

---

# Spotlight Attention: Towards Efficient LLM Generation via Non-linear Hashing-based KV Cache Retrieval

---

**Anonymous Author(s)**

Affiliation

Address

email

## Abstract

1 Reducing the key-value (KV) cache burden in Large Language Models (LLMs)  
2 significantly accelerates inference. Dynamically selecting critical KV caches  
3 during decoding helps maintain performance. Existing methods use random linear  
4 hashing to identify important tokens, but this approach is inefficient due to the  
5 orthogonal distribution of queries and keys within two narrow cones in LLMs. We  
6 introduce Spotlight Attention, a novel method that employs non-linear hashing  
7 functions to optimize the embedding distribution of queries and keys, enhancing  
8 coding efficiency and robustness. We also developed a lightweight, stable training  
9 framework using a Bradley-Terry ranking-based loss, enabling optimization of the  
10 non-linear hashing module on GPUs with 16GB memory in 8 hours. Experimental  
11 results show that Spotlight Attention drastically improves retrieval precision while  
12 shortening the length of the hash code at least 5 $\times$  compared to traditional linear  
13 hashing. Finally, we exploit the computational advantages of bitwise operations  
14 by implementing specialized CUDA kernels, achieving hashing retrieval for 512K  
15 tokens in under 100 $\mu$ s on a single A100 GPU, with end-to-end throughput up to  
16 3 $\times$  higher than vanilla decoding. All the training and evaluation stuff can be found  
17 at Anonymous/Spotlight.

## 18 1 Introduction

19 Large Language Models (LLMs) are propelling groundbreaking advancements in various natural  
20 language tasks, significantly enhancing applications such as content creation and chat assistance.  
21 Generally, the inference process of LLMs can be divided into (1) the pre-filling phase calculates the  
22 key-value (KV) cache for input tokens in the prompt prior to autoregressive generation, and (2) the  
23 decoding phase auto-regressively generates tokens, producing one token per forward pass based on  
24 the KV cache. Among them, the decoding phase serves as the primary inference bottleneck due to  
25 the frequent exchanges between on-board and on-chip memory for model parameters and KV cache,  
26 which limits GPU scalability [3] more so than the pre-filling phase that processes input prompts in  
27 parallel. For example, deploying LLaMA2-7B [21] on an A100 GPU for a single request achieves  
28 nearly 100% GPU utilization during the pre-filling phase but drops to below 10% on average during  
29 decoding, which largely restrains the inference efficiency of LLMs.

30 To alleviate this inference bottleneck, extensive research has focused on heuristically eliminating  
31 the KV cache burden based on attention scores [25, 23, 17]. While effective for short sequences,  
32 such irreversible removal of KV cache can significantly degrade performance on long-sequence  
33 tasks, especially in Needle-in-a-Haystack scenarios [14]. To explain, tokens initially considered  
34 unimportant and removed might later attain higher attention scores during the prolonged decoding  
35 phase, which is crucial for output quality [20, 15]. To overcome this limitation, recent works have

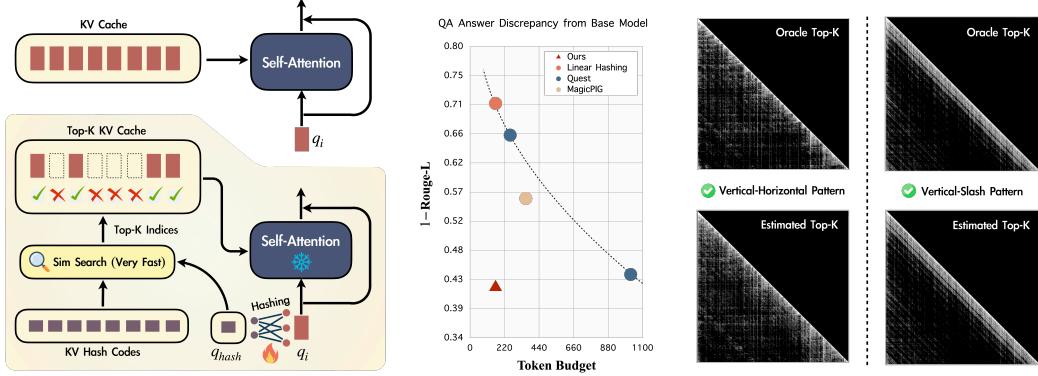


Figure 1: **Overview.** (Left) **Architecture.** Comparison of our Spotlight Attention versus normal attention. Spotlight Attention adds an additional hash code-based retrieval mechanism for each layer. (Middle) **Performance.** Spotlight attention achieves the most accurate retrieval and generates the closest response compared to the original model on QA datasets. (Right) **Visualization.** For arbitrarily complex attention patterns, our method estimates the top-k sequences well, with an average correctness rate of more than half for different models.

36 turned to retaining all KV cache tokens and dynamically selecting important tokens for computation  
37 during decoding [20, 10], which is the focus of this paper.

38 Despite the convincing performance of such on-the-fly KV cache selection, how to effectively pick  
39 up those important tokens remains challenging. As a pioneering effort, Quest [20] selects tokens by  
40 matching queries and keys in a block-wise manner. While effective, such coarse-grained selection  
41 hardly guarantees precise localization of important tokens. MagicPIG [10] advances Quest by  
42 implementing token-level cache retrieval, specifically utilizing Local Sensitive Hashing (LSH) to  
43 encode queries and keys into hash codes and pinpointing the best matches as the selected tokens.  
44 However, as depicted in Figure 2a, the efficacy of such linear hashing heavily depends on the hash  
45 code length, *e.g.*, a hash code length of 1,024 bits per query/key is necessitated to achieve promising  
46 token retrieval. Considering the already substantial size of the KV cache, storing these lengthy hash  
47 codes markedly impairs deployment efficiency. Moreover, employing a linear projection with a large  
48 output dimension for key processing incurs significant computational overhead.

49 Delving deeper, prior work [10] has discovered that the queries and keys within LLMs typically form  
50 nearly orthogonal cones within the embedding space, as depicted in Figure 2b. Given the truth that  
51 linear hashing function partitions the embedding space using random hyperplanes, such orthogonal  
52 distribution of queries and keys barely lead to satisfying encoding quality, which can even result  
53 in a collapse of the hashing outcomes, *i.e.*, identical codes for all queries and keys, as shown in  
54 Figure 2c. Therefore, extremely long hash codes are necessary to mine meaningful information and  
55 accurately match essential tokens. MagicPIG attempted to mitigate this issue by normalizing the keys  
56 before retrieval. However, this approach remains suboptimal as it overlooks the query distribution  
57 and introduces bias to the retrieval process.

58 To address the aforementioned limitation, we propose Spotlight Attention, a novel method that  
59 replaces random hyperplanes with curved surfaces for space partitioning via a non-linear MLP  
60 hashing function. As depicted in Figure 2d, this non-linearity can better fit skewed distributions,  
61 thereby improving code quality. In particular, we utilize the Bradley-Terry ranking objective [7] to  
62 optimize the non-linear MLP layer, wherein the learning target involves minimizing the difference  
63 between the estimated top- $k$  indices and the ground truth top- $k$  indices obtained via the vanilla  
64 attention scores. This learning process is exceptionally efficient, with the LLM backbone remaining  
65 frozen and requiring only a minimal amount of calibration data. As a result, the optimized non-linear  
66 hashing function can match the performance of linear hashing while using 5× shorter hash codes,  
67 achieving higher efficiency than MagicPIG. We further implemented CUDA kernels for the hash code  
68 processing, including bit-packing and bitwise NXOR GEMM operators, achieving significant latency  
69 reductions in practice. For example, our method achieves up to 3× increase in Qwen2.5-7B [24]  
70 inference throughput for both 32K and 128K sequences, with only ~2% performance degradation on  
71 the LLaMA3 [13] series and no loss on Qwen2.5 [24] series.

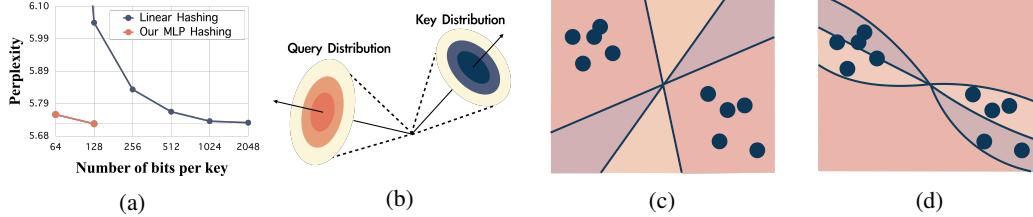


Figure 2: **Motivation.** (a) The empirical evaluation shows that upgrading the hashing function from linear to MLP can bring a huge improvement, (b) this is because query and key are usually distributed in two small cones in the space [10]. (c) In this situation, it is difficult for the space to be uniformly partitioned by linear boundaries, (d) which can be well solved by using an MLP hashing function.

72 Our contributions are threefold:

- 73     • We propose Spotlight Attention for accelerating LLM inference, which employs non-linear  
74       hashing function to encode and match queries and key values within LLMs, thereby effi-  
75       ciently selecting critical KV cache for model inference.  
76     • We develop a lightweight and robust training framework based on the Bradley-Terry ranking  
77       objective, which effectively optimizes the non-linear hashing function using only a small  
78       amount of calibration data.  
79     • Extensive experiments demonstrate that Spotlight Attention can drastically reduce LLM  
80       inference latency while maintaining the strongest performance retention in comparison with  
81       state-of-the-art methods.

## 82 2 Related Work

83 This section covers the spectrum of studies on LLM KV cache pruning that are closely related to  
84 our work, which we heuristically categorize into static pruning, dynamic pruning with permanent  
85 eviction, and dynamic pruning without permanent eviction.

86 **Static KV cache pruning.** These methods compress the KV cache once after the pre-filling phase,  
87 using the compressed cache for subsequent decoding. For example, FastGen [12] introduces a pattern-  
88 aware approach by identifying five fundamental attention structures and applying targeted selection  
89 strategies. SnapKV [16] further simplifies FastGen by focusing solely on retrieving tokens based on  
90 their importance scores, showing that only a subset of prompt tokens carry critical information for  
91 response generation and retain their significance during the whole decoding phase. However, without  
92 pruning during decoding, these methods are primarily suited for scenarios with long prompts and  
93 relatively short responses.

94 **Dynamic pruning with permanent eviction.** This category of methods performs dynamic KV  
95 cache pruning during the decoding phase, permanently removing pruned KV cache tokens from  
96 memory. For example, H2O [25] leverages cumulative attention scores to retain high-impact tokens.  
97 NACL [9] identifies a fundamental limitation in H2O, namely their dependence on potentially  
98 biased local attention statistics. To overcome this issue, they develop an alternative approach,  
99 implementing a diversified random eviction strategy. Keyformer [2] highlights that token removal  
100 distorts the underlying softmax probability distribution. Considering the pivotal role of softmax  
101 distributions in token significance evaluation, they incorporate regularization techniques to mitigate  
102 these distributional perturbations. Unlike static KV cache selection, these methods enable dynamic  
103 pruning during decoding, making them better suited for tasks requiring extensive generation. However,  
104 they assume that critical information is concentrated in a small subset of KV cache tokens, a condition  
105 that does not always hold. As MagicPIG [10] points out, token importance can vary significantly  
106 across tasks, leading to premature eviction of tokens before they are needed. For example, H2O may  
107 fail to answer questions like *a is b, c is d, a is ?* due to forgetting earlier facts.

