# KV-DICT: SPARSE KV CACHE COMPRESSION WITH UNIVERSAL DICTIONARIES

Anonymous authors

Paper under double-blind review

#### Abstract

Transformer has become the de facto architecture for Large Language Models (LLMs), yet its substantial memory required for long contexts makes it costly to deploy. Managing the memory usage of the key-value (KV) cache during inference has become a pressing challenge, as the cache grows with both model size and input length, consuming significant GPU memory.

We introduce a novel post-training KV cache compression method using KV-Dict, a universal dictionary that can accurately decompose and reconstruct key-value states. Unlike traditional quantization methods, KV-Dict leverages sparse dictionary learning, allowing for flexible memory usage with minimal performance loss through fine-grained controls of sparsity levels. Moreover, we retain competitive performance in the low memory regimes that 2-bit compression struggles to offer. KV-Dict is remarkably universal, as it uses a small, input-agnostic dictionary that is shared across tasks and batches without scaling memory. This universality, combined with the ability to control sparsity for different memory requirements, offers a flexible and efficient solution to the KV cache bottleneck, maintaining strong performance on complex reasoning tasks, such as LongBench and GSM8K.

029

037

024

004

010 011

012

013

014

015

016

017

018

019

021

#### 1 INTRODUCTION

Transformers (Vaswani et al., 2017) have become the backbone of frontier Large Language Models
 (LLMs), driving progress in domains beyond natural language processing. However, Transformers are often limited by their significant memory demands. This stems not only from the large number of model parameters, but also from the need to maintain the key-value (KV) cache that grows proportional to the model size and length of the input. Additionally, each input requires its own KV cache, limiting opportunities for reuse across different user inputs. This creates a bottleneck in generation speed for GPUs with limited memory (Yu et al., 2022) and thus, it has become crucial to alleviate KV cache memory usage while preserving its original performance across domains.

KV cache optimization research has explored both training-stage optimizations (Shazeer, 2019; Dai et al., 2024; Sun et al., 2024) and deployment-focused strategies (Kwon et al., 2023; Lin et al., 2024; Ye et al., 2024) to improve the efficiency of serving LLMs. Architectural approaches such as Grouped Query Attention (GQA) (Ainslie et al., 2023) reduce the number of KV heads, effectively reducing KV cache. However, these methods are not applicable as off-the-shelf methods to reduce KV cache for pretrained LLMs, leading to post-training compression efforts.

Post-training approaches include selectively retaining certain tokens (Beltagy et al., 2020; Xiao et al., 2023; Zhang et al., 2024) and quantization methods, which have had empirical success when quantizing KV cache into 2 or 4 bits (Liu et al., 2024b; He et al., 2024; Kang et al., 2024). However, eviction strategies have limitations on long-context tasks that require retaining a majority of previous tokens, while quantizations to 2 or 4 bits have clear upper bounds on compression rates.

In this paper, we focus on utilizing low-dimensional structures for efficient KV cache compression. Prior work reports that each key state lies in a low-rank subspace (Singhania et al., 2024; Wang et al., 2024b; Yu et al., 2024). Yet, it is unclear if *all states* lie in the same subspace; if so, such redundancy remains to be taken advantage of. Thus, we naturally ask the following questions:

050

Can we find substantial redundancy among keys and values, universal of all inputs? How can we leverage this for efficient KV cache compression?

069

071

072 073

075

081

083 084

085

087

090

092

093

094

095

096 097

098

099

102

103



Figure 1: (a) Prefilling: Following attention computation, KV-Dict uses OMP to find sparse representations of the KV states  $(3-8 \times \text{ smaller})$ . (b) Decoding: Key cache consists of the compressed sparse key cache,  $K_{csr}$ , and an full-precision buffer,  $K_{buffer}$ , for the most recent tokens.  $q_t$ ,  $k_t$  represent the query, key vectors for the newly generated token. Computation is reduced by computing the query-dictionary product,  $q_t D_k$ , then multiplying  $K_{csr}$ , to get the pre-softmax attention score.

Towards this end, we propose KV-Dict, a universal dictionary that serves as an overcomplete basis, 074 which can sparsely decompose and reconstruct KV cache with sufficiently small reconstruction error. As shown in Section 2.2, we observe that a subset of key states cluster near each other, even 076 though the keys are from different inputs, while some cluster on different subspaces. To accommo-077 date for such low-dimensional structures, we use sparse dictionary learning, which has developed 078 algorithms for information compression across various domains, such as signal processing and med-079 ical imaging (Candès et al., 2006; Donoho, 2006; Dong et al., 2014; Metzler et al., 2016).

KV-Dict is simple to learn, can be applied off-the-shelf for KV cache compression, and only occupies small constant memory regardless of input or batch size. Methodologically, KV-Dict utilizes both sparsity-based compression (steps 1 and 2) and quantization (step 3) in three straightforward steps:

- 1. **Dictionary pretraining:** As in Figure 3, we train a dictionary on WikiText-103 (Merity, 2016) for each model. This dictionary is only trained once and used universally across all tasks. It only occupies constant memory and does not increase with batch size.
- 2. Sparse decomposition: During prefilling and decoding (Figure 1), KV-Dict decomposes key-value into a sparse linear combination, which consists of s pairs of sparse coefficient and dictionary index. This step by itself provides high compression rates.
- 3. Lightweight sparse coefficients: We obtain higher KV cache compression rates by representing the sparse coefficients in 8 bits instead of FP16. Lowering precision to 8 bits vields minimal degradation. KV-Dict theoretically allows us to compress more than 2-bit quantization (1/8 of FP16 KV cache size) if s < 10 when head dimension is 128.
- Overall, we make the following contributions:
  - Near-lossless performance: Given similar memory requirements, KV-Dict performs on par with or better than baseline quantization methods on challenging language tasks, such as LongBench (Bai et al., 2023) and GSM8K (Cobbe et al., 2021).
  - Compression rates beyond 2-bits: KV-Dict's sparsity parameter enables both wider and more fine-grained control over desired memory usage. This allows us to explore performance when using under 15-20% of the original KV cache size, a low-memory regime previous compression methods could not explore.
- 105 • Universality: Instead of an input-dependent dictionary, we find a sufficiently small universal dictionary (per model) that can be used for all tasks and across multiple users. Advan-107 tageously, such dictionary does not scale with batch size and can be used off-the-shelf.

## 108 2 KV CACHE COMPRESSION WITH DICTIONARIES

## 110 2.1 BACKGROUND & NOTATION

112 During autoregressive decoding in Transformer, the key and value states for preceding tokens are 113 independent of subsequent tokens. As a result, these key and value states are cached to avoid recom-114 putation, thereby accelerating the decoding process.

Let the input token embeddings be denoted as  $X \in \mathbb{R}^{l_{seq} \times d}$ , where  $l_{seq}$  and d are the sequence length and model hidden dimension, respectively. For simplicity, we focus on a single layer and express the computation of query, key, and value states at each attention head during the prefilling stage as:

$$Q^{(h)} = XW_q^{(h)}, \quad K^{(h)} = XW_k^{(h)}, \quad V^{(h)} = XW_v^{(h)},$$
 (1)

where  $W_q^{(h)}, W_k^{(h)}, W_v^{(h)} \in \mathbb{R}^{d \times m}$  are the model weights with *m* representing the head dimension.

Let t represent the current step in the autoregressive decoding, and let  $x_t \in \mathbb{R}^{1 \times d}$  denote the embedding of the newly generated token. The KV cache up to but not including the current token, are denoted as  $K_{t-1}^{(h)}$  and  $V_{t-1}^{(h)}$ , respectively. The typical output computation for each attention head  $h_t^{(h)}$  using the KV cache can be expressed as:

$$\boldsymbol{h}_{t}^{(h)} = \operatorname{Softmax}\left(\boldsymbol{q}_{t}^{(h)}\left(\boldsymbol{K}_{t-1}^{(h)} \parallel \boldsymbol{k}_{t}^{(h)}\right)^{\top} / \sqrt{m}\right)\left(\boldsymbol{V}_{t-1}^{(h)} \parallel \boldsymbol{v}_{t}^{(h)}\right),$$
(2)

where  $q_t^{(h)}, k_t^{(h)}, v_t^{(h)}$  represent the query, key, and value vectors for the new token embedding  $x_t$ . Here,  $\|$  denotes concatenation along the sequence length dimension.

