

BREADCRUMBS REASONING: MEMORY-EFFICIENT REASONING WITH COMPRESSION BEACONS

Anonymous authors

Paper under double-blind review

ABSTRACT

The scalability of large language models for long-context reasoning is severely constrained by the linear growth of their Transformer key-value cache, which incurs significant memory and computational costs. We posit that as a model generates reasoning tokens, the informational value of past generated tokens diminishes, creating an opportunity for compression. In this work, we propose to periodically compress the generation KV cache with a learned, special-purpose token and evict compressed entries. We train the model to perform this compression via a modified joint distillation and reinforcement learning (RL) framework. Our training method minimizes overhead over the conventional RL process, as it leverages RL outputs for distillation. Empirically, our method achieves a superior memory-accuracy Pareto frontier compared to both the model without cache compression and training-free compression techniques.

1 INTRODUCTION

Reasoning through token generation allows large language models (LLMs) to solve arbitrarily complex problems with a fixed depth architecture (Merrill & Sabharwal, 2023), by scaling the compute invested through the generation of more tokens (i.e., test-time scaling) (Snell et al., 2024). This practice carries high computational costs, because of the self-attention design of Transformers (Bahdanau et al., 2014; Vaswani et al., 2017; Keles et al., 2023). Not only it relies on simply generating many more tokens, but later tokens require computation over the representations of all previous tokens, incurring higher time complexity and necessitating increasing memory costs.

However, not all past representations are equally important. For example, the details of a previously explored attempt at a solution are likely not critical, as long as the model retains some signal that

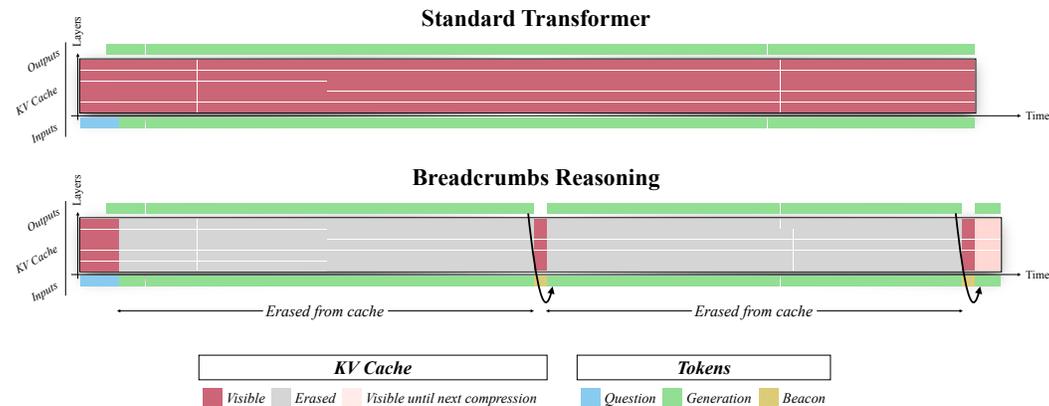


Figure 1: **Breadcrumbs Reasoning**, with a compression ratio $c = 16$. To save memory during inference, a window of c tokens is periodically compressed into a single beacon token. The original KV cache entries for the window are then evicted, leaving only a compact 'breadcrumb' that summarizes the preceding reasoning steps.

054 advises it to avoid exploring this same failed path again. We propose to jointly learn to reason,
055 compress, and discard previously computed representations along reasoning chains.

056
057 The key to our approach is to substitute previously computed key-value (KV) cached representa-
058 tions with significantly more compact representations, inspired by how activation beacons are used
059 for long-context compression (Zhang et al., 2025). We train these representations to contain the
060 information from past tokens that is necessary for continuing the reasoning process to solve the task
061 the model is given, allowing us to evict from the KV cache most previously computed representa-
062 tions. The challenge is to combine the training of these beacons into the reinforcement learning
063 (RL) (Sutton & Barto, 2018) process that makes length-based reasoning possible, a fundamentally
064 different process than the pre-training that enables conventional long-form generation. We design
065 a joint RL-distillation approach, where we train the original non-compression policy using the con-
066 ventional RL process with a verifier for reward computation, and concurrently distill it into a policy
that jointly compresses and reasons.

067 We evaluate our method, Breadcrumbs Reasoning (Figure 1), on the Qwen2.5-1.5B and Phi-4-Mini
068 models across three challenging reasoning benchmarks: Countdown (Gandhi et al., 2024), LinSys,
069 and StarGraph (Bachmann & Nagarajan, 2024). We compare against a strong, uncompressed teacher
070 policy trained with RL, as well as four training-free cache eviction baselines: PyramidKV (Cai
071 et al., 2025), SnapKV (Li et al., 2024), TOVA (Oren et al., 2024), and StreamingLLM (Xiao et al.,
072 2023). Our experiments reveal several key findings. Demonstrating a clear Pareto improvement,
073 Breadcrumbs Reasoning enables effective test-time scaling by generating longer reasoning chains
074 to match or exceed the teacher’s accuracy within a fixed memory budget, while still retaining 65.1–
075 89.8% of the original performance when using 2–32x less memory at a fixed generation length.
076 Notably, at inference time, it is possible to choose the compression ratio for each given input based
077 on the memory budget and required performance. In contrast, the training-free baselines consistently
078 underperform, stressing the necessity of a learned compression scheme for complex reasoning. We
079 also validate our training strategy, showing that our joint RL-distillation approach matches or out-
080 performs a more complex two-stage training pipeline, confirming its efficiency. Our code will be
081 released upon publication.

082 2 RELATED WORK AND BACKGROUND

083
084
085 Generation in Transformers-based LLMs requires reasoning over and storing a key-value (KV)
086 cache. This entails high memory (i.e., space) and time costs, which increase as the context (i.e.,
087 the number of previous tokens) increases. Therefore long-form generation costs suffer not only
088 from the fundamental need to generate more tokens via more steps, but also from the increasing
089 cost of each such step. KV compression is a solution avenue that is receiving significant research
090 attention (Li et al., 2025).

091 An important thread within this compression literature is training models to perform KV cache com-
092 pression. Nawrot et al. (2024) train models to compute importance scores, which are then used to
093 store averaged KV cache entries instead of the original entries. Other methods train the Transformer-
094 based LLMs themselves to summarize past KV entries (Mu et al., 2023; Chevalier et al., 2023; Zhang
095 et al., 2025), so more compact representations can be retained. Our approach is inspired by the ac-
096 tivation beacons method (Zhang et al., 2025), but with significant simplifications and adaptation for
097 reasoning. We do away with the chunk and sub-chunk distinction, and eliminate the addition of spe-
098 cialized attention mechanisms. Rather, we adopt a flat segmentation into blocks, use the standard
099 Transformer attention mechanism, and remove KV cache entries every time a beacon is processed
100 (e.g., ex-ante). These modifications do not only simplify the implementation, but also allow for
101 immediate eviction of cache entries, instead of a delayed one. More broadly, a drawback of these
102 learned methods is that they require fine-tuning on a considerable amount of general-purpose pre-
103 training data. We design a joint reasoning-compression training approach, which adds minimal
overhead over the existing reasoning training processes.