108 **Dynamic pruning without permanent eviction.** The limited applicability of permanent token  
109 eviction methods has led to a shift toward non-permanent eviction approaches. These methods  
110 assume the importance of KV cache tokens varies with each query, requiring importance estimation  
111 at every decoding step. Instead of permanently evicting unimportant tokens, they exclude them

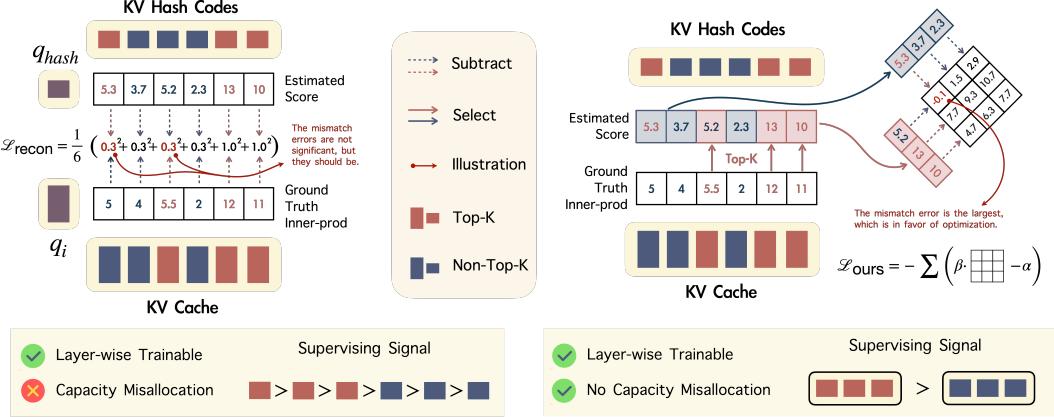


Figure 3: **Optimization.** (Left) *Reconstruction Loss*. This loss minimizes the MSE between the estimated and ground-truth attention scores. It has two main drawbacks. First, it is highly sensitive to score magnitudes and prone to outliers. Second, it wastes most of the hashing function’s capacity on preserving order within the top- $k$  and non-top- $k$  sets. (Right) *Our Ranking Loss*. Our loss adopts the Bradley–Terry ranking objective, which is robust to score magnitude and outliers, and provides supervision focused solely on distinguishing between top- $k$  and non-top- $k$  sets.

from attention calculations for that step only. While this improves accuracy, it demands frequent importance estimation, unlike permanent eviction methods that prune tokens in batches after many steps. Research has therefore focused on optimizing the efficiency and accuracy of these estimations. Quest [20] groups KV cache tokens into blocks, estimating block importance via the dot product between queries and block representations derived from the minimum and maximum key values. Although efficient, this approach suffers from internal fragmentation, as entire blocks are processed even if only a few tokens are important. MagicPIG [10] eliminates this issue by mapping queries and keys to hash codes for token-level retrieval via Hamming distance, avoiding fragmentation but reducing efficiency. Building on MagicPIG, our method significantly shortens hash code, drastically reducing computation while preserving accuracy.

### 3 Methodology

#### 3.1 Preliminary

**Attention computing.** We first present the basic preliminaries for attention computation and KV cache pruning during the decoding phase of LLMs. We define the query, key, and value inputs to the attention module as  $Q, K, V \in \mathbb{R}^{1 \times d}$ , where  $d$  is the embedding dimension. We use  $\oplus$  to denote concatenation, and  $K_{\text{cache}}, V_{\text{cache}} \in \mathbb{R}^{n \times d}$  represent the key-value cache generated during the pre-filling phase and previous decoding steps. With these definitions, the standard attention is calculated as follows:

$$\mathcal{A} = \text{softmax} \left( \frac{f(Q, K_{\text{cache}}) \oplus QK^{\top}}{\sqrt{d}} \right) (V_{\text{cache}} \oplus V), \quad (1)$$

where  $f(X, X') = XX'^{\top}$  is inner-product.

**KV cache pruning.** As the decoding sequence length increases, the size of the KV cache  $\{K, V\}_{\text{cache}}$  can grow exceedingly large, creating an LLM inference bottleneck. KV cache pruning that selectively preserves only essential portions of the cache for computation serves as an efficient way to alleviate this problem. Given a desired cache budget  $K$ , it first identifies the indices  $\mathcal{I}$  of top- $K$  important tokens and then extracts a subset of the KV cache  $\{K, V\}_{\text{subset}}$  for attention computation as

$$\{K, V\}_{\text{subset}} = \text{gather}(\{K, V\}_{\text{cache}}, \mathcal{I}). \quad (2)$$

As previously discussed, existing methods for pruning the KV cache either permanently eliminate cache entries not in the set  $\mathcal{I}$  [25, 16] or retain all caches but dynamically determine  $\mathcal{I}$  during the decoding process [20, 10]. In this paper, we focus on the latter due to its superior performance preservation.

Table 1: **KV retrieval accuracy.** Measured by IoU $\uparrow$  / PPL $\downarrow$  changes before and after training. (1) Training is essential for efficient MLP hashing. (2) Limited improvement from training linear hashing highlights the necessity of MLP hashing. See Appendix B.1 for per-head IoU scores.

Method	Original	Oracle Top-2%	LSH Top-2%		MLP Hashing Top-2% (ours)	
			Before	After	Before	After
LLaMA2-7B	- / 5.58	1.00 / 5.69	0.17 / 5.86	0.20 / 5.84	0.05 / 20.31	0.41 / 5.72
LLaMA2-7B-Chat	- / 7.10	1.00 / 6.87	0.17 / 7.34	0.19 / 7.45	0.05 / 21.34	0.42 / 6.98
LLaMA3-8B	- / 6.45	1.00 / 6.63	0.15 / 7.12	0.18 / 7.07	0.07 / 148.2	0.34 / 6.69
Qwen2.5-7B	- / 7.17	1.00 / 7.28	0.13 / 8.81	0.16 / 8.73	0.09 / 22.07	0.35 / 7.31

### 140 3.2 Revisiting Token-level Cache Retrieval

141 Token-level cache retrieval refers to dynamically selecting cached entries at the granularity of  
 142 individual tokens [10].

143 **Oracle top-k retrieval.** We conducted a preliminary experiment to assess the upper-bound per-  
 144 formance, using full-precision attention scores  $f(Q, K_{\text{cache}})$  to select key tokens.<sup>1</sup> Surprisingly, by  
 145 pruning layers beyond the first two, we discarded up to 98% of the KV cache with only a 0.1 per-  
 146 plexity increase on PG19, revealing significant untapped potential. This suggests that top- $k$  retrieval  
 147 enables near-lossless compression, contrasting sharply with prior findings [10].

148 **LSH top-k retrieval.** Although oracle attention scores accurately identify key tokens, computing  $f$   
 149 is impractical for real-world applications. To address this, MagicPIG approximates  $f$  with  $\tilde{f}$  using  
 150 Locality-Sensitive Hashing (LSH) to efficiently retrieve critical KV entries. Specifically, a linear hash  
 151 function  $\mathcal{H}$  computes:

$$\tilde{f}(X, X') = \mathcal{H}(X) \otimes \mathcal{H}(X'), \quad (3)$$

152 where  $\otimes$  denotes a matrix multiplication-like operation, substituting floating-point multiplication  
 153 with NXOR. The indices  $\mathcal{I}$  of the top- $k$  largest values in  $\tilde{f}(Q, K_{\text{cache}})$  are used to retrieve the most  
 154 relevant KV entries.

155 LSH groups similar vectors into the same bucket with high probability, making it ideal for dense  
 156 vector spaces. A common variant employs random hyperplanes, where distinct hyperplanes create  
 157 linear decision boundaries, assigning data on either side to bit-0 or bit-1. The bits sequence from  
 158 all hyperplanes forms the hash code. In practice, these steps can be simplified to a single matrix  
 159 multiplication followed by a sign operation. For a vector  $x \in \mathbb{R}^d$ , we apply a random projection  
 160 matrix  $R \in \mathbb{R}^{d \times d_H}$  to obtain a hash code by taking the sign of the resulting product:

$$\mathcal{H}(x) = \text{sign}(xR). \quad (4)$$

161 **MLP hashing top-k retrieval.** As shown in Figure 2, queries and keys typically lie within two cone-  
 162 shaped regions in high-dimensional space [10]. This distribution causes uneven partitioning in LSH,  
 163 reducing encoding efficiency. To address this, we propose MLP hashing, a learned non-linear hashing  
 164 network tailored to query and key distributions. This approach enhances hash code information  
 165 density, enabling effective partitioning of skewed data through non-linear decision boundaries.

166 Specifically, we replace the projection matrix  $R$  in Eq. (4) with a two-layer MLP:

$$\text{MLP}(x) = W_2(\text{SiLU}(W_1x + b_1)), \quad (5)$$

167 where  $W_1$ ,  $b_1$ , and  $W_2$  are learnable parameters. Hash codes are then computed as:

$$\mathcal{H}(x) = \text{sign}(\text{MLP}(x)). \quad (6)$$

### 168 3.3 Optimization.

169 Optimizing the MLP hashing function  $\mathcal{H}$  to capture the query and key distribution is more critical  
 170 than the hashing network design and forms the *core contribution* of this work. We next introduce two  
 171 intuitive training objectives and explain their limitations in this context.

<sup>1</sup>See Appendix A.2 for the oracle top- $k$  attention pseudocode.

Table 2: **Perplexity comparison with Quest.** Perplexity evaluation on PG19 (#1), ProofPile (#2), and CodeParrot (#3) datasets. All models truncate inputs to their maximum supported token length. Spotlight Attention achieves performance comparable to Quest with a  $10\times$  smaller token budget.

Method	Configuration	Frozen Layers	LLaMA2-7B			LLaMA2-7B-Chat			LLaMA3-8B			Qwen2.5-7B		
			#1	#2	#3	#1	#2	#3	#1	#2	#3	#1	#2	#3
Vanilla Oracle Top-K	0% Pruned <b>98% Pruned</b>	N/A [0,1]	6.879 <b>6.941</b>	4.277 <b>4.317</b>	3.679 <b>3.729</b>	9.212 <b>9.224</b>	5.943 <b>5.891</b>	4.786 <b>4.795</b>	8.604 <b>8.881</b>	3.517 <b>3.590</b>	5.219 <b>5.353</b>	11.112 <b>10.435</b>	3.833 <b>3.587</b>	4.951 <b>4.733</b>
Quest	75% Pruned	[0,1]	7.116	4.404	3.854	9.282	5.930	4.879	9.912	4.024	5.893	10.485	4.281	5.482
	87.5% Pruned		7.735	4.754	4.054	9.820	6.226	5.084	12.434	4.927	6.646	13.596	5.471	6.213
	93.7% Pruned		9.746	5.775	4.571	12.058	7.440	5.686	17.320	6.749	8.693	19.852	7.274	7.823
	<b>96.9% Pruned</b>		<b>15.494</b>	<b>8.578</b>	<b>5.996</b>	<b>18.719</b>	<b>11.083</b>	<b>7.345</b>	<b>27.510</b>	<b>11.631</b>	<b>14.205</b>	<b>31.583</b>	<b>13.208</b>	<b>12.526</b>
Spotlight (ours)	80% Pruned	[0,1]	6.887	4.278	3.682	9.107	5.860	4.767	8.612	3.519	5.228	9.766	3.465	4.618
	90% Pruned		6.908	4.285	3.689	9.058	5.796	4.754	8.651	3.529	5.239	9.783	3.467	4.621
	95% Pruned <b>98% Pruned</b>		6.959 <b>7.106</b>	4.304 <b>4.364</b>	3.703 <b>3.768</b>	9.067 <b>9.262</b>	5.748 <b>5.770</b>	4.752 <b>4.806</b>	8.734 <b>8.977</b>	3.552 <b>3.621</b>	5.285 <b>5.434</b>	9.825 <b>9.930</b>	3.475 <b>3.497</b>	4.627 <b>4.645</b>

172 **X Language modeling loss.** A natural approach to optimize  $\mathcal{H}$  is to minimize the language modeling  
 173 loss directly through a few hundred post-training steps. However, this method has significant  
 174 limitations. First, it requires full forward and backward passes through the entire LLM, which  
 175 is computationally costly. More critically, differentiating the top- $k$  operator in Eq. (2) requires a  
 176 complex “soft top- $k$ ” approximation, which is challenging to implement.