#### 2.2 SPARSE APPROXIMATION

Given a dictionary, our goal is to efficiently decompose and represent KV cache, *i.e.*, approximate a state  $k \in \mathbb{R}^m$  as a linear combination of a few vectors (atoms) from an overcomplete dictionary  $D \in \mathbb{R}^{m \times N}$ . This reconstruction is given by k = Dy, where  $y \in \mathbb{R}^N$  is the sparse representation vector such that  $s = ||y||_0$ . For implementation, y only requires space proportional to s, not N.

We hypothesize that the KV cache, like other domains where sparse approximation is effective, contains inherent redundancy that can be leveraged for efficient compression. For instance, Figure 2 presents pairwise cosine similarity plots for keys generated during inference on a random subset of the WikiText dataset. Here, we observe that key states cluster in multiple different subspaces. Dictionary learning can take advantage of such redundancy, enabling KV vectors to be represented by a compact set of atoms with only a few active coefficients.





119

123

124

125

132 133

134

145

146

147

148

149

150

151

152

153

Figure 2: Left shows a pairwise cosine similarity matrix between key vectors generated from one input text from all heads in Layer 10 of Llama-3.1-8B-Instruct. Keys are sorted by similarity to demonstrate the clusters. **Right** shows the similarity matrix between key vectors from two *different* input texts. These plots indicate that there may exist a mixture of low-dimensional subspaces in the space of *all* possible keys, a hypothesis that naturally leads to dictionary learning.



Figure 3: Dictionary Learning. We train a linear layer D (our dictionary), that minimizes  $\ell_2$ reconstruction error of KV states. The states of layer *i* are used as training data for dictionary  $D^{(i)}$ . At each step, we apply OMP with fixed D to represent KV as a vector of sparse coefficients; we then perform a step of gradient descent on D and repeat the process. A sparse vector can be efficiently stored as a compressed sparse row (CSR), using a tuple of 16-bit index and 8-bit value.

Sparse approximation, which aims to find y with minimum sparsity given k and D, while ensuring a small reconstruction error, is NP-hard. This optimization problem is typically formulated as:

185 186

195

181 182

162 163

164

165

166

167

169

170

171

172 173

174

175

 $\min_{\boldsymbol{y}} \|\boldsymbol{y}\|_{0} \text{ subject to } \|\boldsymbol{k} - \boldsymbol{D}\boldsymbol{y}\|_{2} \le \delta \|\boldsymbol{k}\|_{2} \text{ for some relative error threshold } \delta > 0$ (3)

187 In this work, we adopt Orthogonal Matching Pursuit (OMP) as the sparse approximation algorithm. 188 Given an input key or value vector k, a dictionary D, and a target sparsity s, OMP iteratively se-189 lects dictionary atoms to minimize the  $\ell_2$ -reconstruction error, with the process continuing until the 190 specified sparsity s is reached. Our implementation of OMP builds on advanced methods that utilize properties of the Cholesky inverse (Zhu et al., 2020) to optimize performance. Additionally, we 191 incorporate implementation details from Lubonja et al. (2024) for efficient batched GPU execution 192 and extend it to include an extra batch dimension, allowing for parallel processing across multiple 193 dictionaries. The full algorithm is detailed in Appendix A. 194

196 2.3 LEARNING LAYER-SPECIFIC KV-DICT

As shown in Figure 3, we train layer-specific KV dictionaries through direct gradient-based optimization. For a given key or value vector, denoted as  $\mathbf{k} \in \mathbb{R}^m$  and a dictionary  $\mathbf{D} \in \mathbb{R}^{m \times N}$ , the OMP algorithm approximates the sparse representation  $\mathbf{y} \in \mathbb{R}^N$ . This process is parallelized across multiple dictionaries, but for simplicity, we present the notation for a single dictionary. The dictionary training objective minimizes the  $\ell_2$  norm of the reconstruction error, with the loss function  $\mathcal{L} = \|\mathbf{k} - \mathbf{D}\mathbf{y}\|_2^2$ . We enforce unit norm constraints on the dictionary atoms by removing any gradient components parallel to the atoms before applying updates.

211

Training. The dictionaries are trained on KV pairs generated from the WikiText-103 dataset using Adam (Kingma & Ba, 2014) with a learning rate of 0.0001 and a cosine decay schedule over 20 epochs. The dictionaries are initialized with a uniform distribution, following the default initialization method for linear layers in PyTorch. For Llama-3.1-8B-Instruct, with a sparsity of s = 32 and a dictionary size of N = 1024, the training process takes about 2 hours on an NVIDIA A100 GPU. We demonstrate our trained dictionaries reconstruct and generalize better than dictionaries trained using sparse autoencoders (similarly to those from Makhzani & Frey (2013); Bricken et al. (2023)) across several corpora in Table 1. Our method consistently achieves lower relative reconstruction errors, such as  $0.19 \pm 0.05$  on out-of-domain dataset CNN/DailyMail, and this trend is consistent across other datasets.

221 Despite being trained only on WikiText-103, KV-Dict demonstrates a degree of universality: our dic-222 tionaries achieve lower test loss on out-of-domain datasets such as TweetEval than the test loss on 223 WikiText-103 for sparse autoencoders, offering significant compression with minimal reconstruc-224 tion error. In the next subsection, we explore how low  $\ell_2$ -reconstruction loss translates to strong 225 performance preservation in language modeling.

226 227

228

235 236

244

245 246

249

250

257

260

#### 2.4 PREFILLING AND DECODING WITH KV-DICT

During the prefilling stage, each layer generates the KV vectors for the input tokens, as illustrated in Figure 1a. KV-Dict uses full-precision KV vectors for attention computation, which are then passed to subsequent layers. Subsequently, OMP finds the sparse representations of the KV vectors using layer-specific key and value dictionaries,  $D_k$  and  $D_v$ .

The compressed key and value caches are denoted as  $K_{csr}, V_{csr} \in \mathbb{R}^{l_{seq} \times N}$  and replace the fullprecision KV cache. The reconstruction of the KV cache is performed as follows:

$$\hat{\boldsymbol{K}} = \boldsymbol{K}_{\rm csr} \boldsymbol{D}_k^{\top}, \quad \hat{\boldsymbol{V}} = \boldsymbol{V}_{\rm csr} \boldsymbol{D}_v^{\top} \tag{4}$$

Recall that at the *t*-th iteration of autoregressive decoding, each layer receives  $q_t$ ,  $k_t$ , and  $v_t$ , the query, key, and value vectors corresponding to the newly generated token. Similarly to prior work (Liu et al., 2024b; Kang et al., 2024), we find that keeping a small number of recent tokens in full precision improves the generative performance of the model. To achieve this, we introduce a buffer that temporarily stores recent tokens in an uncompressed state. The KV vectors stored in the buffer are denoted as  $K_{buffer}$ ,  $V_{buffer} \in \mathbb{R}^{n_b \times m}$ , where  $n_b$  is the number of KV vectors in the buffer. The key cache up to, but not including the new token at iteration t, is then reconstructed as follows:

$$\hat{\boldsymbol{K}}_{t-1} = \boldsymbol{K}_{csr} \boldsymbol{D}_{k}^{\top} \parallel \boldsymbol{K}_{buffer}$$
(5)

Substituting this reconstruction into the Equation 2, the attention weights for each head  $a_t^{(h)}$  are computed as:

$$\boldsymbol{a}_{t}^{(h)} = \operatorname{Softmax}\left(\boldsymbol{q}_{t}^{(h)}(\boldsymbol{K}_{\operatorname{csr}}^{(h)}\boldsymbol{D}_{k}^{\top} \parallel \boldsymbol{K}_{\operatorname{buffer}}^{(h)} \parallel \boldsymbol{k}_{t}^{(h)})^{\top}/\sqrt{m}\right)$$
(6)

A key implementation is that attention for the compressed sparse key cache and the uncompressed key cache is computed separately. For compressed sparse key cache, we first compute the product  $q_t^{(h)}D_k$  before we multiply  $K_{csr}$ , directly calculating the pre-softmax attention scores for compressed tokens. Attention for the buffer tokens is computed as usual. These scores are then concatenated with softmax to produce the final attention weights (Figure 1b). This process is formalized as:

$$\boldsymbol{a}_{t}^{(h)} = \operatorname{Softmax}\left(\left(\boldsymbol{q}_{t}^{(h)}\boldsymbol{D}_{k}\boldsymbol{K}_{\operatorname{csr}}^{(h)\top} \mid \boldsymbol{q}_{t}^{(h)}(\boldsymbol{K}_{\operatorname{buffer}}^{(h)} \parallel \boldsymbol{k}_{t}^{(h)})^{\top}\right)/\sqrt{m}\right),\tag{7}$$

where | represents concatenation along columns for attention scores.