104 An alternative to training-based methods are training-free methods that perform compression at
105 generation time. They can be divided into two main categories. The first is that of sliding-window
106 approaches, which limit the KV cache by only including a sliding window plus an additional subset
107 of tokens. Particularly simple and effective is StreamingLLM (Xiao et al., 2023), which finds that
including a few initial *sink* tokens in addition to the window recovers most of the uncompressed

Algorithm 1 Breadcrumbs Reasoning

Input: Transformer-based policy π_{BR} , beacon token b , prompt tokens \bar{q} , stop token s , compression ratio c . Let $\text{KV}_{\pi_{\text{BR}}}$ be the persistent KV cache of the policy model π_{BR} .

Output: \bar{x}

```

1: Initialize: Encode  $\bar{q}$  through  $\pi$ 
2: for  $i = 0, 1, 2, \dots$  do
3:    $x_i \sim \pi_{\text{BR}}(\cdot | \bar{q}, \bar{x})$   $\triangleright$  Sample the next token.  $\text{KV}_{\pi_{\text{BR}}}$  is updated internally.
4:    $\bar{x} \leftarrow \bar{x} + x_i$   $\triangleright$  Concatenate the sampled token to the end of the output.
5:   if  $x_i = s$  then  $\triangleright$  Check for the generation stopping token.
6:     break
7:   if  $i > 0$  and  $i \bmod c = 0$  then
8:     Encode  $b$  through  $\pi_{\text{BR}}(\cdot | \bar{q}, \bar{x})$   $\triangleright$  Updates  $\text{KV}_{\pi_{\text{BR}}}$  with the entry for the compression token  $b$ .
9:      $\text{KV}_{\pi_{\text{BR}}} \leftarrow \text{KV}_{\pi_{\text{BR}}}[: -c - 1]$   $\triangleright$  Drop the KV cache entries of the most recent  $c$  tokens.
10:     $\text{KV}_{\pi_{\text{BR}}} \leftarrow \text{KV}_{\pi_{\text{BR}}} + \text{KV}_{\pi_{\text{BR}}}[-1]$   $\triangleright$  But, keep the entry of the beacon  $b$ .
11: return  $\bar{x}$ 

```

performance. The second type does not constrain window tokens to remain in the cache, but rather tries to select the empirically more important for attention computations, as in H2O (Zhang et al., 2023) or TOVA (Oren et al., 2024). Other approaches that can be considered part of this category focus on prefill compression, reducing tokens by using the attention scores of the window of most recent token of the entire prefilling prompt (Cai et al., 2025; Li et al., 2024). While these latter methods are designed for prefill and do not compress generation tokens directly, it could be possible to apply them repeatedly to compress at generation time.

An orthogonal approach to reducing the KV cache size is reducing Chain-of-Thought reasoning length (Aggarwal & Welleck, 2025; Kang et al., 2025; Ma et al., 2025; Shen et al., 2025; Yan et al., 2025; Munkhbat et al., 2025; Xia et al., 2025). These methods primarily aim to directly shorten reasoning traces. This is distinct from our objective of dynamic KV cache compression, which focuses on extracting critical information to manage cache sizes effectively during the reasoning process itself. KV cache compression methods, like ours, could be applied in combination with those methods to achieve even higher levels of efficiency.

3 METHODOLOGY

Breadcrumbs Reasoning (BR) periodically computes compressed representations of KV cache entries and evicts them from the KV cache. We design a training process that adds relatively little overhead on top of the conventional reasoning RL process. The learned policy model effectively reasons through token generation and concurrently compresses KV cache representations.

3.1 BREADCRUMBS REASONING

Our compression scheme uses the Transformer architecture itself for both compression and reasoning. We generate tokens following the same procedure as a vanilla Transformer-based language model, but periodically compute compressed representations of past KV cache entries, and evict these entries from the cache. Algorithm 1 describes the process. We add a special token b to the model vocabulary and embedding matrix, to mark when the model should compute compressed representations of past tokens. We input this token every c tokens, where c is the target compression ratio. The KV cache entries for the token b form the compressed representation, and we drop the entries for previous c tokens. Immediately after the beacon b , we force the next input token to be the last sampled token before b was given as input. This is equal to continuing conventional generation. Figure 1 visualizes this process to illustrate the space savings. Roughly speaking, this process leaves a trace of “reasoning breadcrumbs” behind, instead of long, detailed, and eventually irrelevant information.

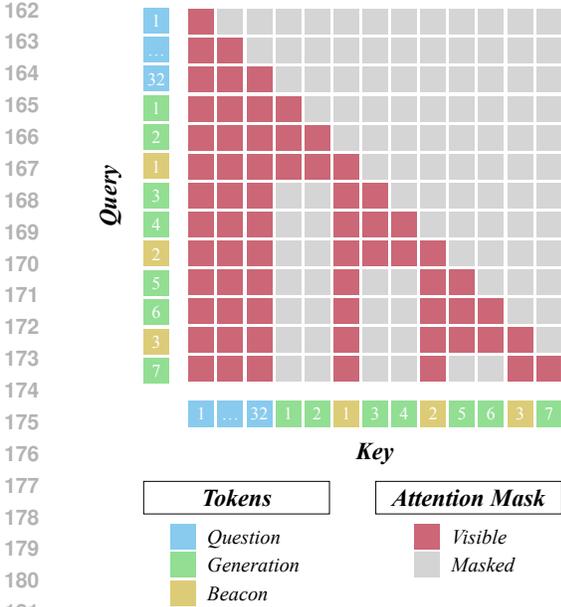


Figure 2: **Attention mask used to enforce compression during training.** Each token after the question can attend to the initial question tokens, all previous beacons, and earlier generation tokens within the same window, i.e., since the most recent preceding beacon. This encourages the model to compress relevant past context into the beacons to support future generations.

3.2 JOINT RL-DISTILLATION TRAINING

Typically, LLMs are trained to solve reasoning tasks through RL (Shao et al., 2024; DeepSeek-AI, 2025; Lambert et al., 2025). Applying this process as is to a breadcrumbs reasoning policy is technically possible, but is unlikely to lead to effective learning. This is because, before training, the model does not have compression ability, so incorporating the compression token and cache eviction will significantly damage its functionality, and it will observe no positive reward during RL.

We use a surrogate teacher policy π_{RL} that does not compress and is trained through RL to perform the reasoning task, and distill it into the breadcrumbs policy π_{BR} . By learning to imitate the surrogate policy π_{RL} , our target breadcrumbs policy π_{BR} simultaneously learns to compress and to perform the new reasoning task. This process relies on the trajectories sampled during RL, so no expensive sampling of trajectories beyond the conventional RL process is needed. This procedure minimizes the overhead over a standard two-steps procedure, where either (a) π_{RL} is completely trained before generating new data for π_{BR} to learn to compress, or (b) π_{BR} is trained to compress through extensive general data and only after that learns the new reasoning task.

Given a token trajectory \bar{x} sampled from π_{RL} , the distillation loss is defined as the average token-level KL divergence between the uncompressed policy π_{RL} and the compressed policy π_{BR} :

$$\begin{aligned}
 L(\bar{x}) &= \frac{1}{|\bar{x}|} \sum_{i=1}^{|\bar{x}|} D_{KL}(\pi_{RL}(x_i|\bar{x}_{<i}) \parallel \pi_{BR}(x_i|\bar{x}_{<i})) \\
 &= \frac{1}{|\bar{x}|} \sum_{i=1}^{|\bar{x}|} \sum_{k=1}^{|\bar{x}|} \pi_{RL}(x_i = k|\bar{x}_{<i}) \log \frac{\pi_{RL}(x_i = k|\bar{x}_{<i})}{\pi_{BR}(x_i = k|\bar{x}_{<i})}
 \end{aligned}
 \tag{1}$$

The total loss is then computed as the average over a batch of B trajectories:

$$L = \mathbb{E}_{\bar{x} \sim \pi_{RL}} [L(\bar{x})] \approx \frac{1}{B} \sum_{b=1}^B L(\bar{x}_b)$$

When training π_{BR} , we only compute the gradient of L with respect to the parameters of π_{BR} . The parameters of the teacher policy π_{RL} are kept frozen during this update step, as π_{RL} should not adapt to π_{BR} .