177 **X Reconstruction loss.** An alternative approach uses MSE to align  $\tilde{f}$  with  $f$ , enabling layer-wise  
 178 optimization and reducing computational cost. However, this method has a key limitation: the  
 179 objective is to select which KV cache entries to retain, not to rank their relative importance. Training  
 180 with this loss misallocates capacity by prioritizing the ranking of excluded entries, deviating from the  
 181 core goal, and causing significant performance degradation. See Figure 3 (*Left*) for its limitations.

182 **✓ Our ranking loss.** To address this issue, we adopt a Bradley-Terry ranking objective, inspired  
 183 by RankNet [8]. As shown in Figure 3 (*Right*), during training, we identify top- $k$  and non-top- $k$   
 184 index sets based on attention scores. These indices split the estimated scores derived from Hamming  
 185 distances of query and key hash codes into sets  $B$  and  $C$ . An optimizer then updates the hashing  
 186 function parameters to ensure each score in  $B$  exceeds every score in  $C$ :

$$\mathcal{L}_{\text{rank}} = -\frac{1}{k(n-k)} \sum_{i,j} \log (\text{sigmoid}(\beta(B_i - C_j) - \alpha)), \quad (7)$$

187 where  $\beta$  and  $\alpha$  are positive constants used to amplify the separation between  $B$  and  $C$ , facilitating  
 188 convergence. The core of this loss design lies in filtering out supervising signals related to internal  
 189 ranking within  $B$  and  $C$ , effectively addressing the issue of capacity misallocation. We provide the  
 190 pseudo-code for our ranking loss in Appendix A.3 to aid readers familiar with code.

191 **Make hashing function differentiable.** After computing the loss, errors can be backpropagated to  
 192 the MLP hashing function. However, the sign function’s non-differentiability blocks gradient flow.  
 193 To resolve this, we substitute the sign function with a soft sign function during training:

$$\text{softsign}(x) = \frac{\gamma x}{1 + \gamma |x|}, \quad (8)$$

194 where  $\gamma \in \mathbb{R}$  is a hyperparameter controlling the extent of smoothing. This soft sign function is used  
 195 only during training. In inference, the non-differentiable sign function is reinstated.

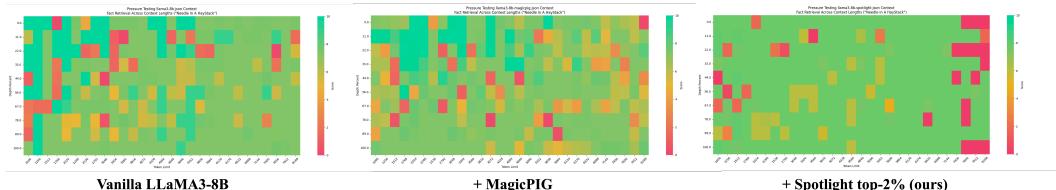


Figure 4: **NIAH results.** Using LLaMA3-8B [13] as the base model, we compared the retrieval accuracy of MagicPIG with our method. Our approach, which relies solely on hash code-based retrieval without local windows or sink tokens, achieves comparable response accuracy.

Table 3: **Perplexity versus MagicPIG.** Comparison of perplexity on PG19 (#1), ProofPile (#2), and CodeParrot (#3). Due to the time-consuming evaluation process of MagicPIG, we sampled only 10 data points from each dataset for testing.

Method	Configuration				LLaMA3-8B		
	Frozen	Local (64)	Sink (4)	Retrieve (2%)	#1	#2	#3
Vanilla Oracle Top-K [0,1]				✓	9.56	2.83	2.26
MagicPIG [0,16]	✓	✓	✓	✓	12.65	3.54	2.91
	✓	✓	✓	✓	16.94	6.27	4.57
			✓	✓	50.96	8.52	6.67
			✓	✓	42.12	8.65	10.43
Spotlight (ours) [0,1]	✓	✓	✓	✓	NaN	NaN	NaN
				✓	9.87	2.89	2.29
					13.89	5.50	3.98
					9.99	2.91	2.30

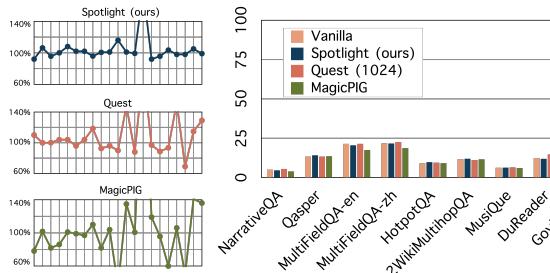


Figure 5: **Downstream QA tasks.** (Left) Relative score of each method compared to the vanilla baseline; each point denotes a subtask. (Right) Absolute score comparison.

## 4 Experimentation

Our evaluation spans multiple dimensions. (1) We first assess the errors introduced by retrieval and sparsification, measured by retrieval accuracy and perplexity score. (2) We then evaluate long-context key information retrieval using the Needle-in-a-Haystack [14] benchmark. (3) Next, we evaluated downstream QA tasks on LongBench [6], comparing response similarity between compressed and original models using Rouge-L. The key insight from this Rouge-L comparison is straightforward: higher-quality compression produces more concise outputs. (4) We also measure end-to-end throughput gains and the execution efficiency of our CUDA ops.

To evaluate generalization, we tested the performance of the Qwen2.5 series [24] across various model sizes and LLaMA3-8B [13] across diverse training corpora. Additionally, we conducted ablation studies to assess the impact of loss functions and attention estimation methods, with details provided in Appendix B.5 and B.6.

**Experimental setups.** We employ LLaMA3-8B [13], LLaMA2-7B, LLaMA2-7B-Chat [21], and Qwen2.5 models (1.5B, 7B, 14B) [24] as base models. The baseline methods compared include (1) oracle top- $k$  retrieval, (2) linear LSH top- $k$ , (3) Quest, and (4) MagicPIG. Their implementation details are provided in Appendix A.2.

Our MLP hashing function employs 128-dimensional input, intermediate, and output layers, with a distinct MLP for each head in every layer, producing a 128-bit hash code—much shorter than MagicPIG’s minimum of 720 bits. Only the hashing functions are trainable. Training data consists of 8,192 samples, evenly drawn from the Book and Arxiv datasets [22].

To improve efficiency, hidden states for all layers are precomputed and stored, enabling independent layer-wise training without joint fine-tuning. Training uses  $\gamma = 64$ , a learning rate of  $1 \times 10^{-3}$ ,  $\beta = 1$ , and  $\alpha = 3$ , for one epoch. Additional details are provided in Appendix A.1. The pruning rate remains fixed at 98% during training, irrespective of evaluation settings.

### 4.1 Main Results

**KV retrieval accuracy.** This experiment compares linear LSH and MLP hashing for KV retrieval. We used the first sample from PG19 as test data and assessed KV retrieval accuracy with the average

Table 4: **QA response fidelity.** On LongBench, output fidelity (measured by Rouge-L between compressed and vanilla model outputs) shows our method achieves performance closest to the vanilla model.

Method	Configuration				Similarity
	Local	Sink	Retrieval	Frozen	
Vanilla	✗	✗	✗	✗	1.00
Oracle Top-K	✗	✗	✗	✗	0.66
	✗	✗	✗	✗	0.37
	✗	✗	✗	[0,1]	0.56
	✗	✗	✗	256	0.34
MagicPIG	64	4	Dynamic	[0,16]	0.44
Spotlight (ours)	✗	✗	✗	[0,1]	<b>0.58</b>

196

Figure 5: **Downstream QA tasks.** (Left) Relative score of each method compared to the vanilla baseline; each point denotes a subtask. (Right) Absolute score comparison.

197

## 4 Experimentation

198

Our evaluation spans multiple dimensions. (1) We first assess the errors introduced by retrieval and sparsification, measured by retrieval accuracy and perplexity score. (2) We then evaluate long-context key information retrieval using the Needle-in-a-Haystack [14] benchmark. (3) Next, we evaluated downstream QA tasks on LongBench [6], comparing response similarity between compressed and original models using Rouge-L. The key insight from this Rouge-L comparison is straightforward: higher-quality compression produces more concise outputs. (4) We also measure end-to-end throughput gains and the execution efficiency of our CUDA ops.

205

To evaluate generalization, we tested the performance of the Qwen2.5 series [24] across various model sizes and LLaMA3-8B [13] across diverse training corpora. Additionally, we conducted ablation studies to assess the impact of loss functions and attention estimation methods, with details provided in Appendix B.5 and B.6.

209

**Experimental setups.** We employ LLaMA3-8B [13], LLaMA2-7B, LLaMA2-7B-Chat [21], and Qwen2.5 models (1.5B, 7B, 14B) [24] as base models. The baseline methods compared include (1) oracle top- $k$  retrieval, (2) linear LSH top- $k$ , (3) Quest, and (4) MagicPIG. Their implementation details are provided in Appendix A.2.

213

Our MLP hashing function employs 128-dimensional input, intermediate, and output layers, with a distinct MLP for each head in every layer, producing a 128-bit hash code—much shorter than MagicPIG’s minimum of 720 bits. Only the hashing functions are trainable. Training data consists of 8,192 samples, evenly drawn from the Book and Arxiv datasets [22].

217

To improve efficiency, hidden states for all layers are precomputed and stored, enabling independent layer-wise training without joint fine-tuning. Training uses  $\gamma = 64$ , a learning rate of  $1 \times 10^{-3}$ ,  $\beta = 1$ , and  $\alpha = 3$ , for one epoch. Additional details are provided in Appendix A.1. The pruning rate remains fixed at 98% during training, irrespective of evaluation settings.