Table 1: Relative reconstruction errors of different methods when training dictionaries of size 1024 and sparsity s = 32 on WikiText-103. Sparse Autoencoder is a two-layer perceptron with hard topk thresholding as activation (encoder as a linear layer + activation in Figure 3). KV-Dict is directly optimized using OMP as encoder. The KV states are generated from Llama-3.1-8B-Instruct model.

Test Dataset	KV-Dict	Sparse Autoencoder	<b>Random Dictionaries</b>
WikiText-103	$0.17 \pm 0.06$	$0.20\pm0.05$	$0.27\pm0.02$
CNN/DailyMail	$0.19 \pm 0.05$	$0.22\pm0.04$	$0.27\pm0.02$
IMDB	$0.18 \pm 0.05$	$0.22\pm0.05$	$0.27\pm0.02$
TweetEval	$0.18 \pm 0.06$	$0.21\pm0.05$	$0.27\pm0.02$

When the buffer reaches capacity, OMP compresses the KV vectors for the oldest  $n_a$  tokens in the buffer. This process is independent of the attention computation for the newest token and can therefore be performed in parallel.

Time and Space Complexity. The sparse representations are stored in CSR format, with values encoded in FP8 (E4M3), and all indices, including offsets, are stored as int16. Each row in CSR corresponds to a single key or value vector. For a given sparsity level *s*, the memory usage includes: nonzero values (*s* bytes), dictionary indices (2*s* bytes), and the offset array (2 bytes), resulting in a total size of 3s + 2 bytes. For a head dimension of 128, a fully uncompressed vector using FP16 takes 256 bytes, yielding a memory usage of  $\frac{3s+2}{256} \times 100 \approx 1.17s\%$  (*e.g.*, 37.5% for s = 32).

In terms of time complexity, computing  $q_t K_t^{\top}$  for a single head requires  $O(l_{seq}m)$  multiplications. On the other hand,  $q_t D_k K_{csr}^{\top}$  needs  $O(Nm + l_{seq}s)$  multiplications. This means that our computation is particularly well-suited for long-context tasks when  $l_{seq} > m$  where m is anywhere between 1024 and 4096. For short contexts when  $l_{seq} < m$ , our method only adds a small overhead to attention computation in actuality.

3 EXPERIMENTS

**Setup.** We evaluate our method on various models (Llama-3-8B, Llama-3.1-8B-Instruct, Llama-3.2-1B-Instruct, Llama-3.2-3B-Instruct, Mistral-7B-Instruct), using dictionaries trained on WikiText-103, as done in Section 2.3. To assess KV-Dict's effectiveness in memory reduction while maintaining long-context understanding, we conduct experiments on selected tasks from Long-Bench (Bai et al., 2023), following the setup of Liu et al. (2024b). See Appendix B for task details.

Additionally, we evaluate generative performance on complex reasoning tasks, such as 294 GSM8K (Cobbe et al., 2021) with 5-shot prompting and MMLU-Pro Engineering/Law (Wang et al., 295 2024a) with zero-shot chain-of-thought. We compare our method against two kinds of KV cache 296 compression methods: namely, quantization-based compression and eviction-based compression. 297 For quantization-based methods, we evaluate KIVI (Liu et al., 2024b), ZipCache (He et al., 2024), 298 and the Hugging Face implementation for per-token quantization. For eviction-based methods, we 299 evaluate PyramidKV (Cai et al., 2024) and SnapKV (Li et al., 2024). We refer to the 4-bit and 2-bit 300 versions of KIVI as KIVI-4 and KIVI-2, respectively, and denote its quantization group size as q. 301

We report KV size as the average percentage of the compressed cache relative to the full cache at the end of generation. KV-Dict's sparsity s is set to match the KV size of the baseline.

- **Hyperparameter Settings.** For both experiments, KV-Dict uses a dictionary size of N = 4096, a buffer size of  $n_b = 128$ , and an approximation window size  $n_a = 1$ , compressing the oldest token with each new token generated. For KIVI-4 and KIVI-2, we use a quantization group size of g = 32and a buffer size of  $n_b = 128$ , as is tested and recommended in Liu et al. (2024b), for LongBench. For GSM8K and MMLU-Pro, we test for stronger memory savings, so we use g = 64 and  $n_b = 64$ for KIVI.
- 310

286

287 288

289

290

291

292

293

311 312 3.1 EXPERIMENTAL RESULTS

313 **LongBench Results.** Table 2 presents the performance of KV-Dict and KIVI on LongBench tasks. 314 KV-Dict demonstrates better performance than KIVI with similar or even smaller KV sizes. Notably, 315 KV-Dict enables exploration of extremely low memory regimes that KIVI-2 cannot achieve. At a 316 memory usage of just 12.4% KV size, KV-Dict maintains reasonable long-context understanding, 317 with only 5.6% p and 4.4% p performance loss on Llama-3.1-8B-Instruct and on Mistral-7B-Instruct-318 v0.3, respectively, compared to the full cache (FP16). The largest performance loss comes from tasks with the lowest full cache accuracy, Qasper, yet there is almost no loss in simpler tasks, such 319 as TriviaQA. This indicates that difficult tasks that require more complex understanding are much 320 more sensitive to performance loss. Hence, it is important to evaluate on GSM8K, one of the harder 321 natural language reasoning tasks, as we do next. 322

323

**GSM8K Results.** The performance of KV-Dict on GSM8K compared to KIVI is shown in Table 3.

324 Table 2: Experimental results on LongBench. For KV-Dict, we use N = 4096 as the dictionary 325 size and  $n_b = 128$  as the buffer size. For KIVI, we use q = 32 (group size for quantization) and 326  $n_b = 128$  (buffer size). Sparsity level s is set to match average KV size of KIVI, while s = 8corresponds to cache size unattainable by common 2-bit quantizations. Full cache is in FP16. 327