For parallel and efficient training, at training time π_{BR} does not execute the KV cache removal strategy (Section 3.1), but simulates it by masking compressed tokens. Figure 2 visualizes the attention mask that results from the training masking pattern. Therefore, π_{BR} learns to use the

216 beacon token b and to compress by aligning its next token distribution to the next token distribution
 217 of a policy without compression, π_{RL} . Thanks to the attention mask, the model at training time
 218 needs to learn how to retain useful information through the beacons’ activations, as future tokens
 219 could not otherwise use it.

221 4 EXPERIMENTAL SETUP

222 **Tasks** We use three reasoning tasks that the initial models solve with only a very low success rate:

223
 224 *Countdown* (Gandhi et al., 2024): the task input is a tuple of numbers, and the goal is to create a
 225 sequence of arithmetic operations using a subset of the numbers to equal a target number. While
 226 this task requires the model to avoid repeating previous attempts, or it can enter an endless loop,
 227 the individual guesses are largely independent of one another.
 228

229
 230 *LinSys*: each problem consists of a linear equation system with a unique integer solution. The
 231 coefficients and variables are randomly generated. In contrast to Countdown, which emphasizes
 232 educated trial and error, this task more closely reflects structured reasoning: solving for one
 233 variable often enables solving for the next, resulting in a multi-step deductive process.

234
 235 *StarGraph* (Bachmann & Nagarajan, 2024): the model is given a list of directed edges of a star
 236 graph (i.e., a graph with multiple branches of a fixed length all expanding from a central node) and
 237 a target end node. The model must output the full sequence of edges to get to the target node from
 238 the central node. This is a task auto-regressive models such as Transformers naturally struggle
 239 with, as observed by Bachmann & Nagarajan (2024) and Hu et al. (2025).

240 **Model and Training Details** We utilize two models for our experiments: Qwen2.5-1.5B-
 241 Instruct (Team, 2025) and Phi-4-mini-instruct (Microsoft, 2025). In BR, we add to an additional
 242 token, the beacon b , to the vocabulary of these models and initialize its embeddings with the average
 243 of all other embeddings. For Qwen2.5-1.5B-Instruct, we train for 1000 steps for all tasks, while
 244 for Phi-4-mini-instruct, we train for 200 steps for StarGraph and LinSys and 1000 steps for Count-
 245 down. We use a batch size of 256 for both models. We use PPO (Schulman et al., 2017) as the RL
 246 algorithm for π_{RL} . We experiment with compression ratios c of 2, 4, 8, 16, and 32 for breadcrumbs
 247 reasoning. We instantiate our training process in two different ways. In SR BR (Single Ratio Bread-
 248 crumbs Reasoning), each model is trained on a single compression ratio. In MR BR (Multi Ratio
 249 Breadcrumbs Reasoning), we train a single model on all compression ratios, by repeating each batch
 250 for all compression ratios before each model weights update. Then, the same model can be used at
 251 generation time with any desired ratio. The reward for an incorrect or not present answer is 0.0, a
 252 correct format of the final answer but an incorrect value gets a 0.1 reward, and a correct response
 253 gets a 1.0 reward.

254 **Baselines** We compare our approach against two primary training-free baselines applied to
 255 π_{RL} : PyramidKV (Cai et al., 2025), SnapKV (Li et al., 2024), TOVA (Oren et al., 2024), and
 256 StreamingLLM (Xiao et al., 2023). Similar to our joint training setup, they also do not require more
 257 data than what π_{RL} was trained on, and they also delete entries from the KV cache.¹ For a fair
 258 comparison, we adapt these methods to use an increasingly large KV cache size during generation,
 259 matching the memory footprint of our approach. This is to avoid potentially penalizing them with
 260 a smaller static cache size. In particular, given a compression ratio c , we allow the cache to grow
 261 by 1 every c generated tokens. We set the sliding window of StreamingLLM and the observation
 262 window of SnapKV and PyramidKV to c to match the budget given to our approach. We also set
 263 the sink size of StreamingLLM to the entire question. For SnapKV, PyramidKV, and TOVA, we set
 264 the initial size of the KV cache to the number of question tokens plus c , so that the first compression
 265 would only happen after at least c tokens, similar to our method. For all methods, we test the same
 266 c as for our policies. We also study the effect of two-step training, by first training the RL policy
 267 π_{RL} and only later generating data for distillation. In both cases, we use the same number of data
 268 samples and training steps (256 samples per step, for 1k steps).

269 ¹We omit comparing to methods such as QUEST (Tang et al., 2024), because they are not focused on
 compression, but rather smart management of the GPU memory, an orthogonal approach to compression.

		QWEN																	
		COUNTDOWN						LINSYS						STARGRAPH					
COMPR. RATIO		1	2	4	8	16	32	1	2	4	8	16	32	1	2	4	8	16	32
NO COMPRESSION																			
TEACHER		0.598	-	-	-	-	-	0.918	-	-	-	-	-	0.902	-	-	-	-	-
BREADCRUMBS REASONING																			
SR BR (OURS)	-	<u>0.605</u>	<u>0.613</u>	<u>0.613</u>	<u>0.574</u>	<u>0.535</u>	-	0.730	0.656	<u>0.410</u>	<u>0.367</u>	<u>0.297</u>	-	<u>0.957</u>	0.969	<u>0.973</u>	<u>0.957</u>	<u>0.957</u>	
MR BR (OURS)	-	0.613	0.617	0.629	0.605	0.582	-	<u>0.711</u>	<u>0.629</u>	0.473	0.414	0.328	-	0.965	<u>0.961</u>	0.980	0.961	0.965	
TRAINING-FREE COMPRESSION																			
PYRAMIDKV	-	0.445	0.215	0.117	0.078	0.012	-	0.000	0.000	0.000	0.000	0.000	-	0.605	0.504	0.516	0.477	0.539	
SNAPKV	-	0.449	0.219	0.141	0.102	0.082	-	0.004	0.000	0.000	0.000	0.000	-	0.660	0.523	0.508	0.465	0.500	
TOVA	-	0.574	0.289	0.172	0.188	0.207	-	0.000	0.000	0.000	0.000	0.000	-	0.664	0.457	0.445	0.430	0.441	
STREAMINGLLM	-	0.012	0.023	0.027	0.051	0.094	-	0.000	0.000	0.000	0.000	0.000	-	0.055	0.055	0.031	0.047	0.117	