221

### 4.1 Main Results

222

**KV retrieval accuracy.** This experiment compares linear LSH and MLP hashing for KV retrieval. We used the first sample from PG19 as test data and assessed KV retrieval accuracy with the average

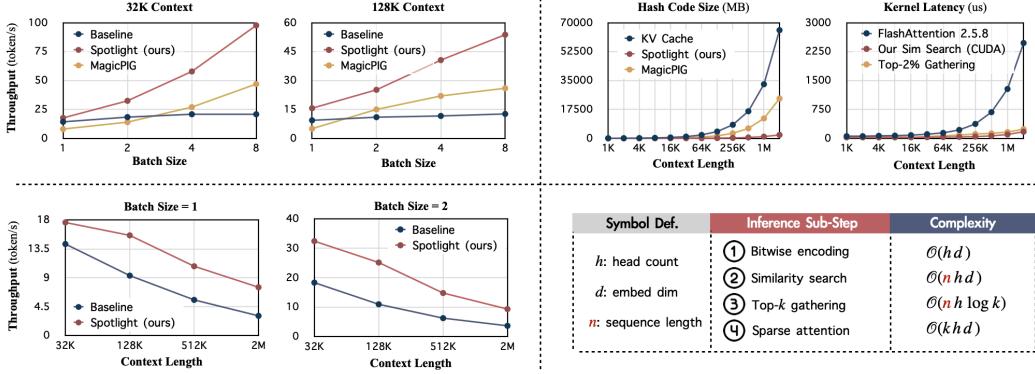


Figure 6: **Efficiency.** (Upper Left) End-to-end throughput comparison with fixed context length across varying batch sizes. (Bottom Left) End-to-end throughput comparison with fixed batch size across different context lengths. (Upper Right) Hash code size comparison between MagicPIG and our method, alongside the execution latency of the two most computationally intensive operations in our method. (Bottom Right) Complexity comparison of different computational steps.

224 IoU score across all heads and layers. IoU is computed as the intersection of the top- $k$  indices  
225 retrieved by the algorithm and the oracle top- $k$  indices, divided by their union.

226 For LSH, we initialized the projection matrix via QR decomposition (see Appendix A.8). For MLP  
227 hashing, parameters were initialized with random Gaussian distributions. Results in Table 1 show  
228 that LSH performs best without training but shows minimal improvement post-training. Conversely,  
229 MLP hashing markedly improves retrieval accuracy after training, achieving the highest performance.

230 **Language modeling.** We evaluated three language modeling benchmarks—PG19 [19], ProofPile [5],  
231 and CodeParrot [18]—with 100, 79, and 100 samples, respectively, using perplexity to detect minor  
232 errors from sparsification. Additional experimental details are in Appendix A.5.

233 Results in Table 2 show our method outperforms Quest with a tenfold reduction in token budget for  
234 most models. Comparisons with MagicPIG (Table 3) indicate that MagicPIG depends heavily on  
235 local windows and sink tokens, failing without them, while our method autonomously identifies these  
236 elements, demonstrating significantly higher retrieval accuracy.

237 MagicPIG excels at retrieving facts in scenarios like Needle-in-a-Haystack [14], its contribution  
238 should not be dismissed solely on its limitations in this language modeling test.

239 **Needle-in-a-Haystack.** We use the offline evaluation version of NIAH [1], which *differs* from  
240 ChatGPT scoring by using the Rouge score to measure output accuracy. We utilize LLaMA3-8B [13]  
241 as the base model, evaluating context lengths ranging from 256 to 8192, with a step size of 256 for  
242 testing. The Needle dataset is highly prompt-sensitive, we provide the needle, retrieval question, and  
243 haystack context in Appendix A.4. Additionally, to reduce output variability, all evaluated methods  
244 use greedy search. As shown in Figure 4, Spotlight Attention achieves performance on par with the  
245 original model.

246 **Downstream QA tasks.** We evaluated the performance of various compression methods on Long-  
247 Bench [6] subtasks. Quest employs a token budget of 1,024, while MagicPIG uses 4 sink tokens,  
248 a 64-token local window, and a 1,500-bit hash code per key, following their default configurations.  
249 Our method uses a token budget of 163. All experimental results reported in the main text use  
250 LLaMA3-8B as the base model, scores of other models are provided in Appendix B.4, and more  
251 experimental setup details are in Appendix A.6. Figure 5 (Right) compares absolute scores across all  
252 subtasks. More importantly, Figure 5 (Left) shows relative scores, revealing that our method’s scores  
253 closely align with those of the vanilla model.

254 We also assessed output fidelity using Rouge-L to measure similarity between the vanilla model’s  
255 outputs and those of these methods. Results are presented in Table 4. Our method’s outputs are the  
256 most similar to the vanilla model, with the longest consecutive subsequence match approaching 60%.  
257 In contrast, Quest requires retrieving six times more tokens to achieve comparable similarity. These  
258 findings hold for nearly all subtasks, with detailed scores of each subtask provided in Appendix B.2.

259 **4.2 Efficiency**

260 All efficiency experiments were performed on Qwen2.5-7B [24] using eight A100 GPUs. For  
 261 enhanced flexibility, experiments utilized the HuggingFace Transformers framework, optimized with  
 262 pipeline parallelization and KV cache pre-allocation to boost throughput.

263 **End-to-end throughput evaluation.** To evaluate model throughput at extended context lengths (e.g.,  
 264 2M tokens), we expanded positional encoding, disregarding output quality. We selected the first  
 265 sample from the PG19 test set [19], repeating it to reach a 2M-token context. GPU execution time  
 266 was measured using CUDA events. We generated eight consecutive tokens, computed throughput by  
 267 dividing by generation time, and averaged three runs per data point. Results in Figure 6 (*Left*) show  
 268 that our approach consistently delivers throughput gains, especially at larger input scales.

269 **Kernel evaluation.** We implemented CUDA kernels for both bit-packing and NXOR GEMM. Bit-  
 270 packing compresses 32 Torch boolean values into a single unsigned 32-bit integer, as detailed in  
 271 Appendix A.7. For NXOR GEMM, we utilized the standard library’s popcount to count bit-1s in  
 272 each NXOR result. For top- $k$  gathering and sparse attention, we employed Torch and FlashAttention  
 273 implementations, respectively. As shown in Figure 6 (*Upper Right*), compared to dense vectors, our  
 274 hashing-based similarity search significantly reduces storage and computation, enabling bit-packing  
 275 and similarity search within 100 $\mu$ s for context lengths up to 512K.

276 **Table 5: Ablation on model size.** Various  
 277 Qwen2.5 model sizes, augmented with Spotlight  
 278 Attention, demonstrated better perplexity across  
 279 diverse language modeling tasks.

Model	Method	PG19	Math	Code
Qwen2.5-1.5B	Vanilla	13.828	4.181	5.081
	Spotlight Top-2% (ours)	13.510	4.143	5.064
Qwen2.5-7B	Vanilla	11.112	3.833	4.951
	Spotlight Top-2% (ours)	9.930	3.497	4.645
Qwen2.5-14B	Vanilla	8.416	3.230	4.472
	Spotlight Top-2% (ours)	8.261	3.196	4.440

276 **Table 6: Ablation on training tasks.** Along-  
 277 side the standard ArXiv+Books training data,  
 278 we trained models on the C4 and GitHub Code  
 279 datasets, achieving comparable perplexity (PPL)  
 280 results.

LLaMA3-8B	Training Corpus	PG19	Math	Code
Oracle Top-2%	-	8.881	3.590	5.353
Spotlight Top-2% (ours)	Arxiv + Book (default)	8.977	3.621	5.434
	C4	8.958	3.631	5.417
	Github Code	<b>8.891</b>	<b>3.611</b>	<b>5.393</b>

277 **4.3 Ablations**

278 In the main text, we present ablation studies on model size and training tasks only. Additional ablation  
 279 experiments, including loss functions and attention estimation methods, are detailed in Appendix B.5  
 280 and B.6, respectively.

281 All ablation experiments use language modeling perplexity as the evaluation metric. We assessed  
 282 performance on PG19 [19], Proof-Pile (Math) [5], and CodeParrot (Code) [18], with sample sizes of  
 283 100, 79, and 100, respectively.

284 **Model size.** We compare Qwen2.5 [24] models of 1.5B, 7B, and 14B parameters, all trained with the  
 285 standard recipe. As shown in Table 5, Spotlight Attention achieves consistently strong performance  
 286 across different model sizes.

287 **Training tasks.** To validate our method’s applicability across diverse tasks, we trained models on  
 288 the GitHub Code [18] and C4 [18] datasets, in addition to the ArXiv and Books [22] datasets used  
 289 previously. As presented in Table 6, training with GitHub Code or C4 unexpectedly outperformed  
 290 the ArXiv+Books combination, demonstrating the robust adaptability of our training framework.  
 291 However, due to the already-completion of the main results, we did not re-evaluate all benchmarks  
 292 with these checkpoints despite their superior performance.

293 **5 Conclusion and Limitation**

294 We introduce Spotlight Attention, an advancement over Quest and MagicPIG, incorporating a non-  
 295 linear hashing function and an optimized framework. This approach addresses the underfitting  
 296 of MagicPIG’s linear hashing while significantly reducing hash code length. Spotlight Attention  
 297 performs well on downstream tasks but has limitations. The IoU remains around 40% despite  
 298 non-linear hashing, indicating potential for improvement.

299 **References**

- 300 [1] Github repository: 66ring/llmtest\_needleinahaystack-local, 2023.
- 301 [2] M. Adnan, A. Arunkumar, G. Jain, P. Nair, I. Soloveychik, and P. Kamath. Keyformer: Kv  
302 cache reduction through key tokens selection for efficient generative inference. *MLSys*, 2024.
- 303 [3] A. Agrawal, A. Panwar, J. Mohan, N. Kwatra, B. S. Gulavani, and R. Ramjee. Sarathi: Efficient  
304 llm inference by piggybacking decodes with chunked prefills. *arXiv*, 2023.
- 305 [4] A. Andoni, P. Indyk, T. Laarhoven, I. Razenshteyn, and L. Schmidt. Practical and optimal lsh  
306 for angular distance. In *NeurIPS*, 2015.
- 307 [5] Z. Azerbayev, E. Ayers, and B. Piotrowski. Github repository: hoskison-center/proof-pile, 2022.
- 308 [6] Y. Bai, X. Lv, J. Zhang, H. Lyu, J. Tang, Z. Huang, Z. Du, X. Liu, A. Zeng, L. Hou, et al.  
309 LongBench: A bilingual, multitask benchmark for long context understanding. In *ACL*, 2024.
- 310 [7] R. A. Bradley and M. E. Terry. Rank analysis of incomplete block designs: I. the method of  
311 paired comparisons. *Biometrika*, 1952.
- 312 [8] Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tsai, and H. Li. Learning to rank: from pairwise approach to  
313 listwise approach. In *ICML*, 2007.
- 314 [9] Y. Chen, G. Wang, J. Shang, S. Cui, Z. Zhang, T. Liu, S. Wang, Y. Sun, et al. NACL: A general  
315 and effective KV cache eviction framework for LLM at inference time. In *NACL*, 2024.
- 316 [10] Z. Chen, R. Sadhukhan, Z. Ye, Y. Zhou, J. Zhang, N. Nolte, Y. Tian, M. Douze, L. Bottou,  
317 Z. Jia, and B. Chen. MagicPIG: LSH sampling for efficient LLM generation. In *NeurIPS*, 2024.
- 318 [11] L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, et al. Github repository:  
319 Eleutherai/Im-evaluation-harness, 2024.
- 320 [12] S. Ge, Y. Zhang, L. Liu, M. Zhang, J. Han, and J. Gao. Model tells you what to discard:  
321 Adaptive KV cache compression for LLMs. In *ICLR*, 2024.
- 322 [13] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, et al. The  
323 llama 3 herd of models. *arXiv*, 2024.
- 324 [14] G. Kamradt. Github repository: gkamradt/llmtest\_needleinahaystack, 2023.
- 325 [15] H. Li, Y. Li, A. Tian, T. Tang, Z. Xu, X. Chen, N. Hu, W. Dong, Q. Li, and L. Chen. A survey  
326 on large language model acceleration based on kv cache management. *arXiv*, 2025.
- 327 [16] Y. Li, Y. Huang, B. Yang, B. Venkitesh, A. Locatelli, H. Ye, T. Cai, P. Lewis, and D. Chen.  
328 SnapKV: LLM knows what you are looking for before generation. In *NeurIPS*, 2024.
- 329 [17] M. Oren, M. Hassid, N. Yarden, Y. Adi, and R. Schwartz. Transformers are multi-state rnns.  
330 *arXiv*, 2024.
- 331 [18] Z. Peitian. Huggingface dataset: namespace-pt/long-llm-data, 2024.
- 332 [19] J. W. Rae, A. Potapenko, S. M. Jayakumar, C. Hillier, and T. P. Lillicrap. Compressive  
333 transformers for long-range sequence modelling. In *ICLR*, 2020.
- 334 [20] J. Tang, Y. Zhao, K. Zhu, G. Xiao, B. Kasikci, and S. Han. QUEST: Query-Aware Sparsity for  
335 Efficient Long-Context LLM Inference. In *ICML*, 2024.
- 336 [21] H. Touvron, L. Martin, K. Stone, T. Scialom, et al. Llama 2: Open foundation and fine-tuned  
337 chat models. *arXiv*, 2023.
- 338 [22] M. Weber, D. Y. Fu, Q. G. Anthony, Y. Oren, S. Adams, A. Alexandrov, X. Lyu, et al. Redpajama:  
339 an open dataset for training large language models. In *NeurIPS*, 2024.
- 340 [23] G. Xiao, Y. Tian, B. Chen, S. Han, and M. Lewis. Efficient streaming language models with  
341 attention sinks. In *ICLR*, 2024.

