in Appendix E.3. Table 7 shows that all three train-time sparsification policies
 979 Attention Matching performance compared to the original base checkpoint for every reported metric
 980 and budget. This further supports the KV-CAT framework, demonstrating that the downstream
 981 compression improvements are not tied exclusively to learned router parameterization. Among the
 982 policies, ROUTER gives the best result for every reported keep ratio and metric. The margin over
 983 RAND is sometimes small, especially at 10% and 40% keep ratio, but ROUTER is consistently at least
 984 as good under compression and is clearly stronger on dense task retention. We therefore use ROUTER
 985 in the main experiments as the default train-time sparsification policy.

Table 7: Attention-matching prefix compression on FineWeb-Edu for the base model and three KV-CAT models with the ROUTER, RAND, and ATTN KV sparsification policies. We report degradation relative to each model’s native dense-prefix forward pass. Lower ΔPPL and KL are better; higher top-1 agreement is better.

| Keep ratio | ΔPPL (\downarrow) | | | | KL (\downarrow) | | | | Top-1 (\uparrow) | | | |
|------------|-------------------------------------|--------------|-------|-------|---------------------|---------------|--------|--------|----------------------|--------------|-------|-------|
| | Base | ROUTER | RAND | ATTN | Base | ROUTER | RAND | ATTN | Base | ROUTER | RAND | ATTN |
| 0.05 | 0.780 | 0.461 | 0.547 | 0.593 | 0.0562 | 0.0385 | 0.0420 | 0.0440 | 88.5% | 91.1% | 90.3% | 90.3% |
| 0.10 | 0.662 | 0.422 | 0.427 | 0.495 | 0.0483 | 0.0321 | 0.0356 | 0.0363 | 89.3% | 91.4% | 90.9% | 90.8% |
| 0.20 | 0.777 | 0.438 | 0.547 | 0.594 | 0.0549 | 0.0374 | 0.0409 | 0.0396 | 88.7% | 90.8% | 90.4% | 90.5% |
| 0.40 | 0.592 | 0.360 | 0.361 | 0.386 | 0.0431 | 0.0272 | 0.0283 | 0.0288 | 90.3% | 92.5% | 92.3% | 92.4% |

986 **F.2 Cross-domain attention matching evaluation**

987 We test whether the Attention Matching optimization gains observed for KV-CAT transfer to prefix-
 988 suffix pair from additional textual domains. We evaluate Attention Matching KV cache compression
 989 on validation examples from three held-out corpora: WikiText-103 [37], PG-19 [45], and arXiv [13].
 990 For each example, we take a 768-token prefix and a 256-token suffix, compress only the prefix KV
 991 cache with the same Attention Matching procedure used in Appendix E.3, and then score suffix
 992 next-token prediction against each model’s own dense full-prefix distribution.

993 Across all three held-out domains, the compression-aware checkpoint remains a better target for the
 994 same post-hoc compressor. It reduces KL and increases top-1 agreement for every corpus and keep
 995 ratio. It also lowers Δ PPL in 11 of 12 settings; the only exception is arXiv at 40% keep ratio, where
 996 Δ PPL increases even though KL and top-1 agreement still improve. We therefore interpret this table
 997 as evidence that KV-CAT improves cross-domain compressibility under Attention Matching.

998 **F.3 Additional KV keep-ratio curves**

999 Figure 4 extends the main-body optimization curves in Figure 2 with three additional KV keep ratios:
 1000 15%, 25%, and 35%.

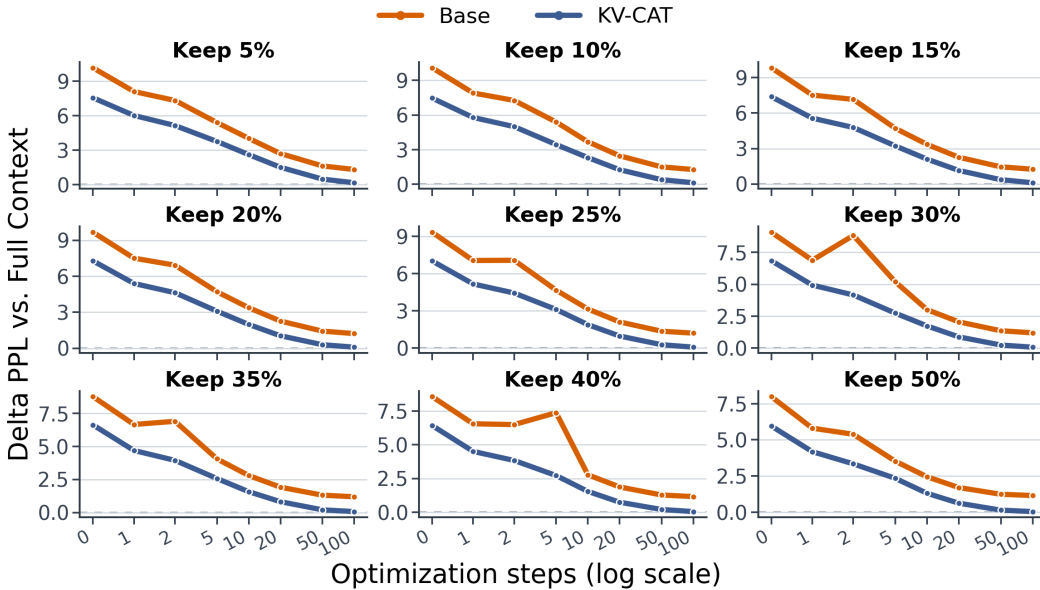


Figure 4: **KV-CAT speeds up gradient-based KV cache compression.** We plot the gap in suffix perplexity under full/compressed-prefix inference throughout gradient-based KV cache optimization. Each panel fixes a different KV keep ratio. Across ratios, the KV-CAT checkpoint achieves a comparable Δ PPL in fewer optimization steps than the base model, yielding up to a $5\times$ speedup.