Method	KV Size	Qasper	QMSum	MultiNews	TREC	TriviaQA	SAMSum	LCC	<b>RepoBench-P</b>	Averag
				Llama	-3.1-8B-I	nstruct				
Full Cache	100%	22.54	24.57	27.44	72.5	91.65	43.47	63.15	56.76	50.20
KIVI-4 KV-Dict <sub>s=24</sub>	$33.2\%\ 30.6\%$	$22.83 \\ 21.68$	$23.72 \\ 24.25$	$27.95 \\ 27.20$	$71.0 \\ 72.5$	$90.39 \\ 91.58$	$44.25 \\ 42.93$	$62.93 \\ 62.92$	$55.48 \\ 56.51$	49.7 <b>49</b> .9
KIVI-2 KV-Dict <sub>s=16</sub>	$21.1\% \\ 21.4\%$	$13.77 \\ 15.45$	22.72 23.13	$27.35 \\ 25.78$	$71.0 \\ 72.5$	$90.85 \\ 92.25$	$43.53 \\ 42.02$	$\begin{array}{c} 62.03\\ 63.01 \end{array}$	$53.00 \\ 55.58$	48.03 <b>48.7</b>
KV-Dict s=8	$\mathbf{12.4\%}$	11.66	21.04	22.35	60.0	91.01	40.30	59.60	51.46	44.6
				Mistral	-7B-Instr	uct-v0.3				
Full Cache	100%	41.58	25.69	27.76	76.0	88.59	47.58	59.37	60.60	53.4
KIVI-4 KV-Dict $_{s=24}$	$33.2\%\ 30.6\%$	$\begin{array}{c} 40.37\\ 41.01 \end{array}$	$24.51 \\ 25.32$	$27.75 \\ 27.51$	$74.0 \\ 76.0$		$47.56 \\ 46.27$	$58.49 \\ 59.98$	$58.31 \\ 59.44$	52.4 <b>53.0</b>
KIVI-2 KV-Dict <sub>s=16</sub>	$21.1\% \\ 21.4\%$	$\begin{array}{c} 38.24 \\ 40.34 \end{array}$	$24.08 \\ 24.97$	$26.99 \\ 26.36$	$74.5 \\ 76.0$		$47.66 \\ 45.84$	$57.51 \\ 59.31$	$56.46 \\ 59.50$	51.7 <b>52</b> .7
KV-Dict s=8	$\mathbf{12.4\%}$	33.03	22.80	22.85	68.5	87.85	43.10	56.66	56.85	48.9

Table 3: Experimental results on GSM8K. For KV-Dict, we use N = 4096 as the dictionary size and  $n_b = 128$  as the buffer size. For KIVI, we use g = 64 (group size for quantization) and  $n_b = 64$ (buffer size). Sparsity level s is set to match the average KV size of KIVI, while s = 4 corresponds to cache size unattainable by common 2-bit quantizations. Full cache is in FP16.

(a) Llama 3.x 8B Models

(b) Mistral 7B v0.3 Model

Method	KV Size	Llama-3-8B	3.1-8B-Instruct	Method	KV Size	7B-Instruct
Full Cache	100%	49.89	79.61	Full Cache	100%	48.60
KIVI-4	38.2%	49.13	78.17	KIVI-4	38.2%	48.52
KV-Dict $_{s=24}$	36.9%	48.29	76.88	KV-Dict $_{s=20}$	32.7%	<b>48.60</b>
KIVI-2	25.7%	40.56	67.93	KIVI-2	25.7%	42.91
KV-Dict $_{s=14}$	26.1%	<b>48.75</b>	<b>75.06</b>	KV-Dict $_{s=10}$	22.0%	44.35
KV-Dict s=4	$\mathbf{15.8\%}$	40.03	51.71	KV-Dict $_{s=4}$	$\mathbf{15.8\%}$	39.20

With a KV size of 36.9%, KV-Dict on Llama 8B models experiences a slight accuracy drop of less 360 than 3%p, underperforming KIVI-4 at a similar KV size. However, in the lower memory regime near 361 25% KV size, KV-Dict significantly outperforms KIVI-2, achieving a higher accuracy by 8.2% p on 362 the Llama-3-8B model and 7.1%p on the Llama-3.1-8B-Instruct model. These results highlight the robustness of KV-Dict in low-memory settings, demonstrating that low reconstruction error can be 364 achieved using only a few atoms from our universal dictionary. To further test the resilience of KV-Dict, we set the sparsity to s = 4, observing a noticeable drop in accuracy on the Llama-3.1-8B-Instruct model. Despite this, both Llama models maintain an accuracy above 40%, which is 366 remarkable given that only 4 atoms from KV-Dict were used for each key-value state, utilizing just 367 15.8% of the full cache, including the buffer. 368

369 The performance of KV-Dict on the Mistral-7B-Instruct model is even more impressive. We demon-370 strate that for Mistral, KV-Dict not only outperforms KIVI-4 and KIVI-2 but also achieves higher 371 accuracy with even less memory usage. We also evaluate KV-Dict with s = 4 on the Mistral model and observe an accuracy of 39.2%, further demonstrating robustness in low-memory settings. 372

373

343

344

345

346

347 348

358 359

Results across model sizes and baselines. We illustrate the trade-off between memory usage and 374 375 performance across six different KV cache compression methods on Llama models (1B, 3B, and 8B) in Figure 4. For all three model sizes, KV-Dict consistently lies on the Pareto frontier, achieving 376 higher scores than other compression methods at similar KV cache budget sizes. Notably, KV-Dict 377 demonstrates greater robustness at smaller model scales, with larger performance gaps observed



Figure 4: Memory usage vs. performance of KV-Dict compared to other KV cache compression methods on GSM8K. The figure illustrates the relationship between KV cache size and the performance of KV-Dict on Llama models on GSM8K. For KV-Dict, we use N = 4096 as the dictionary size and  $n_b = 128$  as the buffer size. KV-Dict consistently outperforms both eviction-based methods (SnapKV, PyramidKV) and quantization-based methods (per-token quantization, KIVI, ZipCache).



Figure 5: Memory usage vs. performance of KV-Dict on MMLU-Pro. The figure illustrates the relationship between KV cache size and the performance of KV-Dict on Llama models on MMLU-Pro Engineering/Law. For KV-Dict, we use N = 4096 as the dictionary size and  $n_b = 128$  as the buffer size. KV-Dict often outperforms both eviction-based methods (SnapKV, PyramidKV) and quantization-based methods (per-token quantization, KIVI, ZipCache). For Law, our method slightly underperforms around 25%, but in lower memory regimes, our method still outperforms any other baseline.

425

393

394

395

396

397 398

399

400

401

402

403

404

405 406

407

408

409

410

411

for the 1B and 3B models. In the extremely low-memory regime below 20%, where quantization methods such as KIVI and ZipCache cannot achieve feasible cache sizes, KV-Dict achieves superior performance. Furthermore, while eviction-based methods (SnapKV, PyramidKV) can operate in these extremely low-memory settings, their performance lags significantly behind due to their incompatibility with GQA, making KV-Dict the effective choice for stringent memory constraints.

426 MMLU-Pro Results. Figure 5 illustrates the trade-offs between memory usage and performance
 427 for KV-Dict on the MMLU-Pro Engineering and Law subjects using the Llama-3.1-8B-Instruct
 428 model. KV-Dict outperforms eviction-based methods like SnapKV and PyramidKV across all mem 429 ory settings, though its performance is comparable to quantization-based methods such as KIVI
 430 and ZipCache. In a low memory regime below 20% cache, however, our method still outperforms
 431 any other baseline. This highlights that KV-Dict supports a wide range of compression ratios quite effectively and that our dictionary is generalizable across input distribution.

## 432 3.2 ABLATION STUDY 433

Error thresholding in sparse approximation. KV-Dict also supports a quality-controlled method
 for memory saving by allowing early termination of the sparse approximation process when a pre defined error threshold is met. This approach conserves memory that would otherwise be used for
 marginal improvements in approximation quality. Detailed descriptions and results of this ablation
 study are provided in the Appendix C.1.

**Performance without buffer.** To evaluate the impact of the buffer, we first conducted experiments with varying sparsity without the buffer, with the results shown by the dashed lines in Figure 6 in Appendix C.2. The comparison shows that removing the buffer results in a more pronounced decline in performance, especially at lower KV sizes.

Balancing memory between buffer and sparse representation. Additionally, as shown in Table 4, we examine how balancing memory allocation between the buffer and the sparse representation affects performance. By fixing the total KV cache size at 25% of the original, we vary the
memory distribution between the buffer and the sparse representation across three LongBench tasks:
Qasper, MultiNews, and TREC. The results demonstrate that long-context understanding ability
while using KV-Dict is not solely reliant on the buffer or the sparse representation. Rather, there
exist optimal balance points where performance is maximized for each task.

Table 4: Balancing memory between buffer and sparse representation. This table shows the performance of KV-Dict with the Llama-3.1-8B-Instruct model on LongBench tasks (Qasper, Multi-News, TREC) while varying the memory allocation between the buffer and the sparse representation, with the total KV cache size fixed at 25% of the original size.