342 [24] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, et al. Qwen2.5  
343 technical report. *arXiv*, 2025.

344 [25] Z. Zhang, Y. Sheng, T. Zhou, T. Chen, L. Zheng, R. Cai, Z. Song, Y. Tian, C. Re, C. Barrett,  
345 Z. Wang, and B. Chen. H2o: Heavy-hitter oracle for efficient generative inference of large  
346 language models. In *NeurIPS*, 2023.

347 **A Implementation Details**

348 **A.1 Training Setup**

349 Table 7 summarizes the training configuration. For further details, including low-level specifics and weight files, refer to our open-source code.

Table 7: Detailed training configuration.

General			Learning Rate					Gradient	
Precision	Num Iters	Batch Size	Max LR	Min LR	Warm Up Iters	Warm Up Method	Annealing	Accumulation	Clipping
<hr/>									
Optimizer	$\beta_1$	$\beta_2$	Weight Decay	Corpus	Arxiv Samples	Book Samples	LLaMA2 Trunc	LLaMA3/Qwen Trunc	Trunc Side
<hr/>									
adamw	0.9	0.98	0.1	arxiv, book	4,096	4,096	4,096	8,192	right

350

351 **A.2 Baseline Methods**

352 **Oracle top-k.** As introduced in Section 3.2, this baseline employs original attention to identify top- $k$  indices, serving as a theoretical upper bound on performance. We provide pseudo-code below to illustrate a concrete implementation of the oracle top- $k$  algorithm.

355 **LSH top-k.** This method uses training-free angular LSH [4] with the same hash code length as 356 Spotlight Attention, providing a reference for MagicPIG’s linear hashing under our framework. 357 Initialization details for angular LSH are in Appendix A.8.

358 **Quest.** [20] We adopt Quest’s official implementation with default hyperparameters, modifying only 359 the token budget. Beyond the default 1024-token budget, we test ultra-low budgets to align with 360 Spotlight Attention. For instance, for LLaMA2-7B and LLaMA3-8B, we use budgets of 128 and 256 361 tokens, compared to Spotlight Attention’s 81 and 163 tokens.

362 **MagicPIG.** [10] We adopt MagicPIG’s official implementation with default hyperparameters ( $K = 363 15, L = 100$ ), retaining 64 local tokens and 4 initial tokens. For all experiments except the efficiency 364 test, we use the official Python evaluation code to conduct the experiments.

365 **A.3 Ranking Loss**

366 We provide the ranking loss calculation in the following pseudo-code. To support longer sequence 367 lengths during training, we employ three optimization techniques: (1) Random query selection, where 368 only queries specified by `query_index` are optimized, rather than all queries. (2) Random top- $k$  369 selection, where `max_top` is randomly sampled from the top- $k$  set for optimization. (3) Random 370 non-top- $k$  selection, where `max_oth` is randomly sampled from the non-top- $k$  set for optimization. 371 These techniques enhance training efficiency in long-context scenarios.

```

372   1 def ranking_loss(
373   2     draft_attn,
374   3     true_attn,
375   4     query_index,
376   5     max_top,
377   6     max_oth,
378   7     maskout,
379   8     beta: float = 1.0,
380   9     alpha: float = 0.0):
381   10
382   11     loss = torch.tensor(0, dtype=torch.float32)
383   12     criterion = torch.nn.BCEWithLogitsLoss()
384   13
385   14     # PREPARE & APPLY MASK
386   15     num_kv = true_attn.shape[-1]
387   16     mask = torch.triu(torch.ones((num_kv, num_kv), dtype=torch.bool, device=true_attn.device), diagonal=1)[None, None,
388   17     :, :]
389   18     if query_index is not None:
390   19         mask = mask[:, query_index, :]
391   20     true_attn = torch.masked_fill(true_attn, mask, value=torch.finfo(true_attn.dtype).min)
392   21
393   22     indices = torch.argsort(true_attn, dim=-1, descending=True)
394   23     top_cnt = int(indices.shape[-1] * (1 - maskout))
395   24     top_indices = indices[:, :top_cnt]
396   25     oth_indices = indices[:, top_cnt:]
397   26

```

```

400    27      if max_top is not None:
401    28          top_rnd_indices = torch.randperm(top_cnt, dtype=torch.int64, device=indices.device)[:,max_top]
402    29          top_indices = top_indices[..., top_rnd_indices]
403    30      if max_oth is not None:
404    31          oth_rnd_indices = torch.randperm(indices.shape[-1] - top_cnt, dtype=torch.int64, device=indices.device)[:,max_oth]
405    32          oth_indices = oth_indices[..., oth_rnd_indices]
406
407    33
408    34      top_mask = torch.gather(mask.expand_as(true_attn), dim=-1, index=top_indices)[..., :, None]
409    35      oth_mask = torch.gather(mask.expand_as(true_attn), dim=-1, index=oth_indices)[..., None, :]
410    36
411    37      top_draft_attn = torch.gather(draft_attn, dim=-1, index=top_indices)[..., :, None]
412    38      oth_draft_attn = torch.gather(draft_attn, dim=-1, index=oth_indices)[..., None, :]
413    39
414    40      residual = top_draft_attn - oth_draft_attn
415    41      residual_mask = (top_mask | oth_mask).expand_as(residual).flatten(-3)
416    42
417    43      logits = residual.flatten(-3)[~residual_mask.bool()]
418    44      labels = torch.ones_like(logits)
419    45      loss += criterion(logits * beta - alpha, labels).cpu()
420    46
421    47      diff = torch.count_nonzero(logits < 0) / logits.numel()
422    48
423    49      return diff, loss

```

#### 425 A.4 Needle-in-a-Haystack

426 NIAH is a prompt-sensitive test focusing on relative performance changes before and after applying  
 427 Spotlight Attention, rather than absolute performance. The haystack, needle, and question used in  
 428 our evaluation are depicted in Figure 7. The prompt, selected from a set proven effective in practice,  
 is shown in Figure 8.

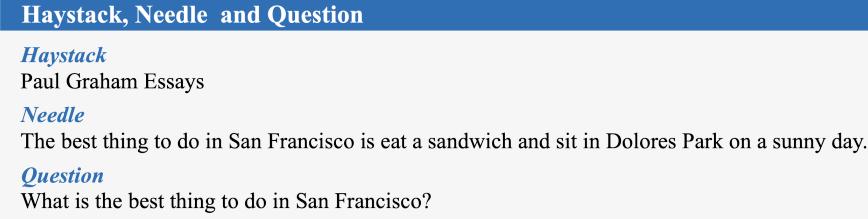


Figure 7: The haystack, needle, and retrieval question for the NIAH.

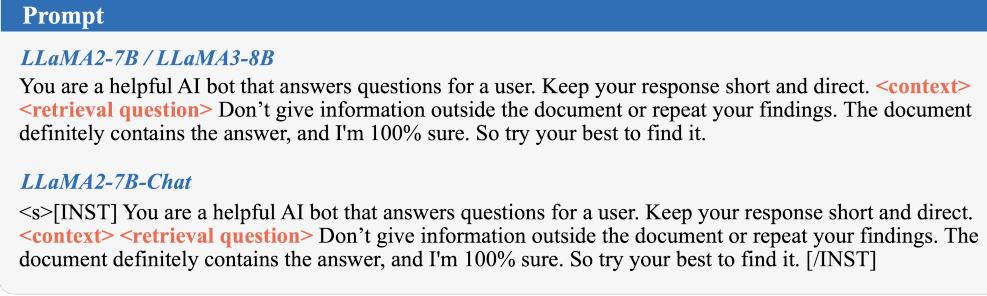


Figure 8: Prompts used in NIAH test.

429

#### 430 A.5 Language Modeling Perplexity

431 For Quest, we used the first 128 tokens for pre-filling with full attention. For MagicPIG, we pre-filled  
 432 the first 1024 tokens with full attention to compute the key's mean value. Our Spotlight Attention  
 433 method employed sparse KV retrieval throughout without pre-filling. For LLaMA2 models, we  
 434 evaluated the first 4K tokens per sample; for LLaMA3 and Qwen2.5 models, we used the first 8K  
 435 tokens. We compared Quest and MagicPIG separately, as Quest and our method use a fixed token  
 436 budget, while MagicPIG dynamically selects tokens and uses local windows and sink tokens.

437 In comparison with Quest, we calculated the token budget by multiplying the token sequence length  
 438 by the pruning rate, with a minimum budget of 20.

#### 439 A.6 Downstream QA Tasks

440 For contexts exceeding 8K tokens, we truncated 8K tokens from right to left. We evaluated all  
 441 LongBench [6] subdatasets using the official test scripts. The LLaMA2 chat model employs the  
 442 official chat template, whereas other models do not.

#### 443 A.7 Bit-Packing CUDA Kernel

444 Bit-packing (Figure 9) is crucial because PyTorch lacks a native bit type, and boolean values are  
 445 stored as full bytes. Without compaction, storage usage would increase substantially. The bit-packing  
 446 program groups 32 boolean values and iteratively packs them into a single `uint_32t`, as detailed in  
 the pseudocode accompanying Figure 9.

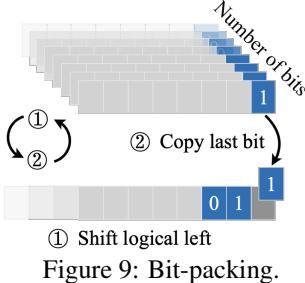


Figure 9: Bit-packing.