		PHI																	
		COUNTDOWN						LINSYS						STARGRAPH					
COMPR. RATIO		1	2	4	8	16	32	1	2	4	8	16	32	1	2	4	8	16	32
NO COMPRESSION																			
TEACHER		0.633	-	-	-	-	-	0.898	-	-	-	-	-	0.848	-	-	-	-	-
BREADCRUMBS REASONING																			
SR BR (OURS)	-	0.609	0.625	0.613	0.625	<u>0.613</u>	-	<u>0.652</u>	<u>0.539</u>	<u>0.363</u>	<u>0.219</u>	<u>0.195</u>	-	0.812	0.832	0.836	0.812	0.816	
MR BR (OURS)	-	<u>0.586</u>	<u>0.594</u>	0.613	<u>0.586</u>	0.625	-	0.668	0.582	0.461	0.297	0.270	-	<u>0.809</u>	<u>0.805</u>	<u>0.812</u>	<u>0.758</u>	<u>0.766</u>	
TRAINING-FREE COMPRESSION																			
PYRAMIDKV	-	0.484	0.238	<u>0.184</u>	0.020	0.000	-	0.016	0.000	0.000	0.000	0.000	-	0.637	0.402	0.305	0.336	0.367	
SNAPKV	-	0.469	0.230	0.152	0.062	0.012	-	0.008	0.000	0.000	0.000	0.000	-	0.652	0.465	0.344	0.340	0.312	
TOVA	-	0.230	0.066	0.051	0.047	0.098	-	0.000	0.000	0.000	0.000	0.000	-	0.801	0.562	0.457	0.453	0.395	
STREAMINGLLM	-	0.016	0.031	0.109	0.156	0.312	-	0.000	0.000	0.000	0.000	0.000	-	0.180	0.020	0.031	0.098	0.102	

Table 1: Model performance on long-context reasoning tasks for a fixed generation cache size. We report accuracy given a maximum cache size of 1,000 entries. **Bold** indicates the best result; underlined indicates the second best.

Evaluation We generate evaluation answers for 256 held-out test examples for each task. In all settings, we fix the maximum KV cache size to 1,000 entries (i.e., tokens), which is also the maximum response length permitted during training. Generation is interrupted if this limit is reached.

5 RESULTS

We compare the Breadcrumbs Reasoning policy to an uncompressed RL policy and to the baselines in two different configurations. We first compare results with a set maximum cache size (Table 1, Table 6, and Figure 6). We also compare with a fixed maximum number of generation tokens (Table 2), the same used during training of π_{RL} .

Test-Time Scaling with Breadcrumbs Reasoning Table 1 shows accuracy performance with a fixed maximum cache budget of 1,000 steps, while we refer to Table 6 in Appendix E for the accuracy area under the curve (AUAC) with varying the maximum cache size up to 1,000. Figure 6 shows the curves. Across most settings, BR recovers most of the performance of the teacher at a much lower memory cache budget, both in the single ratio and multi ratio settings. Except for LinSys, which remains challenging, all compression ratios outperform the teacher for most of the budgets. This is because the compressed models are able to effectively accommodate more reasoning steps within the same cache budget, allowing for more aggressive test-time compute scaling. On Countdown, both BR models outperform even the teacher at maximum KV cache budget.

We measure the trade-off between accuracy and KV cache consumption by measuring the area under the accuracy-cache size curve (Table 6). This metric integrates accuracy over the full range of cache sizes, and thus captures a model’s overall robustness to a variety of cache constraints. A larger AUAC indicates that a method maintains high accuracy across a wider range of cache budgets, rather than only performing well with high memory usage. We can observe that BR outperforms even the teacher policy by 3.6–92.8% for Qwen and 10.9–219.7% for Phi across the different ratios and tasks. For high compression ratios, the number of generated tokens is much higher than during training (e.g., 32,000 tokens vs 1,000 for 32x compression ratio). This is evidence that compression

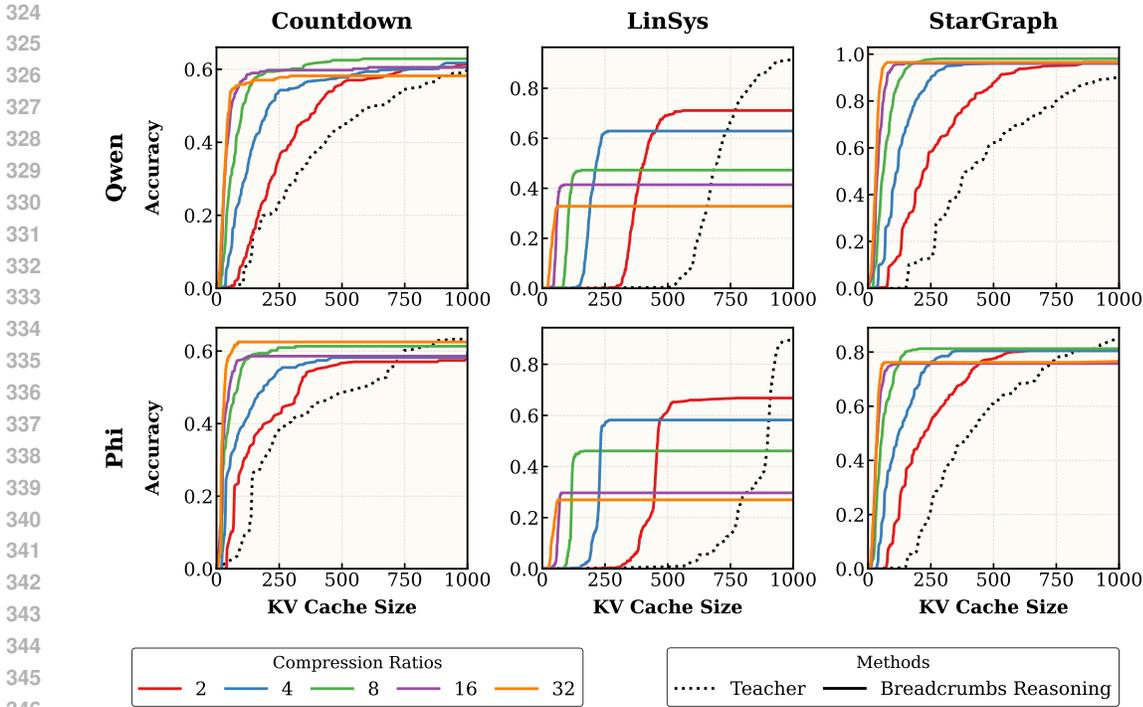


Figure 3: **Accuracy vs. KV Cache Size** Multi Ratio Breadcrumbs Reasoning retains most of the teacher’s performance while using significantly fewer KV cache entries, even outperforming the teacher when setting a fixed maximum KV cache size in Countdown and StarGraph.

generalizes to much longer sequences than those observed during training, except for LinSys (and we refer to Figure 7 in the Appendix to visualize this effect). Crucially, despite potentially generating significantly more tokens (up to $32\times$), BR maintains an inference time comparable to the teacher and on average faster than all training-free baselines, as reported in Table 3.

Beyond matching teacher accuracy at full context, BR is far more efficient with cache usage. In Countdown and StarGraph, Figure 6 shows that BR nearly reaches or even exceeds the teacher’s peak accuracy with fewer than 250 tokens in cache - less than a quarter of the teacher’s required 1,000 tokens. For instance, in Phi-Countdown, BR surpasses 0.60 accuracy by cache size 100, while the teacher requires more than 800 tokens to reach the same level. A similar pattern holds in StarGraph, where BR consistently reaches >0.80 accuracy with only a few hundred tokens, whereas the teacher climbs much more gradually. The only clear exception is LinSys, where BR improves more slowly and is not able to completely close the gap to the teacher.

Breadcrumbs Reasoning Retains Most of Uncompressed Performance We compare the BR policies across different compression ratios to the π_{RL} that we use as the distillation source. We maintain the same source policy for all our experiments to make this comparison valid. Table 2 reports the results for a fixed maximum generation length (i.e., time steps) of 1,000.