Qasper			Mult	iNews	TREC			
s	$n_b$	F1 Score	s	$n_b$	ROUGE-L	s	$n_b$	Accuracy
1	862	6.38	1	503	17.20	1	1232	58.5
4	724	8.36	4	423	20.21	4	1035	63.5
8	517	14.58	8	302	21.27	8	<b>739</b>	65.0
12	<b>278</b>	17.84	12	163	22.81	12	398	63.5
16	0	8.27	16	0	10.70	16	0	54.5

Adaptive dictionary learning. While our universal dictionaries demonstrate strong performance, we explore an adaptive learning method to better incorporate input-specific context. This adaptive approach improves performance by adding new dictionary atoms during generation when the predefined reconstruction error threshold is not met. These atoms, tailored to the input prompt, improve performance but cannot be shared across batches, requiring them to be included in the KV size calculation. While this approach boosts accuracy, it increases memory usage, limiting its ability to achieve low-memory regimes. Detailed methods and results are provided in Appendix C.3.

476 3.3 LATENCY ANALYSIS

In this section, we present latency measurements of the forward pass and OMP portion of KV-Dict during decoding stage in Table 5. We run simple generation tests on a 1000 token input to Llama-3.1-8B-Instruct model and generate up to 250 tokens to measure and aggregate latency metrics. We compare both dictionary sizes N = 1024 and 4096, which primarily affects OMP computation time. We set the sparsity level to s = 24, and process OMP in batches of  $n_a = 8$ .

Although we list the forward pass and OMP separately, the processes are implemented to run in parallel such that the one generation step takes the maximum of the two durations plus some overhead.
However, with parallelization, there exists a time versus space complexity tradeoff, since running OMP also consumes GPU memory.

Table 5: Latency measurements. The following latencies measure the total time it takes for the
respective computation to process when generating one new token. We use Llama-3.1-8B-Instruct
and sum up the time each operation takes in total across all 32 layers. Latencies when dictionary
sizes are 1024 and 4096 are measured. Detailed settings are described in Section 3.3.

Computation Type	Latency (per token)			
-	N = 1024	N = 4096		
Standard forward pass $(\boldsymbol{q}\boldsymbol{K}^{\top})$	48.39 ms	_		
KV-Dict: forward pass using $\boldsymbol{q}(\boldsymbol{K}_{csr}\boldsymbol{D}_k^{\top})^{\top}$	55.56  ms	56.35  ms		
KV-Dict: sparse approximation via OMP	26.57  ms	$40.58 \mathrm{~ms}$		

### 4 RELATED WORK

Prior work on KV cache optimization has explored both training-stage and deployment-focused strategies to improve the efficiency of LLMs. On the deployment side, Kwon et al. (2023) introduce a Paged Attention mechanism and the popular vLLM framework, which adapt CPU-style page memory to map KV caches onto GPU memory using a mapping table, thereby minimizing memory fragmentation and leveraging custom CUDA kernels for efficient inference. While there is a significant and important line of research in this direction (Lin et al., 2024; Qin et al., 2024), this direction is orthogonal to our work and can often be used in tandem with quantization.

507 Current post-training KV cache compression methods can broadly be categorized into eviction, 508 quantization, and merging. Zhang et al. (2024) introduced H2O, which uses attention scores to se-509 lectively retain tokens while preserving recent ones that are strongly correlated with current tokens. 510 Multiple works discuss various heuristics and algorithms to find which tokens can be discarded, 511 while some works find how to complement evictions methods (Ge et al., 2023; Li et al., 2024; Liu 512 et al., 2024a; Devoto et al., 2024; Dong et al., 2024). For this line of work, there is a chance that 513 evictions can work well together with KV-Dict, as Liu et al. (2024a) have successfully combined 514 quantization and eviction.

515 Quantization methods have also played a crucial role in reducing KV cache size without compro-516 mising model performance. Although there is a flurry of work, we only mention those that are most 517 relevant to our discussion and methodology. Hooper et al. (2024) identified outlier channels in key 518 matrices and developed KVQuant, while Liu et al. (2024b) pursue a similar per-channel strategy 519 in KIVI. Further extending these ideas, Yue et al. (2024) presented WKVQuant, which quantizes 520 model weights as well as KV cache using two-dimensional quantization. Kang et al. (2024) follow similar per-channel key and per-token value quantization as KIVI, but with additional low-rank and 521 sparse structures to manage quantization errors. 522

523 524

499 500

#### 5 CONCLUSION

525

526 In conclusion, our proposed method, KV-Dict, offers a novel approach to compressing KV cache 527 for transformers by leveraging low-dimensional structures and sparse dictionary learning. Through 528 this method, we demonstrate that substantial redundancy exists among key states across various in-529 puts, allowing us to compress the KV cache efficiently while maintaining near-lossless performance. 530 Furthermore, KV-Dict enables compression rates that surpass traditional quantization techniques, of-531 fering fine-grained and wide control over memory usage. Importantly, our universal dictionary is both compact and scalable, making it applicable across tasks and user inputs without increasing 532 memory demands. This approach provides strong memory savings, particularly for long-context 533 tasks, by alleviating the memory bottlenecks associated with KV cache storage without dropping 534 any previous tokens. 535

Future research directions based on our work include optimizing CSR tensors through customized
quantizations and improving latency tradeoffs that occur due to the use of OMP during prefilling
and decoding. It would also be interesting to apply "soft-eviction strategies" for sparse tensors in
which sparsity level s is determined or dropped later on based on the estimated importance of the
token. A dynamic allocation of sparsity can further improve our compression method.

#### 540 REFERENCES 541

551 552

553

554

559

560

561

562

576

577

578

579

580 581

582

583

584

585

586

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebron, and Sumit 542 Sanghai. Gqa: Training generalized multi-query transformer models from multi-head check-543 points. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Pro-544 cessing, pp. 4895–4901, 2023. 545
- 546 Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, 547 Xiao Liu, Aohan Zeng, Lei Hou, et al. Longbench: A bilingual, multitask benchmark for long 548 context understanding. arXiv preprint arXiv:2308.14508, 2023. 549
- Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. 550 arXiv preprint arXiv:2004.05150, 2020.
- Trenton Bricken, Adly Templeton, Joshua Batson, Brian Chen, Adam Jermyn, Tom Conerly, Nick Turner, Cem Anil, Carson Denison, Amanda Askell, et al. Towards monosemanticity: Decomposing language models with dictionary learning. Transformer Circuits Thread, 2, 2023. 555
- Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Tianyu Liu, Keming Lu, Wayne Xiong, Yue Dong, 556 Baobao Chang, Junjie Hu, et al. Pyramidkv: Dynamic kv cache compression based on pyramidal information funneling. arXiv preprint arXiv:2406.02069, 2024. 558
  - Emmanuel J Candès, Justin Romberg, and Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. IEEE Transactions on information theory, 52(2):489-509, 2006.
- 563 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve math word problems. arXiv preprint arXiv:2110.14168, 2021. 565
- 566 Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu, Huazuo Gao, Deli Chen, Jiashi Li, Wangding 567 Zeng, Xingkai Yu, Y Wu, et al. Deepseekmoe: Towards ultimate expert specialization in mixture-568 of-experts language models. arXiv preprint arXiv:2401.06066, 2024. 569
- 570 Alessio Devoto, Yu Zhao, Simone Scardapane, and Pasquale Minervini. A simple and effective  $l_2$ norm-based strategy for ky cache compression. arXiv preprint arXiv:2406.11430, 2024. 571
- 572 Harry Dong, Xinyu Yang, Zhenyu Zhang, Zhangyang Wang, Yuejie Chi, and Beidi Chen. Get more 573 with less: Synthesizing recurrence with kv cache compression for efficient llm inference. In 574 Forty-first International Conference on Machine Learning, 2024. 575
  - Weisheng Dong, Guangming Shi, Xin Li, Yi Ma, and Feng Huang. Compressive sensing via nonlocal low-rank regularization. *IEEE transactions on image processing*, 23(8):3618–3632, 2014.
  - David L Donoho. Compressed sensing. IEEE Transactions on information theory, 52(4):1289-1306, 2006.
  - Suyu Ge, Yunan Zhang, Liyuan Liu, Minjia Zhang, Jiawei Han, and Jianfeng Gao. Model tells you what to discard: Adaptive ky cache compression for llms. In The Twelfth International Conference on Learning Representations, 2023.
  - Yefei He, Luoming Zhang, Weijia Wu, Jing Liu, Hong Zhou, and Bohan Zhuang. Zipcache: Accurate and efficient kv cache quantization with salient token identification. arXiv preprint arXiv:2405.14256, 2024.
- 588 Coleman Hooper, Schoon Kim, Hiva Mohammadzadeh, Michael W Mahoney, Yakun Sophia Shao, 589 Kurt Keutzer, and Amir Gholami. Kvquant: Towards 10 million context length llm inference with 590 kv cache quantization. arXiv preprint arXiv:2401.18079, 2024. 591
- Hao Kang, Qingru Zhang, Souvik Kundu, Geonhwa Jeong, Zaoxing Liu, Tushar Krishna, and Tuo 592 Zhao. Gear: An efficient ky cache compression recipefor near-lossless generative inference of llm. arXiv preprint arXiv:2403.05527, 2024.