---

```

1 def bit_packing(tensor):
2     # WE PRESENT THE LOGIC HERE,
3     # WITH THE ACTUAL IMPLEMENTATION WRITTEN IN CUDA.
4     n, d = tensor.shape
5     assert d % 32 == 0
6
7     tensor = tensor.chunk(32, dim=-1)
8     output = torch.zeros((n, d // 32), dtype=uint32)
9     for x in tensor:
10         output <= 1
11         output |= x & 0x01
12
13     return output

```

---

447

#### 448 A.8 LSH Weight Initialization

449 For the linear hashing function, we employ angular LSH with a rotation matrix as the initial value,  
 450 outperforming standard random initialization. To generate a  $d$ -dimensional rotation matrix, we use  
 451 QR decomposition. We first create a random matrix  $R \in \mathbb{R}^{d \times d}$ , with each element independently  
 452 sampled from a standard normal distribution:

$$R \in \mathbb{R}^{n \times n}, R_{i,j} \sim \mathcal{N}(0, 1). \quad (9)$$

453 We then perform QR decomposition on  $R$ , yielding an orthogonal matrix  $\mathbf{Q}$  and an upper triangular  
 454 matrix  $\mathbf{R}$ :

$$R = \mathbf{Q}\mathbf{R}, \quad \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}. \quad (10)$$

455 The matrix  $\mathbf{Q}$  is not necessarily in the special orthogonal group  $\text{SO}(d)$ , as its determinant can be  
 456 either +1 or -1. To ensure  $\mathbf{Q} \in \text{SO}(d)$ , we flip the sign of the first column of  $\mathbf{Q}$  if  $\det(\mathbf{Q}) < 0$ .

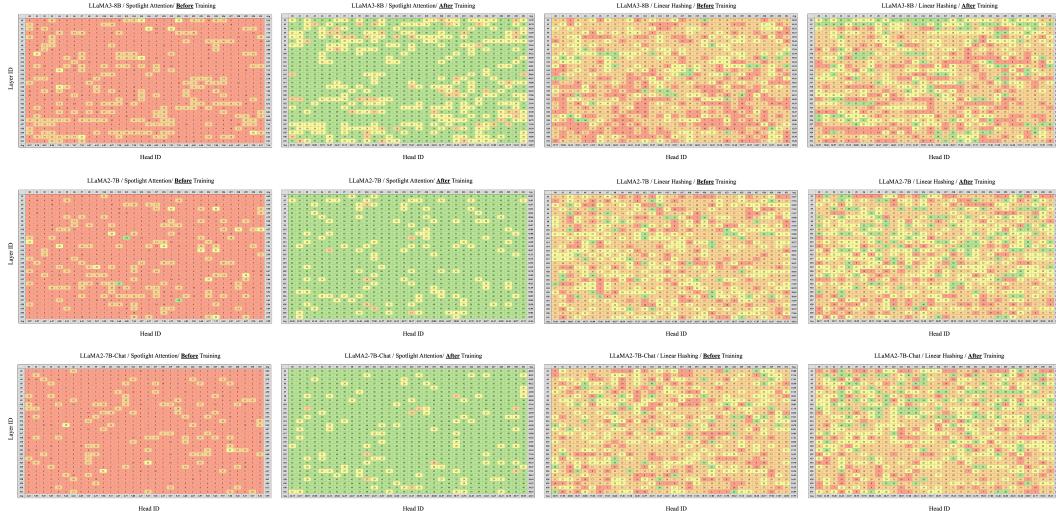
## 457 B More Experimental Results

### 458 B.1 Per-Head Retrieval Accuracy

459 The IoU reported in the experiments section is averaged across all heads and layers. Given the  
 460 varied behavior of LLM heads, individual IoU scores differ significantly. Figure 10 presents detailed  
 461 per-head IoU comparisons for various models, using Spotlight Attention and linear hashing functions,  
 462 both before and after training.

### 463 B.2 Detailed QA Response Fidelity Scores

464 In the main text, we report the average Rouge-L score across all subdatasets as the similarity  
 465 metric. However, scores vary significantly across individual subdatasets. Figure 11 provides detailed  
 466 similarity scores for each method compared to the original LLaMA3-8B model. To avoid confusion,  
 467 we assign a Rouge-L score of 1 to identical outputs, except in special cases like outputs containing  
 468 only \n characters, where tokenization issues may prevent a perfect score.



**Figure 10: Per-head retrieval accuracy.** Measured by Intersection over Union (IoU) of top- $k$  sets predicted by oracle and Spotlight. (1) The results of linear hashing reveal that most heads maintain low IoU both pre- and post-training, suggesting their latent distributions are challenging to approximate with linear functions. (2) Random parameter initialization results in low before-training IoU for Spotlight Attention, which improves significantly after training.

### 469 B.3 Few-Shot Learning

- 470 To evaluate Spotlight Attention with short context lengths, we restrict the KV cache to 20 tokens  
 471 and assess performance on 5-shot learning datasets from LM-Eval-Harness [11], including GLUE,  
 472 SuperGLUE, OpenBookQA, HellaSwag, PiQA, Winogrande, ARC-E, ARC-C, MathQA, and MMLU.  
 473 As shown in Table 8, Spotlight Attention delivers strong performance. Comparisons with Quest and  
 474 MagicPIG are omitted, as their large local windows or budgets consume most of the prompt length  
 475 on most datasets, rendering such comparisons less meaningful. Instead, we emphasize the relative  
 476 performance between Spotlight Attention and the original model.

Table 8: A 5-shot learning comparison between original model and Spotlight Attention under a fixed budget of 20 tokens was conducted across: GLUE (1), SuperGLUE (2), OpenBookQA (3), HellaSwag (4), PiQA (5), Winogrande (6), ARC-E (7), ARC-C (8), MathQA (9), and MMLU (10).

	#1	#2	#3	#4	#5	#6	#7	#8	#9	#10	Win Rate
<b>LLaMA2-7B</b>	0.489	0.646	0.342	0.604	0.776	0.733	0.793	0.478	0.262	0.468	56%
+Spotlight	0.494	0.632	0.336	0.581	0.786	0.748	0.793	0.469	0.284	0.466	44%
<b>LLaMA2-7B-Chat</b>	0.625	0.717	0.346	0.598	0.763	0.729	0.811	0.474	0.295	0.482	89%
+Spotlight	0.582	0.664	0.35	0.578	0.759	0.703	0.809	0.467	0.275	0.482	11%
<b>LLaMA3-8B</b>	0.618	0.727	0.378	0.631	0.804	0.772	0.85	0.534	0.418	0.663	38%
+Spotlight	0.620	0.736	0.378	0.624	0.804	0.773	0.851	0.533	0.411	0.664	62%

### 477 B.4 Downstream QA Tasks

- 478 The main text primarily reports LLaMA3-8B scores on LongBench. Here, we additionally present  
 479 results for LLaMA2-7B and LLaMA2-7B-Chat, alongside comparisons of our method and Quest  
 480 under varying token budgets, as shown in Table 9.

### 481 B.5 Ablation on Loss Function

- 482 In the previous section, we analyzed the limitations of reconstruction loss. Here, we compare its  
 483 empirical performance with our proposed ranking loss. As shown in Table 10, our ranking loss  
 484 significantly outperforms reconstruction loss, which provides minimal to no benefit.

	<b>LLaMA3-8B</b>	<b>Upper Bound</b> 98% Pruned	<b>Linear Hashing</b> 98% Pruned	<b>Spotlight</b> 98% Pruned	<b>Quest</b> 1024 Budget	<b>Quest</b> 256 Budget	<b>MagicPIG</b>
NarrativeQA	0.95	0.67	0.43	0.59	0.53	0.36	0.58
Qasper	1.00	0.57	0.34	0.54	0.57	0.37	0.48
MultiFieldQA-en	0.93	0.71	0.43	0.62	0.62	0.43	0.49
MultiFieldQA-zh	0.64	0.58	0.24	0.47	0.52	0.33	0.33
HotpotQA	1.00	0.76	0.48	0.66	0.70	0.49	0.69
2WikiMultihopQA	1.00	0.83	0.53	0.75	0.77	0.61	0.69
MuSiQue	1.00	0.78	0.50	0.69	0.72	0.56	0.73
DuReader	0.43	0.69	0.18	0.55	0.37	0.16	0.25
GovReport	1.00	0.58	0.26	0.55	0.51	0.23	0.34
QMSum	1.00	0.48	0.33	0.44	0.39	0.28	0.48
MultiNews	0.28	0.82	0.60	0.79	0.84	0.05	0.03
VCSUM	0.21	0.46	0.19	0.39	0.41	0.06	0.10
TREC	1.00	0.71	0.60	0.63	0.64	0.51	0.70
TriviaQA	1.00	0.66	0.31	0.53	0.55	0.34	0.51
SAMSum	1.00	0.54	0.30	0.48	0.45	0.28	0.34
LSHT	0.07	0.39	0.02	0.26	0.18	0.01	0.02
PassageCount	1.00	0.82	0.58	0.80	0.65	0.59	0.70
PassageRetrieval-en	1.00	0.71	0.41	0.55	0.60	0.42	0.41
PassageRetrieval-zh	1.00	0.66	0.33	0.60	0.61	0.36	0.46
LCC	0.81	0.67	0.34	0.60	0.70	0.40	0.46
RepoBench-P	0.70	0.71	0.31	0.59	0.47	0.30	0.41
<b>Average</b>	0.81	0.66	0.37	0.58	0.56	0.34	0.44

Figure 11: Detailed similarity between (i) the outputs of different models (including LLaMA3-8B itself) and (ii) those of LLaMA3-8B.

Table 9: **Absolute scores on LongBench.** Performance comparison of different methods on LongBench’s long-text downstream tasks: NarrativeQA (1-1), Qasper (1-2), MultiFieldQA-en (1-3), MultiFieldQA-zh (1-4), HotpotQA (2-1), 2WikiMultihopQA (2-2), MuSiQue (2-3), DuReader (2-4), GovReport (3-1), QMSum (3-2), MultiNews (3-3), VCSUM (3-4), TREC (4-1), TriviaQA (4-2), SAMSum (4-3), LSHT (4-4), PassageCount (5-1), PassageRetrieval-en (5-2), PassageRetrieval-zh (5-3), LCC (6-1), and RepoBench-P (6-2). (1) With 98% tokens pruned and *no local window or global sink tokens*, Spotlight Attention outperforms Quest and MagicPIG. (2) Spotlight Attention achieves performance on par with the original model, even in tasks like summarization and few-shot learning. (3) On subsets of Chinese (1-4, 2-4, 3-4), other LLaMA2-7B-Chat-based models generated answers in English, while Quest produced responses in Chinese, achieving overwhelmingly high scores.