While on Countdown and StarGraph all compression ratios work well, performance on LinSys varies greatly across compression ratios. As expected, higher compression ratios lead to worse performance. For LinSys, BR lags behind the teacher for both models. We speculate this is due to differences in the reasoning challenge compared to Countdown and StarGraph. Overall, given a fixed response length, SR BR preserves performance across tasks by 67.1–94.0% for Qwen and 65.1–84.5% for Phi while using only 2–32x fewer KV cache entries at generation time.

Multi-Ratio Training Robustness We analyze the performance of MR BR, where a single model is trained to handle all compression ratios simultaneously. One might expect that learning a policy capable of varying compression granularities would be a harder optimization problem, potentially

QWEN																		
COMPR. RATIO	COUNTDOWN						LINSYS						STARGRAPH					
	1	2	4	8	16	32	1	2	4	8	16	32	1	2	4	8	16	32
NO COMPRESSION																		
TEACHER	0.598	-	-	-	-	-	0.918	-	-	-	-	-	0.902	-	-	-	-	-
BREADCRUMBS REASONING																		
SR BR (OURS)	-	0.559	0.578	<u>0.508</u>	<u>0.430</u>	<u>0.438</u>	-	0.707	0.648	<u>0.395</u>	<u>0.359</u>	<u>0.289</u>	-	0.906	<u>0.883</u>	<u>0.895</u>	<u>0.867</u>	<u>0.891</u>
MR BR (OURS)	-	0.559	<u>0.539</u>	0.551	0.527	0.531	-	<u>0.691</u>	<u>0.625</u>	0.457	0.406	0.328	-	<u>0.875</u>	0.898	0.898	0.875	0.895
TRAINING-FREE COMPRESSION																		
PYRAMIDKV	-	0.438	0.211	0.113	0.078	0.012	-	0.000	0.000	0.000	0.000	0.000	-	0.598	0.504	0.516	0.477	0.539
SNAPKV	-	0.441	0.219	0.141	0.102	0.082	-	0.004	0.000	0.000	0.000	0.000	-	0.660	0.523	0.508	0.465	0.500
TOVA	-	<u>0.535</u>	0.289	0.172	0.188	0.207	-	0.000	0.000	0.000	0.000	0.000	-	0.648	0.457	0.441	0.430	0.441
STREAMINGLLM	-	0.008	0.012	0.012	0.051	0.094	-	0.000	0.000	0.000	0.000	0.000	-	0.043	0.016	0.012	0.012	0.047

PHI																		
COMPR. RATIO	COUNTDOWN						LINSYS						STARGRAPH					
	1	2	4	8	16	32	1	2	4	8	16	32	1	2	4	8	16	32
NO COMPRESSION																		
TEACHER	0.633	-	-	-	-	-	0.895	-	-	-	-	-	0.848	-	-	-	-	-
BREADCRUMBS REASONING																		
SR BR (OURS)	-	0.590	0.574	<u>0.574</u>	0.594	<u>0.570</u>	-	0.641	0.516	<u>0.352</u>	<u>0.207</u>	<u>0.195</u>	-	0.777	0.797	0.793	0.785	0.781
MR BR (OURS)	-	<u>0.566</u>	<u>0.547</u>	0.578	<u>0.555</u>	0.582	-	<u>0.617</u>	0.574	0.434	0.285	0.270	-	<u>0.773</u>	<u>0.762</u>	<u>0.770</u>	<u>0.734</u>	<u>0.734</u>
TRAINING-FREE COMPRESSION																		
PYRAMIDKV	-	0.477	0.234	0.184	0.020	0.000	-	0.016	0.000	0.000	0.000	0.000	-	0.613	0.402	0.305	0.336	0.367
SNAPKV	-	0.469	0.230	0.152	0.062	0.012	-	0.008	0.000	0.000	0.000	0.000	-	0.652	0.465	0.344	0.340	0.312
TOVA	-	0.219	0.066	0.051	0.047	0.098	-	0.000	0.000	0.000	0.000	0.000	-	<u>0.773</u>	0.559	0.453	0.453	0.395
STREAMINGLLM	-	0.016	0.012	0.039	0.137	0.297	-	0.000	0.000	0.000	0.000	0.000	-	0.156	0.016	0.031	0.078	0.102

Table 2: **Model accuracy on long-context reasoning tasks for a fixed generation length.** The metric shown is accuracy at a sequence length of $L = 1,000$. **Bold** indicates the best result; underlined indicates the second best.

leading to lower performance compared to the specialized SR BR models. However, our results indicate the opposite. MR BR not only retains the flexibility to operate at any ratio at inference time but also outperforms SR BR on average. For example, on the Qwen model in Table 1, MR BR achieves an average accuracy of 69.6% across all tasks and compression ratios, compared to 68.1% for SR BR. This suggests that the joint training facilitates information sharing, enabling the model to transfer effective compression strategies across varying ratios.

Training-Free Methods Underperform The training-free methods TOVA and StreamingLLM consistently underperform across tasks. On Countdown, their performance drops dramatically with higher compression (e.g., TOVA falls from 0.574 at 2x to 0.172 at 8x for Qwen; StreamingLLM remains below 0.32 across all settings). On StarGraph, the gap between the baselines and our approach is even more severe, with StreamingLLM dropping below 0.1 accuracy in nearly every configuration. In LinSys, both StreamingLLM and TOVA fail to reach even a single correct output, while SnapKV and PyramidKV are still below 2% accuracy even with the lowest compression ratio. These results highlight the limitation of training-free cache eviction methods: for tasks requiring long coherent reasoning chains, simply truncating past tokens eliminates essential intermediate steps, from which the model is unable to recover.

Why Does BR Struggle with Linear Systems? While BR retains most of the accuracy in Countdown and StarGraph, LinSys seems to plateau earlier, especially for larger compression ratios. The main difference between this task and the other two is that it requires careful arithmetic and algebraic manipulation. However, it could also be that beacons forget more in this task than in the others. We test the two possibilities (BR struggles with arithmetic, beacons fail to memorize) with an LLM-as-a-Judge (Zheng et al., 2023) pipeline (detailed in Appendix D) and report our findings in Figure 4. Interestingly, we find that the vast majority of mistakes are attributed to arithmetic mistakes, and their number increases significantly with increased compression. One potential factor is that compressing numbers might be harder, and also that our compression method might break some arithmetic circuits of the model, and that our training is not long enough to alleviate these issues.

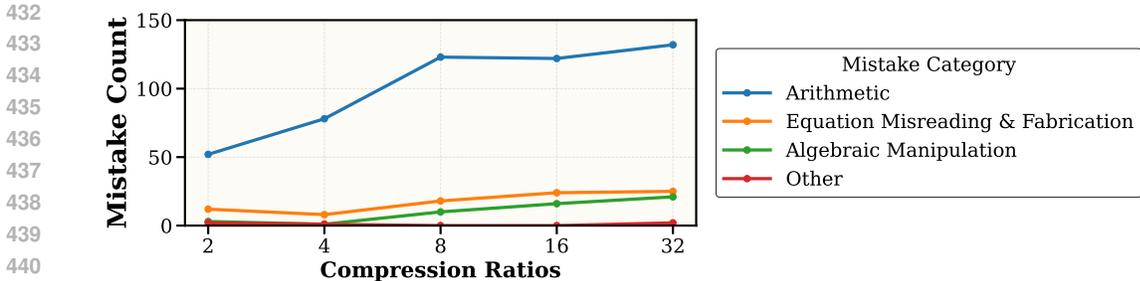


Figure 4: **Failure Mode Analysis on Linear Systems.** We classify error types of Qwen with SR BR across increasing compression ratios. The results indicate that the performance drop in LinSys is primarily driven by arithmetic errors, which scale significantly with compression, rather than a failure to memorize or retrieve information.