- <sup>594</sup> Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 611–626, 2023.
- Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye, Tianle
   Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are looking for before
   generation. *arXiv preprint arXiv:2404.14469*, 2024.
- Bin Lin, Tao Peng, Chen Zhang, Minmin Sun, Lanbo Li, Hanyu Zhao, Wencong Xiao, Qi Xu, Xiafei Qiu, Shen Li, et al. Infinite-Ilm: Efficient Ilm service for long context with distattention and distributed kvcache. *arXiv preprint arXiv:2401.02669*, 2024.
- Zichang Liu, Aditya Desai, Fangshuo Liao, Weitao Wang, Victor Xie, Zhaozhuo Xu, Anastasios Kyrillidis, and Anshumali Shrivastava. Scissorhands: Exploiting the persistence of importance hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing Systems*, 36, 2024a.
- Zirui Liu, Jiayi Yuan, Hongye Jin, Shaochen Zhong, Zhaozhuo Xu, Vladimir Braverman, Beidi
  Chen, and Xia Hu. Kivi: A tuning-free asymmetric 2bit quantization for kv cache. *arXiv preprint arXiv:2402.02750*, 2024b.
- Ariel Lubonja, Sebastian Kazmarek Præsius, and Trac Duy Tran. Efficient batched cpu/gpu implementation of orthogonal matching pursuit for python. *arXiv preprint arXiv:2407.06434*, 2024.
- 618 Alireza Makhzani and Brendan Frey. K-sparse autoencoders. arXiv preprint arXiv:1312.5663, 2013.
- Stephen Merity. The wikitext long term dependency language modeling dataset. Salesforce Metamind, 9, 2016.
- 622 Christopher A Metzler, Arian Maleki, and Richard G Baraniuk. From denoising to compressed
   623 sensing. *IEEE Transactions on Information Theory*, 62(9):5117–5144, 2016.
- Ruoyu Qin, Zheming Li, Weiran He, Mingxing Zhang, Yongwei Wu, Weimin Zheng, and Xinran Xu. Mooncake: A kvcache-centric disaggregated architecture for llm serving, 2024. *arXiv* preprint arxiv:2407.00079, 2024.
- Noam Shazeer. Fast transformer decoding: One write-head is all you need. *arXiv preprint arXiv:1911.02150*, 2019.
- Prajwal Singhania, Siddharth Singh, Shwai He, Soheil Feizi, and Abhinav Bhatele. Loki: Low-rank
   keys for efficient sparse attention. *arXiv preprint arXiv:2406.02542*, 2024.
- Yutao Sun, Li Dong, Yi Zhu, Shaohan Huang, Wenhui Wang, Shuming Ma, Quanlu Zhang, Jianyong Wang, and Furu Wei. You only cache once: Decoder-decoder architectures for language models. *arXiv preprint arXiv:2405.05254*, 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, pp. 5998–6008, 2017.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming
  Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging
  multi-task language understanding benchmark. *arXiv preprint arXiv:2406.01574*, 2024a.
- <sup>643</sup>
   <sup>644</sup>
   <sup>645</sup>
   <sup>646</sup> Zheng Wang, Boxiao Jin, Zhongzhi Yu, and Minjia Zhang. Model tells you where to merge: Adaptive kv cache merging for llms on long-context tasks. *arXiv preprint arXiv:2407.08454*, 2024b.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
   language models with attention sinks. In *The Twelfth International Conference on Learning Representations*, 2023.

648 649	Lu Ye, Ze Tao, Yong Huang, and Yang Li. Chunkattention: Efficient self-attention with prefix-aware ky cache and two-phase partition. <i>arXiv preprint arXiv:2402.15220</i> , 2024.
650	
651	Gyeong-In Yu, Joo Seong Jeong, Geon-Woo Kim, Soojeong Kim, and Byung-Gon Chun. Orca: A
652	distributed serving system for {Transformer-Based} generative models. In 16th USENIX Sympo-
653	sium on Operating Systems Design and Implementation (OSDI 22), pp. 521–538, 2022.
654	Hao Yu, Zelan Yang, Shen Li, Yong Li, and Jianxin Wu. Effectively compress kv heads for llm.
655 656	arXiv preprint arXiv:2406.07056, 2024.
657	Yuxuan Yue, Zhihang Yuan, Haojie Duanmu, Sifan Zhou, Jianlong Wu, and Liqiang Nie. Wkvquant:
658	Quantizing weight and key/value cache for large language models gains more. arXiv preprint
659	arXiv:2402.12065, 2024.
660	Zhanyu Zhang, Ving Shang, Tianyi Zhau, Tianlang Chan, Lianmin Zhang, Duisi Cai, Zhao Sang
661	Zhenyu Zhang, Ting Sheng, Hanyi Zhou, Hanong Chen, Liamini Zheng, Kuisi Cai, Zhao Song, Yuandong Tian, Christopher Pá, Clark Parratt, et al. 120: Heavy hitter oragle for efficient gen
662	erative inference of large language models. Advances in Neural Information Processing Systems
663	36 2024
664	50, 2021.
665	Hufei Zhu, Wen Chen, and Yanpeng Wu. Efficient implementations for orthogonal matching pursuit.
666	<i>Electronics</i> , 9(9):1507, 2020.
667	
669	
000	
669	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

## 702 APPENDIX

704 705

706

707

708

709

710

711 712

713

714

715

716

717

718

719

720

721

722 723 724

725

#### A IMPLEMENTATION DETAILS

Algorithm 1 illustrates a naive implementation of OMP for understanding. In KV-Dict, we adopt the implementation of OMP v0 proposed by (Zhu et al., 2020), which minimizes computational complexity using efficient inverse Cholesky factorization. Additionally, we integrate methods from (Lubonja et al., 2024) for batched GPU execution and extend the implementation to handle multiple dictionaries in parallel.