Method	Configuration	Single-Doc.			Multi-Doc.			Summarization			Few-Shot			Synthetic			Code											
		#1-1	#1-2	#1-3	#1-4	Avg.	#2-1	#2-2	#2-3	#2-4	Avg.	#3-1	#3-2	#3-3	#3-4	Avg.	#4-1	#4-2	#4-3	#4-4	Avg.	#5-1	#5-2	#5-3	Avg.	#6-1	#6-2	Avg.
<b>LLaMA2-7B</b>																												
+Quest	1024 Token Budget	8.73	7.18	15.42	13.84	11.29	6.55	8.27	2.91	11.38	7.27	15.88	19.80	6.03	9.30	12.54	68.00	30.62	30.83	18.25	36.92	1.26	6.97	8.00	5.41	63.66	56.63	60.14
+Quest	128 Token Budget	8.78	9.78	17.95	14.86	12.84	8.91	8.51	2.95	12.53	8.22	18.38	20.84	9.63	8.40	14.31	66.00	67.06	30.35	17.50	42.22	1.69	6.47	7.38	5.18	63.88	58.82	61.35
+MagicPIG	128 Token Budget	8.78	9.78	17.95	14.86	12.84	8.91	8.51	2.95	12.53	8.22	18.38	20.84	9.63	8.40	14.31	66.00	67.06	30.35	17.50	42.22	1.69	6.47	7.38	5.18	63.88	58.82	61.35
+Spotlight	Default	11.85	7.20	20.01	14.23	13.32	8.30	8.23	4.76	12.39	8.42	16.13	20.71	2.21	8.34	11.84	66.50	88.48	35.04	19.50	52.38	0.99	7.78	8.52	5.76	64.95	58.63	61.79
+Spotlight	90% Pruned (< 409)	10.39	7.18	15.86	14.21	11.90	6.58	8.44	3.09	11.61	7.43	16.32	19.31	10.24	8.51	13.59	67.50	38.08	30.80	18.50	38.67	2.07	8.27	6.79	5.71	64.10	56.76	60.43
+Spotlight	98% Pruned (< 81)	18.06	25.36	30.99	27.90	20.58	30.95	12.61	12.51	3.79	18.34	27.58	19.44	17.25	20.22	26.36	0.17	18.53	63.50	25.72	49.29	2.35	3.75	4.25	3.08	59.07	53.44	55.75
<b>LLaMA2-7B-Chat</b>																												
+Quest	1024 Token Budget	19.14	17.78	27.12	21.53	35.99	26.67	12.96	12.33	21.98	27.21	20.62	26.21	14.17	22.05	62.50	78.24	40.53	15.25	49.13	2.00	11.00	6.84	6.61	53.71	50.46	52.08	
+Quest	128 Token Budget	14.11	12.78	17.59	13.47	28.38	21.03	8.99	8.53	16.73	11.55	18.44	20.91	3.99	14.82	38.50	66.44	31.69	10.50	36.78	0.76	7.58	3.30	3.70	43.21	39.55	41.38	
+MagicPIG	128 Token Budget	14.11	12.78	17.59	13.47	28.38	21.03	8.99	8.53	16.73	11.55	18.44	20.91	3.99	14.82	38.50	66.44	31.69	10.50	36.78	0.76	7.58	3.30	3.70	43.21	39.55	41.38	
+Spotlight	90% Pruned (< 409)	18.43	24.76	32.33	7.76	20.81	37.16	29.88	12.91	2.80	19.44	27.55	20.22	26.16	0.17	18.53	63.50	25.72	49.29	2.35	3.75	4.25	3.08	59.07	53.44	55.75		
+Spotlight	98% Pruned (< 81)	18.06	25.36	30.99	7.90	20.58	30.95	12.61	12.51	3.79	18.34	27.24	20.58	26.36	0.32	18.63	59.70	41.27	16.50	47.89	2.97	5.50	5.75	4.74	57.72	52.62	55.17	
<b>LLaMA3-8B</b>																												
+Quest	1024 Token Budget	5.47	13.29	21.03	23.81	10.65	6.61	12.40	9.84	28.80	23.01	3.78	5.56	14.79	71.00	28.48	35.75	36.50	42.83	2.00	6.72	27.61	12.11	49.69	48.18	48.93		
+Quest	256 Token Budget	5.40	12.80	20.03	14.10	9.13	11.79	13.88	10.30	20.41	17.66	12.24	4.51	11.23	47.50	55.10	35.75	26.75	42.47	2.05	6.82	10.61	8.89	53.13	53.07	59.60		
+MagicPIG	256 Token Budget	3.90	13.53	17.51	18.71	13.41	9.16	6.03	13.70	10.12	23.58	23.90	1.37	4.80	13.41	71.50	90.37	44.02	33.50	59.84	1.20	7.15	13.87	7.40	53.36	47.62	50.49	
+Spotlight	90% Pruned (< 819)	4.87	13.63	21.54	20.91	15.23	9.16	11.83	6.54	11.89	9.80	27.09	23.35	3.36	3.37	14.29	70.50	36.89	33.46	45.17	2.03	6.16	24.51	10.90	53.36	47.62	50.49	
+Spotlight	98% Pruned (< 163)	4.59	14.14	20.48	21.64	15.21	9.82	11.97	6.31	11.89	9.99	28.97	23.28	4.37	3.58	14.20	70.50	31.33	34.06	33.50	46.84	2.07	6.59	27.11	9.42	52.26	47.83	48.54

## 485 B.6 Ablation on Attention Estimation Methods

486 We evaluated two approaches: hashing with Hamming distance (our default choice) and down-  
487 projection with inner product. At a 16x compression rate, results in Table 11 show that hashing  
488 outperforms down-projection, demonstrating greater efficiency when the dimensionality is low.

489

**Table 10: Ablation on training loss.** Compared to attention reconstruction loss, our proposed ranking loss yields significantly improved training outcomes.

Loss	PG19	Math	Code
Attn. Recon. Loss	21.341	8.856	11.573
Ranking Loss (ours)	<b>8.977</b>	<b>3.621</b>	<b>5.434</b>

**Table 11: Ablation on estimation method.** We compared two attention estimation methods: hashing with Hamming distance (our default choice) and down-projection with inner product. Results indicate that hashing with Hamming distance performs better.

Estimation Method	Training	PG19	Math	Code
Down Proj. ( $16 \times$ ) + Inner Prod.	✓	14.743	5.425	7.619
Hashing + Hamming Dist. (ours)		<b>8.977</b>	<b>3.621</b>	<b>5.434</b>

490 **NeurIPS Paper Checklist**

491 The checklist is designed to encourage best practices for responsible machine learning research,  
492 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove  
493 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should  
494 follow the references and follow the (optional) supplemental material. The checklist does NOT count  
495 towards the page limit.

496 Please read the checklist guidelines carefully for information on how to answer these questions. For  
497 each question in the checklist:

- 498 • You should answer [Yes] , [No] , or [NA] .  
499 • [NA] means either that the question is Not Applicable for that particular paper or the  
500 relevant information is Not Available.  
501 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

502 **The checklist answers are an integral part of your paper submission.** They are visible to the  
503 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it  
504 (after eventual revisions) with the final version of your paper, and its final version will be published  
505 with the paper.

506 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.  
507 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a  
508 proper justification is given (e.g., "error bars are not reported because it would be too computationally  
509 expensive" or "we were unable to find the license for the dataset we used"). In general, answering  
510 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we  
511 acknowledge that the true answer is often more nuanced, so please just use your best judgment and  
512 write a justification to elaborate. All supporting evidence can appear either in the main paper or the  
513 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification  
514 please point to the section(s) where related material for the question can be found.

515 **IMPORTANT**, please:

- 516 • **Delete this instruction block, but keep the section heading “NeurIPS Paper Checklist”,**  
517 • **Keep the checklist subsection headings, questions/answers and guidelines below.**  
518 • **Do not modify the questions and only use the provided macros for your answers.**

519 **1. Claims**

520 Question: Do the main claims made in the abstract and introduction accurately reflect the  
521 paper's contributions and scope?

522 Answer: [Yes]

523 Justification: The experimental section validates all ideas introduced in the Introduction,  
524 confirming the accuracy of all stated contributions.

525 Guidelines:

- 526 • The answer NA means that the abstract and introduction do not include the claims  
527 made in the paper.  
528 • The abstract and/or introduction should clearly state the claims made, including the  
529 contributions made in the paper and important assumptions and limitations. A No or  
530 NA answer to this question will not be perceived well by the reviewers.  
531 • The claims made should match theoretical and experimental results, and reflect how  
532 much the results can be expected to generalize to other settings.  
533 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
534 are not attained by the paper.

535 **2. Limitations**

536 Question: Does the paper discuss the limitations of the work performed by the authors?

537 Answer: [Yes]

538 Justification: In the Conclusion and Limitations section, we outline the limitations of this  
539 study.

540 Guidelines:

- 541 • The answer NA means that the paper has no limitation while the answer No means that  
542 the paper has limitations, but those are not discussed in the paper.
- 543 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 544 • The paper should point out any strong assumptions and how robust the results are to  
545 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
546 model well-specification, asymptotic approximations only holding locally). The authors  
547 should reflect on how these assumptions might be violated in practice and what the  
548 implications would be.
- 549 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
550 only tested on a few datasets or with a few runs. In general, empirical results often  
551 depend on implicit assumptions, which should be articulated.
- 552 • The authors should reflect on the factors that influence the performance of the approach.  
553 For example, a facial recognition algorithm may perform poorly when image resolution  
554 is low or images are taken in low lighting. Or a speech-to-text system might not be  
555 used reliably to provide closed captions for online lectures because it fails to handle  
556 technical jargon.
- 557 • The authors should discuss the computational efficiency of the proposed algorithms  
558 and how they scale with dataset size.
- 559 • If applicable, the authors should discuss possible limitations of their approach to  
560 address problems of privacy and fairness.
- 561 • While the authors might fear that complete honesty about limitations might be used by  
562 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
563 limitations that aren't acknowledged in the paper. The authors should use their best  
564 judgment and recognize that individual actions in favor of transparency play an impor-  
565 tant role in developing norms that preserve the integrity of the community. Reviewers  
566 will be specifically instructed to not penalize honesty concerning limitations.

### 567 3. Theory assumptions and proofs

568 Question: For each theoretical result, does the paper provide the full set of assumptions and  
569 a complete (and correct) proof?

570 Answer: [NA]

571 Justification: Our experiments do not contain theoretical results, only some empirical  
572 conclusions.

573 Guidelines:

- 574 • The answer NA means that the paper does not include theoretical results.
- 575 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
576 referenced.
- 577 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 578 • The proofs can either appear in the main paper or the supplemental material, but if  
579 they appear in the supplemental material, the authors are encouraged to provide a short  
580 proof sketch to provide intuition.
- 581 • Inversely, any informal proof provided in the core of the paper should be complemented  
582 by formal proofs provided in appendix or supplemental material.
- 583 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 584 4. Experimental result reproducibility

585 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
586 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
587 of the paper (regardless of whether the code and data are provided or not)?

588 Answer: [Yes]

589 Justification: The Experimental Setup section provides detailed descriptions of the experi-  
590 ments, with comprehensive details included in the Appendix A.