METHOD	QWEN					PHI				
	2	4	8	16	32	2	4	8	16	32
BREADCRUMBS REASONING										
SR BR (OURS)	1.990	1.546	1.593	1.646	1.468	2.029	1.486	1.518	1.506	1.164
MR BR (OURS)	1.999	<u>1.665</u>	1.507	1.396	1.366	1.821	<u>1.574</u>	1.307	1.306	<u>1.220</u>
TRAINING-FREE COMPRESSION										
PYRAMIDKV	1.844	2.212	1.858	<u>1.495</u>	1.649	1.823	4.817	10.341	7.669	10.687
SNAPKV	2.607	2.364	1.748	1.664	1.504	<u>1.738</u>	4.347	14.438	10.575	15.234
TOVA	<u>1.880</u>	2.353	2.751	2.563	2.470	1.655	2.184	3.188	3.850	3.890
STREAMINGLLM	2.326	3.210	4.988	5.281	6.349	2.191	3.815	7.354	3.485	3.857

Table 3: **Latency increase across tasks.** The values represent the relative slowdown compared to the Teacher model (lower is better), given a maximum generation KV cache size of 1,000 entries (i.e., up to 32x tokens more with a ratio of 32). Columns represent the compression ratio. **Bold** indicates the best result; underlined indicates the second best.

Joint RL-Distillation Training Matches or Outperforms a Two-Steps Training We compare two strategies for distilling from π_{RL} in Figure 5 (for SR BR). In our joint RL-distillation training, BR is distilled online from the rollouts of the teacher π_{RL} as it learns; in late distillation, compression policies are trained only on trajectories sampled from a final checkpoint of π_{RL} . BR achieves equivalent or superior performance in 26 of the 30 configurations tested, and very close performance on the other four. This demonstrates that it effectively piggybacks on the same data used to train π_{RL} . This eliminates the need for additional distillation data, minimizes training overhead, and does not impose the need to decide a priori a number of training steps for π_{RL} . We hypothesize that the superior performance may derive from the distributional shift of π_{RL} during RL training.

6 DISCUSSION

We present a training-based approach to compress reasoning chains: Breadcrumbs Reasoning. Our empirical results demonstrate that indeed there is significant room for compression in reasoning chains, and not all information or tokens are equally important for downstream reasoning and even task completion. Our approach shows effective compression while retaining much of the reasoning performance. In contrast, training-free methods are significantly less effective.

We propose a joint RL-distillation training scheme, which provides an efficient way to teach a model to reason while also learning to compress. Training is necessary in any case to develop a policy that can successfully solve a reasoning task. Our approach integrates this requirement into a more effective procedure.

Although at a fixed generation length our method shows a small performance loss, we demonstrate that when generating many more tokens under the same memory budget, performance often surpasses that of a non-compressing teacher. This is enabled by effective test-time scaling: when

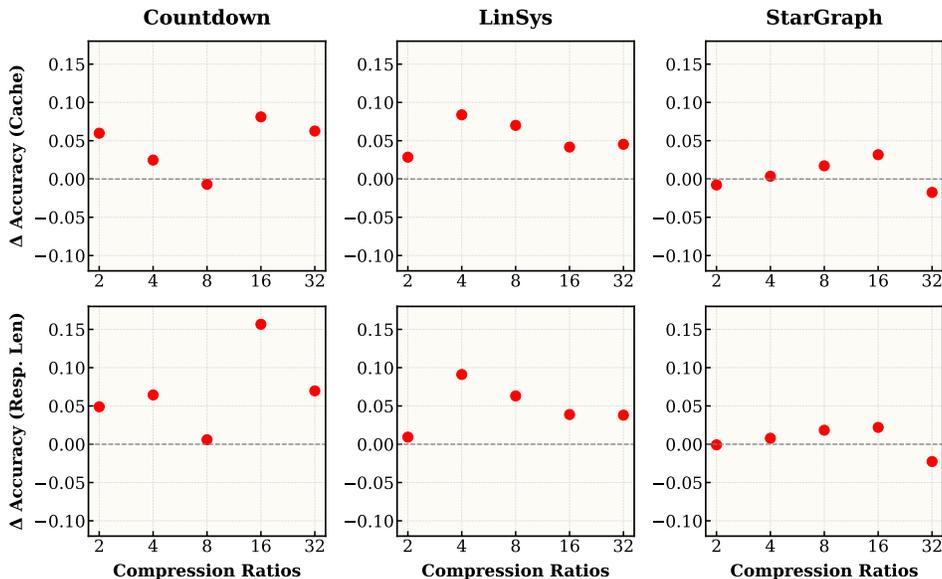


Figure 5: **Joint RL-distillation vs. Two-step Training on Qwen.** We compare our joint approach to the standard two-step process, where compression policies are trained on trajectories from a final π_{RL} checkpoint. Each point shows the accuracy difference (Joint – Two-step) at a given compression ratio. The top row fixes cache size; the bottom row fixes response length. Positive values favor joint training, which outperforms or matches in 26/30 settings, showing BR can learn compression online during RL without separate distillation data.

matching the KV-cache size, we are able to generate significantly more tokens. To some extent, these results show that we trade memory for time. Such trade-offs are common in efficiency-oriented methods — for example, speculative decoding (Leviathan et al., 2023; Chen et al., 2023) reduces latency but requires more memory, since multiple models must be loaded simultaneously.

There are several directions for future work, entailed from several areas where our approach can be improved. Foremost, while Multi Ratio Breadcrumbs Reasoning supports various compression rates, it lacks dynamic adaptivity. It does not automatically select the most appropriate ratio, and once a ratio is chosen, it remains fixed for the duration of the generation rather than adjusting dynamically across the sequence. This is an important direction for future work, and one that is relatively understudied in the compression literature. Our work also charts the direction for future work to improve test-time scaling and compression tradeoffs. Ideally, one can compress aggressively without needing to increase the number of reasoning steps. Our work exposes this tradeoff in reasoning models, and outlines the methodology to analyze it. Moreover, it would be interesting to explore how well our method combines with the orthogonal CoT shortening approach (Aggarwal & Welleck, 2025; Kang et al., 2025; Ma et al., 2025; Shen et al., 2025; Yan et al., 2025; Munkhbat et al., 2025; Xia et al., 2025). Methods in this direction train models to output shorter reasoning traces, but our method can still improve them by saving a significant fraction of memory. Even more, it is important to clarify that these methods do not claim that long reasoning chains are not necessary in general. In fact, the opposite has been theoretically proved by Merrill & Sabharwal (2023): as long as the Transformer architecture is used, long reasoning chains are still required as tasks get more complex. Finally, as the space of accessible benchmarking scenarios in reasoning research develops, it is important to understand the behavior of our approach across a broader set of domains.

REFERENCES

Pranjal Aggarwal and Sean Welleck. L1: Controlling how long a reasoning model thinks with reinforcement learning. In *Second Conference on Language Modeling*, 2025. URL <https://openreview.net/forum?id=4jdIxXBNve>.