#### Algorithm 1 OMP

**Require:** Signal  $\mathbf{k} \in \mathbb{R}^m$ , dictionary  $\mathbf{D} \in \mathbb{R}^{m \times N}$ , sparsity s **Ensure:** Sparse representation  $\mathbf{y} \in \mathbb{R}^n$ 1: Initialize  $\mathbf{r}^{(0)} \leftarrow \mathbf{k}$ ,  $\mathbb{I}^{(0)} \leftarrow \emptyset$ ,  $\mathbf{y}^{(0)} \leftarrow \mathbf{0}$ 2: for i = 1 to s do 3:  $n^{(i)} \leftarrow \arg \max_{1 \le n \le N} \left\{ \left| \left( \mathbf{D}^\top \left( \mathbf{k} - \mathbf{D} \mathbf{y}^{(i)} \right) \right)_n \right| \right\}$ 4:  $\mathbb{I}^{(i)} \leftarrow \mathbb{I}^{(i-1)} \cup \left\{ n^{(i)} \right\}$ 5:  $\mathbf{y}^{(i+1)} \leftarrow \arg \min_{\mathbf{y} \in \mathbb{R}^N} \left\{ \left\| \mathbf{k} - \mathbf{D} \mathbf{y} \right\|_2$ , Supp  $(\mathbf{y}) \subset \mathbb{I}^{(i)} \right\}$ 6: end for 7: return  $\mathbf{y}$ 

#### Algorithm 2 Prefilling and decoding with KV-Dict

726 1: **Parameter:** sparsity s, buffer length  $n_b$ , approximation length  $n_a$ 727 2: procedure PREFILLING 728 Input:  $oldsymbol{X} \in \mathbb{R}^{l_{ ext{seq}} imes d}$ 3: 729 4:  $K \leftarrow XW_k, V \leftarrow XW_v$ 730  $\boldsymbol{K}_{\text{csr}} \leftarrow \text{OMP}\left(\boldsymbol{K}\left[:l_{ ext{seq}} - n_b
ight], \boldsymbol{D}_k, s
ight)$ 5: 731  $V_{\text{csr}} \leftarrow \text{OMP}\left(\boldsymbol{V}\left[:l_{\text{seq}} - n_b\right], \boldsymbol{D}_v, s\right)$ 6: 732  $\boldsymbol{K}_{\text{buffer}} \leftarrow \boldsymbol{K}[l_{\text{seq}} - n_b :], \boldsymbol{V}_{\text{buffer}} \leftarrow \boldsymbol{V}[l_{\text{seq}} - n_b :]$ 7: 733 8: KV cache  $\leftarrow K_{csr}, K_{buffer}, V_{csr}, V_{buffer}$ 734 9: return K, V 735 10: end procedure 736 11: procedure DECODING Input: KV cache,  $oldsymbol{x}_t \in \mathbb{R}^{1 imes d}$ 738 12: 739  $oldsymbol{q}_t \leftarrow oldsymbol{x}_t W_q, oldsymbol{k}_t \leftarrow oldsymbol{x}_t W_k, oldsymbol{v}_t \leftarrow oldsymbol{x}_t W_v$ 13:  $K_{\text{csr}}, K_{\text{buffer}}, V_{\text{csr}}, V_{\text{buffer}} \leftarrow \text{KV}$  cache 740 14: 15:  $\boldsymbol{K}_{\text{buffer}} \leftarrow \text{Concat}\left(\left[\boldsymbol{K}_{\text{buffer}}, \boldsymbol{k}_{t}\right], \text{dim} = \text{token}\right)$ 741 16:  $V_{\text{buffer}} \leftarrow \text{Concat}\left(\left[V_{\text{buffer}}, v_t\right], \text{dim} = \text{token}\right)$ 742  $\boldsymbol{a}_t \leftarrow \operatorname{Concat}\left(\left[\boldsymbol{q}_t \boldsymbol{D}_k \boldsymbol{K}_{\operatorname{csr}}, \boldsymbol{q}_t \boldsymbol{K}_{\operatorname{buffer}}\right], \operatorname{dim} = \operatorname{token}\right)$ 17: 743 18:  $\boldsymbol{a}_{t} \leftarrow \text{Softmax}\left(\boldsymbol{a}_{t}\right)$ 744 19:  $V \leftarrow \text{Concat}([D_v V_{csr}, V_{buffer}], \text{dim} = \text{token})$ 745 20:  $o_t \leftarrow a_t V$ 746 21: if len  $(\mathbf{K}_{buffer}) > n_b$  then 747 22:  $\boldsymbol{K}'_{csr} \leftarrow OMP\left(\boldsymbol{K}_{buffer}\left[:n_a\right], \boldsymbol{D}_k, s\right)$ 748  $\boldsymbol{V}'_{\text{csr}} \leftarrow \text{OMP}\left(\boldsymbol{V}_{\text{buffer}}\left[:n_a\right], \boldsymbol{D}_v, s\right)$ 23: 749  $\boldsymbol{K}_{csr} \leftarrow Concat\left(\left[\boldsymbol{K}_{csr}, \boldsymbol{K'}_{csr}\right], dim = token\right)$ 24:  $\begin{array}{l} \mathbf{V}_{\mathrm{csr}} \leftarrow \mathrm{Concat}\left([\mathbf{V}_{\mathrm{csr}}, \mathbf{V}'_{\mathrm{csr}}], \mathrm{dim} = \mathrm{token}\right) \\ \mathbf{K}_{\mathrm{buffer}} \leftarrow \mathbf{K}_{\mathrm{buffer}}\left[n_a:], \mathbf{V}_{\mathrm{buffer}} \leftarrow \mathbf{V}_{\mathrm{buffer}}\left[n_a:] \end{array}$ 750 25: 26: 751 end if 27: 752 28: KV cache  $\leftarrow K_{csr}, K_{buffer}, V_{csr}, V_{buffer}$ 753 29: return o<sub>t</sub> 754 30: end procedure 755

### **B** LONGBENCH TASK STATISTICS

#### Table 6: Details of LongBench tasks used in experiments.

Task	Task Type	<b>Evaluation Metric</b>	Average Length	# of Samples
Qasper	Single-doc QA	F1	3619	200
QMSum	Summarization	ROUGE-L	10614	200
MultiNews	Summarization	ROUGE-L	2113	200
TREC	Few-shot information retrieval	Accuracy	5177	200
TriviaQA	Few-shot reading comprehension	F1	8209	200
SAMSum	Few-shot dialogue summarization	ROUGE-L	6258	200
LCC	Code completion	Edit Similarity	1235	500
RepoBench-P	Code completion	Edit Similarity	4206	500

769 770

771

772 773

774

787 788 789

796

798 799

756

758

759

#### C ABLATION STUDY: EXPERIMENTAL RESULTS

#### C.1 ERROR THRESHOLDING IN SPARSE APPROXIMATION

For the error thresholding ablation study, detailed results are provided in Table 7. We set a maximum sparsity of 32, corresponding to the maximum number of iterations for the OMP algorithm. However, if the reconstruction error at any iteration falls below a predefined error threshold, we let the OMP terminate early, saving memory that would otherwise be used for minor approximation improvements. This approach is particularly compatible with OMP, as its greedy nature ensures that early termination yields the same results as using higher sparsity (less non-zero elements). Additionally, OMP inherently computes the residual at each iteration, allowing for continuous evaluation of the relative reconstruction error without requiring any additional computation.

Table 7: Impact of error thresholding on LongBench performance and memory usage. The table presents the performance of KV-Dict on the Llama-3.1-8B-Instruct model at various reconstruction error thresholds ( $\delta$ ) for early termination of the sparse approximation algorithm. A dictionary size of N = 1024 and FP16 precision for the values of the CSR tensors are used.

Threshold $(\delta)$	KV Size	Qasper	QMSum	MultiNews	TREC	TriviaQA	SAMSum	LCC	RepoBench-P	Average
				Llama	-3.1-8B-I	nstruct				
Full Cache	100%	22.54	24.57	27.44	72.5	91.65	43.47	63.15	56.76	50.26
0.2	50.6%	20.03	23.65	26.44	72.5	91.61	43.47	62.72	56.63	49.63
0.3	41.1%	16.49	23.35	25.34	72.5	91.34	43.02	62.53	56.65	48.90
0.4	30.9%	16.08	22.91	23.77	69.5	90.79	42.70	61.28	54.82	47.73
0.5	22.8%	12.43	21.75	21.29	57.5	88.56	41.04	58.85	53.19	44.33

#### C.2 PERFORMANCE WITHOUT BUFFER

In this section, we assess the effect of the buffer by comparing results with and without its use. The results for LongBench and GSM8K are presented in Table 8 and Table 9, respectively.