591 Guidelines:

- 592 • The answer NA means that the paper does not include experiments.
- 593 • If the paper includes experiments, a No answer to this question will not be perceived
- 594 well by the reviewers: Making the paper reproducible is important, regardless of
- 595 whether the code and data are provided or not.
- 596 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 597 to make their results reproducible or verifiable.
- 598 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 599 For example, if the contribution is a novel architecture, describing the architecture fully
- 600 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 601 be necessary to either make it possible for others to replicate the model with the same
- 602 dataset, or provide access to the model. In general, releasing code and data is often
- 603 one good way to accomplish this, but reproducibility can also be provided via detailed
- 604 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 605 of a large language model), releasing of a model checkpoint, or other means that are
- 606 appropriate to the research performed.
- 607 • While NeurIPS does not require releasing code, the conference does require all submissions
- 608 to provide some reasonable avenue for reproducibility, which may depend on the
- 609 nature of the contribution. For example
  - 610 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
  - 611 to reproduce that algorithm.
  - 612 (b) If the contribution is primarily a new model architecture, the paper should describe
  - 613 the architecture clearly and fully.
  - 614 (c) If the contribution is a new model (e.g., a large language model), then there should
  - 615 either be a way to access this model for reproducing the results or a way to reproduce
  - 616 the model (e.g., with an open-source dataset or instructions for how to construct
  - 617 the dataset).
  - 618 (d) We recognize that reproducibility may be tricky in some cases, in which case
  - 619 authors are welcome to describe the particular way they provide for reproducibility.
  - 620 In the case of closed-source models, it may be that access to the model is limited in
  - 621 some way (e.g., to registered users), but it should be possible for other researchers
  - 622 to have some path to reproducing or verifying the results.

623 **5. Open access to data and code**

624 Question: Does the paper provide open access to the data and code, with sufficient instruc-

625 tions to faithfully reproduce the main experimental results, as described in supplemental

626 material?

627 Answer: [Yes]

628 Justification: We provide open access to data and code in the Abstract.

629 Guidelines:

- 630 • The answer NA means that paper does not include experiments requiring code.
- 631 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 632 • While we encourage the release of code and data, we understand that this might not be
- 633 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
- 634 including code, unless this is central to the contribution (e.g., for a new open-source
- 635 benchmark).
- 636 • The instructions should contain the exact command and environment needed to run to
- 637 reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 638 • The authors should provide instructions on data access and preparation, including how
- 639 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 640 • The authors should provide scripts to reproduce all experimental results for the new
- 641 proposed method and baselines. If only a subset of experiments are reproducible, they
- 642 should state which ones are omitted from the script and why.

- 645           • At submission time, to preserve anonymity, the authors should release anonymized  
646            versions (if applicable).  
647           • Providing as much information as possible in supplemental material (appended to the  
648            paper) is recommended, but including URLs to data and code is permitted.

649       **6. Experimental setting/details**

650       Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
651           parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
652           results?

653       Answer: [Yes]

654       Justification: All experimental details are provided in the Experiment Setup chapter, with  
655           implementation details included in the Appendix A.

656       Guidelines:

- 657           • The answer NA means that the paper does not include experiments.  
658           • The experimental setting should be presented in the core of the paper to a level of detail  
659            that is necessary to appreciate the results and make sense of them.  
660           • The full details can be provided either with the code, in appendix, or as supplemental  
661            material.

662       **7. Experiment statistical significance**

663       Question: Does the paper report error bars suitably and correctly defined or other appropriate  
664           information about the statistical significance of the experiments?

665       Answer: [Yes]

666       Justification: For experiments involving randomness (Sec 4.2), such as the efficiency tests,  
667           we conducted multiple sampling iterations and reported the mean values.

668       Guidelines:

- 669           • The answer NA means that the paper does not include experiments.  
670           • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
671            dence intervals, or statistical significance tests, at least for the experiments that support  
672            the main claims of the paper.  
673           • The factors of variability that the error bars are capturing should be clearly stated (for  
674            example, train/test split, initialization, random drawing of some parameter, or overall  
675            run with given experimental conditions).  
676           • The method for calculating the error bars should be explained (closed form formula,  
677            call to a library function, bootstrap, etc.)  
678           • The assumptions made should be given (e.g., Normally distributed errors).  
679           • It should be clear whether the error bar is the standard deviation or the standard error  
680            of the mean.  
681           • It is OK to report 1-sigma error bars, but one should state it. The authors should  
682            preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
683            of Normality of errors is not verified.  
684           • For asymmetric distributions, the authors should be careful not to show in tables or  
685            figures symmetric error bars that would yield results that are out of range (e.g. negative  
686            error rates).  
687           • If error bars are reported in tables or plots, The authors should explain in the text how  
688            they were calculated and reference the corresponding figures or tables in the text.

689       **8. Experiments compute resources**

690       Question: For each experiment, does the paper provide sufficient information on the com-  
691           puter resources (type of compute workers, memory, time of execution) needed to reproduce  
692           the experiments?

693       Answer: [Yes]

694       Justification: We give the relevant information in Sec 4.2.

695       Guidelines:

- 696           • The answer NA means that the paper does not include experiments.  
697           • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
698            or cloud provider, including relevant memory and storage.  
699           • The paper should provide the amount of compute required for each of the individual  
700            experimental runs as well as estimate the total compute.  
701           • The paper should disclose whether the full research project required more compute  
702            than the experiments reported in the paper (e.g., preliminary or failed experiments that  
703            didn't make it into the paper).

704           **9. Code of ethics**

705           Question: Does the research conducted in the paper conform, in every respect, with the  
706           NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

707           Answer: [Yes]

708           Justification: Our research conducted in the paper conform, in every respect, with the  
709           NeurIPS Code of Ethics.

710           Guidelines:

- 711           • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.  
712           • If the authors answer No, they should explain the special circumstances that require a  
713            deviation from the Code of Ethics.  
714           • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
715            eration due to laws or regulations in their jurisdiction).

716           **10. Broader impacts**

717           Question: Does the paper discuss both potential positive societal impacts and negative  
718           societal impacts of the work performed?

719           Answer: [NA]

720           Justification: Our approach does not involve societal impact.

721           Guidelines:

- 722           • The answer NA means that there is no societal impact of the work performed.  
723           • If the authors answer NA or No, they should explain why their work has no societal  
724            impact or why the paper does not address societal impact.  
725           • Examples of negative societal impacts include potential malicious or unintended uses  
726            (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
727            (e.g., deployment of technologies that could make decisions that unfairly impact specific  
728            groups), privacy considerations, and security considerations.  
729           • The conference expects that many papers will be foundational research and not tied  
730            to particular applications, let alone deployments. However, if there is a direct path to  
731            any negative applications, the authors should point it out. For example, it is legitimate  
732            to point out that an improvement in the quality of generative models could be used to  
733            generate deepfakes for disinformation. On the other hand, it is not needed to point out  
734            that a generic algorithm for optimizing neural networks could enable people to train  
735            models that generate Deepfakes faster.  
736           • The authors should consider possible harms that could arise when the technology is  
737            being used as intended and functioning correctly, harms that could arise when the  
738            technology is being used as intended but gives incorrect results, and harms following  
739            from (intentional or unintentional) misuse of the technology.  
740           • If there are negative societal impacts, the authors could also discuss possible mitigation  
741            strategies (e.g., gated release of models, providing defenses in addition to attacks,  
742            mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
743            feedback over time, improving the efficiency and accessibility of ML).

744           **11. Safeguards**

745           Question: Does the paper describe safeguards that have been put in place for responsible  
746           release of data or models that have a high risk for misuse (e.g., pretrained language models,  
747           image generators, or scraped datasets)?

748                  Answer: [NA]

749                  Justification: Our work does not include such risks.

750                  Guidelines:

- 751                  • The answer NA means that the paper poses no such risks.
- 752                  • Released models that have a high risk for misuse or dual-use should be released with  
753                  necessary safeguards to allow for controlled use of the model, for example by requiring  
754                  that users adhere to usage guidelines or restrictions to access the model or implementing  
755                  safety filters.
- 756                  • Datasets that have been scraped from the Internet could pose safety risks. The authors  
757                  should describe how they avoided releasing unsafe images.
- 758                  • We recognize that providing effective safeguards is challenging, and many papers do  
759                  not require this, but we encourage authors to take this into account and make a best  
760                  faith effort.

## 761                  12. Licenses for existing assets

762                  Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
763                  the paper, properly credited and are the license and terms of use explicitly mentioned and  
764                  properly respected?

765                  Answer: [Yes]

766                  Justification: All assets used are well-known, properly referenced, and their information is  
767                  readily accessible.

768                  Guidelines:

- 769                  • The answer NA means that the paper does not use existing assets.
- 770                  • The authors should cite the original paper that produced the code package or dataset.
- 771                  • The authors should state which version of the asset is used and, if possible, include a  
772                  URL.
- 773                  • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 774                  • For scraped data from a particular source (e.g., website), the copyright and terms of  
775                  service of that source should be provided.
- 776                  • If assets are released, the license, copyright information, and terms of use in the  
777                  package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)  
778                  has curated licenses for some datasets. Their licensing guide can help determine the  
779                  license of a dataset.
- 780                  • For existing datasets that are re-packaged, both the original license and the license of  
781                  the derived asset (if it has changed) should be provided.
- 782                  • If this information is not available online, the authors are encouraged to reach out to  
783                  the asset's creators.

## 784                  13. New assets

785                  Question: Are new assets introduced in the paper well documented and is the documentation  
786                  provided alongside the assets?

787                  Answer: [NA]

788                  Justification: Our paper does not release new assets.

789                  Guidelines:

- 790                  • The answer NA means that the paper does not release new assets.
- 791                  • Researchers should communicate the details of the dataset/code/model as part of their  
792                  submissions via structured templates. This includes details about training, license,  
793                  limitations, etc.
- 794                  • The paper should discuss whether and how consent was obtained from people whose  
795                  asset is used.
- 796                  • At submission time, remember to anonymize your assets (if applicable). You can either  
797                  create an anonymized URL or include an anonymized zip file.

## 798                  14. Crowdsourcing and research with human subjects

799 Question: For crowdsourcing experiments and research with human subjects, does the paper  
800 include the full text of instructions given to participants and screenshots, if applicable, as  
801 well as details about compensation (if any)?

802 Answer: [NA]

803 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

804 Guidelines:

- 805 • The answer NA means that the paper does not involve crowdsourcing nor research with  
806 human subjects.
- 807 • Including this information in the supplemental material is fine, but if the main contribu-  
808 tion of the paper involves human subjects, then as much detail as possible should be  
809 included in the main paper.
- 810 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,  
811 or other labor should be paid at least the minimum wage in the country of the data  
812 collector.

813 **15. Institutional review board (IRB) approvals or equivalent for research with human  
814 subjects**

815 Question: Does the paper describe potential risks incurred by study participants, whether  
816 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
817 approvals (or an equivalent approval/review based on the requirements of your country or  
818 institution) were obtained?

819 Answer: [NA]

820 Justification: Our paper does not involve crowdsourcing nor research with human subjects.

821 Guidelines:

- 822 • The answer NA means that the paper does not involve crowdsourcing nor research with  
823 human subjects.
- 824 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
825 may be required for any human subjects research. If you obtained IRB approval, you  
826 should clearly state this in the paper.
- 827 • We recognize that the procedures for this may vary significantly between institutions  
828 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
829 guidelines for their institution.
- 830 • For initial submissions, do not include any information that would break anonymity (if  
831 applicable), such as the institution conducting the review.

832 **16. Declaration of LLM usage**

833 Question: Does the paper describe the usage of LLMs if it is an important, original, or  
834 non-standard component of the core methods in this research? Note that if the LLM is used  
835 only for writing, editing, or formatting purposes and does not impact the core methodology,  
836 scientific rigorousness, or originality of the research, declaration is not required.

837 Answer: [NA]

838 Justification: Our study focuses on open-source LLMs, but LLMs are not involved in the  
839 research methodology.

840 Guidelines:

- 841 • The answer NA means that the core method development in this research does not  
842 involve LLMs as any important, original, or non-standard components.
- 843 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
844 for what should or should not be described.