- 540 Gregor Bachmann and Vaishnavh Nagarajan. The pitfalls of next-token prediction. In *Proceed-*
541 *ings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of*
542 *Machine Learning Research*, pp. 1547–1574. PMLR, 2024. URL [https://proceedings.](https://proceedings.mlr.press/v235/bachmann24a.html)
543 [mlr.press/v235/bachmann24a.html](https://proceedings.mlr.press/v235/bachmann24a.html).
- 544 Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by
545 jointly learning to align and translate. *CoRR*, abs/1409.0473, 2014. URL [https://api.](https://api.semanticscholar.org/CorpusID:11212020)
546 [semanticscholar.org/CorpusID:11212020](https://api.semanticscholar.org/CorpusID:11212020).
- 547 Zefan Cai, Yichi Zhang, Bofei Gao, Yuliang Liu, Yucheng Li, Tianyu Liu, Keming Lu, Wayne
548 Xiong, Yue Dong, Junjie Hu, and Wen Xiao. PyramidKV: Dynamic KV cache compression
549 based on pyramidal information funneling. In *Second Conference on Language Modeling*, 2025.
550 URL <https://openreview.net/forum?id=ayi7qezU87>.
- 551 Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John
552 Jumper. Accelerating large language model decoding with speculative sampling, 2023.
- 553 Alexis Chevalier, Alexander Wettig, Anirudh Ajith, and Danqi Chen. Adapting language models
554 to compress contexts. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of*
555 *the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 3829–3846,
556 Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
557 emnlp-main.232. URL <https://aclanthology.org/2023.emnlp-main.232>.
- 558 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning,
559 2025. URL <https://arxiv.org/abs/2501.12948>.
- 560 Alessio Devoto, Maximilian Jeblick, and Simon Jégou. Expected attention: Kv cache compression
561 by estimating attention from future queries distribution. *arXiv preprint arXiv:2510.00636*, 2025.
562 URL <https://arxiv.org/abs/2510.00636>.
- 563 Kanishk Gandhi, Denise Lee, Gabriel Grand, Muxin Liu, Winson Cheng, Archit Sharma, and
564 Noah D. Goodman. Stream of search (sos): Learning to search in language, 2024. URL
565 <https://arxiv.org/abs/2404.03683>.
- 566 Edward S. Hu, Kwangjun Ahn, Qinghua Liu, Haoran Xu, Manan Tomar, Ada Langford, Dinesh
567 Jayaraman, Alex Lamb, and John Langford. The belief state transformer. In *The Thirteenth*
568 *International Conference on Learning Representations*, 2025. URL [https://openreview.](https://openreview.net/forum?id=ThRMTCgpvo)
569 [net/forum?id=ThRMTCgpvo](https://openreview.net/forum?id=ThRMTCgpvo).
- 570 Yu Kang, Xianghui Sun, Liangyu Chen, and Wei Zou. C3ot: Generating shorter chain-of-thought
571 without compromising effectiveness. In *Proceedings of the AAAI Conference on Artificial Intelli-*
572 *gence*, volume 39, pp. 24312–24320, 2025.
- 573 Feyza Duman Keles, Pruthvi Mahesakya Wijewardena, and Chinmay Hegde. On the computational
574 complexity of self-attention. In *Proceedings of the 34th International Conference on Algorithmic*
575 *Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pp. 1–23. PMLR,
576 2023. URL <https://proceedings.mlr.press/v201/duman-keles23a.html>.
- 577 Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brah-
578 man, Lester James Validad Miranda, Alisa Liu, Nouha Dziri, Xinxin Lyu, Yuling Gu, Saumya
579 Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Christo-
580 pher Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Ha-
581 jishirzi. Tulu 3: Pushing frontiers in open language model post-training. In *Second Conference on*
582 *Language Modeling*, 2025. URL <https://openreview.net/forum?id=iluGbfHHpH>.
- 583 Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative
584 decoding. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202
585 of *Proceedings of Machine Learning Research*, pp. 19274–19298. PMLR, 2023. URL <https://proceedings.mlr.press/v202/leviathan23a.html>.
- 586 Haoyang LI, Yiming Li, Anxin Tian, Tianhao Tang, Zhanchao Xu, Xuejia Chen, Nicole HU, Wei
587 Dong, Li Qing, and Lei Chen. A survey on large language model acceleration based on KV
588 cache management. *Transactions on Machine Learning Research*, 2025. ISSN 2835-8856. URL
589 <https://openreview.net/forum?id=z3JZzu9EA3>.
- 590

- 594 Yuhong Li, Yingbing Huang, Bowen Yang, Bharat Venkitesh, Acyr Locatelli, Hanchen Ye,
595 Tianle Cai, Patrick Lewis, and Deming Chen. Snapkv: Llm knows what you are look-
596 ing for before generation. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Pa-
597 quet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Sys-*
598 *tems*, volume 37, pp. 22947–22970. Curran Associates, Inc., 2024. doi: 10.52202/
599 079017-0722. URL [https://proceedings.neurips.cc/paper_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2024/file/28ab418242603e0f7323e54185d19bde-Paper-Conference.pdf)
600 [2024/file/28ab418242603e0f7323e54185d19bde-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2024/file/28ab418242603e0f7323e54185d19bde-Paper-Conference.pdf).
- 601 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- 602
603
- 604 Xinyin Ma, Guangnian Wan, Runpeng Yu, Gongfan Fang, and Xinchao Wang. Cot-valve: Length-
605 compressible chain-of-thought tuning. *arXiv preprint arXiv:2502.09601*, 2025.
- 606
- 607 William Merrill and Ashish Sabharwal. The expressive power of transformers with chain of thought,
608 2023.
- 609 Microsoft. Phi-4-mini technical report: Compact yet powerful multimodal language models via
610 mixture-of-loras, 2025.
- 611
- 612 Jesse Mu, Xiang Li, and Noah Goodman. Learning to compress prompts with gist tokens. In
613 A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in*
614 *Neural Information Processing Systems*, volume 36, pp. 19327–19352. Curran Associates, Inc.,
615 2023. URL [https://proceedings.neurips.cc/paper_files/paper/2023/](https://proceedings.neurips.cc/paper_files/paper/2023/file/3d77c6dcc7f143aa2154e7f4d5e22d68-Paper-Conference.pdf)
616 [file/3d77c6dcc7f143aa2154e7f4d5e22d68-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/3d77c6dcc7f143aa2154e7f4d5e22d68-Paper-Conference.pdf).
- 617 Tergel Munkhbat, Namgyu Ho, Seo Hyun Kim, Yongjin Yang, Yujin Kim, and Se-Young Yun. Self-
618 training elicits concise reasoning in large language models. *arXiv preprint arXiv:2502.20122*,
619 2025.
- 620
- 621 Piotr Nawrot, Adrian Łańcucki, Marcin Chochowski, David Tarjan, and Edoardo M Ponti. Dy-
622 namic memory compression: Retrofitting llms for accelerated inference. *arXiv preprint*
623 *arXiv:2403.09636*, 2024.
- 624
- 625 Matanel Oren, Michael Hassid, Nir Yarden, Yossi Adi, and Roy Schwartz. Transformers are multi-
626 state rnns. *arXiv preprint arXiv:2401.06104*, 2024.
- 627
- 628 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
629 optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- 630
- 631 Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang,
632 Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. Deepseekmath: Pushing the limits of mathe-
633 matical reasoning in open language models, 2024.
- 634
- 635 Yi Shen, Jian Zhang, Jiayun Huang, Shuming Shi, Wenjing Zhang, Jiangze Yan, Ning Wang, Kai
636 Wang, and Shiguo Lian. Dast: Difficulty-adaptive slow-thinking for large reasoning models.
637 *arXiv preprint arXiv:2503.04472*, 2025.
- 638
- 639 Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng,
640 Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. *arXiv preprint*
641 *arXiv: 2409.19256*, 2024.
- 642
- 643 Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally
644 can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- 645
- 646 Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford
647 Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- 648
- 649 Jiaming Tang, Yilong Zhao, Kan Zhu, Guangxuan Xiao, Baris Kasikci, and Song Han. Quest:
650 Query-aware sparsity for efficient long-context llm inference, 2024.
- 651
- 652 Qwen Team. Qwen2.5 technical report, 2025. URL [https://arxiv.org/abs/2412.](https://arxiv.org/abs/2412.15115)
653 [15115](https://arxiv.org/abs/2412.15115).