#### 800 C.3 Adaptive KV-Dict

801 While we observe some degree of universality in our dictionaries, as shown in Table 1, their perfor-802 mance is particularly strong on WikiText-103, the dataset they were trained on. To better incorporate 803 input context information, we propose an extension that adaptively learns the dictionary during gen-804 eration. In this framework, we begin with a pre-trained universal dictionary as the initial dictionary. 805 If, during the generation process, the sparse approximation fails to meet the predefined relative re-806 construction error threshold, the problematic uncompressed key or value vector is normalized and 807 added to the dictionary. The sparse representation of this vector is then stored with a sparsity of s = 1, where its index corresponds to the newly added atom and its value is the  $\ell_2$  norm of the un-808 compressed vector. The updated dictionary is subsequently used for further sparse approximations 809 during the generation task. In this way, the adaptive learning framework incrementally refines the



Figure 6: Memory usage vs. performance of KV-Dict with and without buffer on LongBench and GSM8K. The figure illustrates the impact of removing the buffer on the performance of KV-Dict when evaluated on the Llama 3.1-8B-Instruct and Mistral-7B-Instruct models for LongBench (left) and GSM8K (right) tasks. Solid lines represent configurations with a buffer, while dashed lines represent configurations without a buffer. We use a dictionary size of N = 1024 and FP16 precision for the values of CSR tensors to vary sparsity and explore a wide range of KV sizes.

Table 8: LongBench performance without buffer. This table shows the impact of removing the buffer of KV-Dict on the performance of the Llama-3.1-8B-Instruct and Mistral-7B-Instruct models at varying sparsity levels. A dictionary size of N = 1024 and FP16 precision for the values of the CSR tensors are used.

Sparsity	KV Size	Qasper	QMSum	MultiNews	TREC	TriviaQA	SAMSum	LCC	RepoBench-P	Average
-				Llam	a-3.1-8B-	Instruct				
Full Cache	100%	13.10	23.46	26.94	72.5	91.65	43.47	63.15	56.76	48.88
s = 32	50.8%	14.87	26.51	26.57	71.5	92.48	42.88	61.54	54.04	48.80
s = 24	38.2%	13.37	25.02	22.54	65.0	91.75	39.71	52.21	46.48	44.51
s = 16	25.8%	8.27	13.74	10.70	54.5	77.51	20.45	26.53	22.46	29.27
s = 12	19.5%	6.31	10.15	5.66	39.0	53.70	6.83	22.18	19.46	20.41
s = 8	13.3%	2.74	8.05	4.17	36.5	34.45	4.27	18.24	18.32	15.84
				Mistra	al-7B-Inst	truct-v0.3				
Full Cache	100%	41.58	25.69	27.76	76.0	88.59	47.58	59.37	60.60	53.40
s = 32	50.8%	40.27	25.21	27.53	76.5	89.01	45.77	58.64	59.07	52.75
s = 24	38.2%	37.46	24.41	27.34	75.5	88.66	43.87	48.55	49.50	49.41
s = 16	25.8%	25.57	18.49	15.19	71.5	81.91	27.90	19.39	21.45	35.18
s = 12	19.5%	18.59	13.11	5.95	58.0	50.13	2.86	13.38	12.60	21.83
s = 8	13.3%	10.32	6.98	2.67	31.5	20.01	2.27	10.18	8.11	11.51

Table 9: **GSM8K performance without buffer.** This table shows the impact of removing the buffer of KV-Dict on the performance of the Llama-3.1-8B-Instruct and Mistral-7B-Instruct models at varying sparsity levels. A dictionary size of N = 1024 and FP16 precision for the values of the CSR tensors are used.

Sparsity	KV Size	Llama-3.1-8B-Instruct	Mistral-7B-Instruct-v0.3
Full Cache	100%	79.61	48.60
s = 32	50.8%	69.07	43.97
s = 24	38.2%	32.75	23.20
s = 16	25.8%	1.97	1.29
s = 12	19.6%	1.36	0.76

dictionaries, tailoring them to the specific generative task and enhancing overall performance at the cost of additional memory usage.

In our experiment, we use a universal dictionary of size 1024, allowing up to 1024 additional atoms to be added during generation. The maximum sparsity of 16 is used, with a buffer size of  $n_b = 128$ , and FP16 precision for the values of the CSR tensors. Results of this experiment are presented in

Table 10: **GSM8K performance of adaptive KV-Dict.** The table shows the GSM8K performance and KV cache sizes of adaptive KV-Dict on the Llama-3.1-8B-Instruct and Mistral-7B-Instruct-v0.3 models at varying reconstruction error thresholds ( $\delta$ ). A universal dictionary of size 1024 is used, with up to 1024 additional atoms added during generation. The maximum sparsity of s = 16, buffer size of  $n_b = 128$ , and FP16 precision for CSR tensor values are applied.

<b>Threshold</b> ( $\delta$ )	Llama-3.1-8B-Instruct		Mistral-7B-Instruct-v0.3		
	KV Size	GSM8K Score	KV Size	GSM8K Score	
Full Cache	100%	79.61	100%	48.60	
0.25	N/A	N/A	42.1%	48.07	
0.30	43.5%	77.41	41.3%	48.14	
0.35	42.0%	76.80	39.8%	47.76	

#### D QUALITATIVE RESULTS

#### D.1 GSM8K EXAMPLE GENERATION

The boxes highlight the correct intermeidate answers that must be derived. The highlighted colors in red show incorrect reasoning and answers. We provide the original problem and solution below, along with outputs from Llama-3.1-8B-Instruct with full KV cache (FP16), KV-Dict (sparsity s = 14), and KIVI-2.

#### **Original Problem & Solution**

**Problem:** The Kennel house keeps 3 German Shepherds and 2 Bulldogs. If a German Shepherd consumes 5 kilograms of dog food and a bulldog consumes 3 kilograms of dog food per day, how many kilograms of dog food will they need in a week?

**Solution:** The German Shepherds consume  $3 \times 5 = 15$  kilograms of dog food per day.

The Bulldogs consume  $2 \times 3 = 6$  kilograms of dog food per day.

The dogs consume a total of 15 + 6 = 21 kilograms of dog food per day.

Therefore, they need  $21 \times 7 = 147$  kilograms of dog food in a week.

Final Answer: 147

#### Llama-3.1-8B-Instruct with Full Cache

First, find the total amount of food the German Shepherds consume in a day:  $3 \text{ dogs} \times 5 \text{ kg/dog} = 15 \text{ kg}$ .

Then find the total amount of food the Bulldogs consume in a day:  $2 \text{ dogs} \times 3 \text{ kg/dog} = 6 \text{ kg}$ . Add the two amounts together to find the total amount of food the dogs consume in a day: 15 kg + 6 kg = 21 kg.

Multiply that number by 7 to find the total amount of food the dogs consume in a week: 21 kg/day  $\times$ 7 days = 147 kg.

913 914 915

882 883

885

887 888

889

890

891 892 893

894

895

896

897

899

900

901

902 903

904 905 906

907

908

909

910

911

912

916 917

Llama-3.1-8B-Instruct with KV-Dict (s = 14)

The kennel has 3 German Shepherds and 2 Bulldogs, so they have $3 + 2 = 5$ dogs.	
The German Shepherds consume $5 \times 3 = 15$ kilograms of dog food per day.	
The Bulldogs consume $2 \times 3 = 6$ kilograms of dog food per day.	
So, the kennel needs $15 + 6 = 21$ kilograms of dog food per day.	
In a week, they will need $21 \times 7 = 147$ kilograms of dog food.	
The answer is 147.	

#### Llama-3.1-8B-Instruct with KIVI-2

The Kennel house has 3 German Shepherds and 2 Bulldogs, so they have a total of 3 + 2 = 5 dogs.

Each German Shepherd consumes 5 kilograms of dog food per day, so the total amount of dog food consumed by the German Shepherds is  $5 \times 5 = 25$  kilograms per day.

Each Bulldog consumes 3 kilograms of dog food per day, so the total amount of dog food consumed by the Bulldogs is  $2 \times 3 = 6$  kilograms per day.

The total amount of dog food consumed per day is 25 + 6 = 31 kilograms.

The Kennel house will need  $31 \times 7 = 217$  kilograms of dog food in a week.