- 648 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
649 Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von
650 Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Ad-
651 vances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.,
652 2017. URL [https://proceedings.neurips.cc/paper_files/paper/2017/
653 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 654 Heming Xia, Yongqi Li, Chak Tou Leong, Wenjie Wang, and Wenjie Li. Tokenskip: Controllable
655 chain-of-thought compression in llms. *arXiv preprint arXiv:2502.12067*, 2025.
656
- 657 Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming
658 language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- 659 Yuchen Yan, Yongliang Shen, Yang Liu, Jin Jiang, Mengdi Zhang, Jian Shao, and Yueting Zhuang.
660 Infythink: Breaking the length limits of long-context reasoning in large language models. *arXiv
661 preprint arXiv:2503.06692*, 2025.
662
- 663 Peitian Zhang, Zheng Liu, Shitao Xiao, Ninglu Shao, Qiwei Ye, and Zhicheng Dou. Long context
664 compression with activation beacon. In *The Thirteenth International Conference on Learning
665 Representations*, 2025. URL <https://openreview.net/forum?id=1eQT9OzfNQ>.
- 666 Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song,
667 Yuandong Tian, Christopher Ré, Clark Barrett, et al. H2o: Heavy-hitter oracle for efficient gener-
668 ative inference of large language models. *Advances in Neural Information Processing Systems*,
669 36:34661–34710, 2023.
- 670 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
671 Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Sto-
672 ica. Judging llm-as-a-judge with mt-bench and chatbot arena. In A. Oh, T. Naumann,
673 A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Informa-
674 tion Processing Systems*, volume 36, pp. 46595–46623. Curran Associates, Inc., 2023.
675 URL [https://proceedings.neurips.cc/paper_files/paper/2023/file/
676 91f18a1287b398d378ef22505bf41832-Paper-Datasets_and_Benchmarks.
677 pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/91f18a1287b398d378ef22505bf41832-Paper-Datasets_and_Benchmarks.pdf).
- 678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A IMPLEMENTATION DETAILS

We use VeRL (Sheng et al., 2024) as the backbone for our training code. We adapt kvpress (Devoto et al., 2025), a recent library for Transformers KV cache compression for inference (e.g., for evaluation). We decouple RL and distillation during our experiments, so that the teacher data is the same for all BR policies and results can be compared more fairly. Otherwise, different experimental results may depend on the different quality of π_{RL} training, rather than on the methods being tested. We save all π_{RL} sampled trajectories to decouple training, and the top-100 logits for each time step, because saving all logits is extremely storage-intensive. In all training setups, this leads to consistently distilling more than 90% of the probability mass for each token, although it frequently reaches even higher levels throughout training.

Overview of Beacon Learning Before training, we add to the vocabulary of the model a new token, the beacon. In the embedding matrix of the model, the beacon token is initialized as the average of all other embeddings. The model learns to compress normal tokens’ key/value activations because, thanks to the attention mask in Figure 2, otherwise those activations’ information would not be accessible from future tokens.

B TASK DETAILS

Countdown: The task is to combine 3–4 numbers using arithmetic operations to reach a target value. All target values are less than or equal to 100.

Countdown Example Prompt

Using the numbers 25, 10, 7, 3, create an equation that equals 96. You can use basic arithmetic operations (+, -, *, /) and each number can only be used once. Make sure to solve it by thinking step by step. Return the final answer in `<answer>` `</answer>` tags, for example `<answer> (1 + 2) / 3 </answer>`.

StarGraph: We use star graphs (Bachmann & Nagarajan, 2024) with branches of length 5 and up to 25 branches per graph. When creating the dataset, the number of branches for each graph is sampled uniformly from the range [2, 25]. We empirically found that this variable complexity helps the reinforcement learning policy, π_{RL} , by allowing it to gradually learn to solve more complex graphs.

StarGraph Example Prompt

You are given a star graph with the following nodes: 34, 72, [...], 304.

The graph has the following directed edges:

34 -> 72
[...]

Find the path from the center node 34 to the target node 304.

Think step by step about the graph structure and trace the path from the center node to the target node.

Return your answer as a list of nodes representing the path from center to target.

Return the final answer in `<answer>` `</answer>` tags, for example `<answer> [1, 3, 7, 12] </answer>`.

LinSys: We generate systems of linear equations that have a single, unique solution. To ensure an appropriate level of difficulty for each model family, we used two distinct configurations:

- For Qwen models: We generated systems of 4 equations in 4 variables. Each equation has at most two non-zero coefficients, and the maximum absolute value for any coefficient is 20.

PPO Hyperparameter	Value	AdamW Hyperparameter	Value
Actor Learning rate	1×10^{-6}	Weight decay	0.01
Critic Learning rate	1×10^{-5}	β_1	0.9
Clip ratio (ϵ)	0.2	β_2	0.999
Number of epochs	1	Epsilon (ϵ)	1×10^{-8}
Mini-batch size	256		
Discount factor (γ)	1.0		
GAE parameter (λ)	1.0		
Entropy coefficient	0.001		
Value loss coefficient	0.5		
Clip range value	0.5		
Max gradient norm	1.0		

Table 4: Hyperparameters for PPO (left) and AdamW (right) for π_{RL} .

Hyperparameter	Value
Weight decay	0.01
β_1	0.9
β_2	0.999
Epsilon (ϵ)	1×10^{-8}

Table 5: Hyperparameters for AdamW optimizer for π_{BR} .

- For Phi models: We found the 4x4 configuration to be too simple (over 40% accuracy before training). We therefore created a more challenging setup: systems of 3 equations in 3 variables, with no restrictions on the number of non-zero coefficients and a maximum absolute coefficient value of 20.

LinSys Example Prompt

Solve the following system of linear equations:

$$2 * x_1 - x_2 + 3 * x_3 + 4 * x_4 = 10$$

$$-x_1 + 4 * x_2 - 2 * x_3 + x_4 = 5$$

$$3 * x_1 + x_2 + x_3 - x_4 = 12$$

$$x_1 + 2 * x_2 + 5 * x_3 + 2 * x_4 = 20$$

Find the values for x_1 , x_2 , ..., x_4 . Make sure to solve it by thinking step by step, and do not assume access to any external tools.

Return the final answer as a list of numbers in `<answer>` `</answer>` tags, for example `<answer>[1, -2, 3, 7]</answer>`.

C TRAINING DETAILS

C.1 π_{RL} TRAINING

Table 4 provides the hyperparameters for PPO Schulman et al. (2017) and AdamW (Loshchilov & Hutter, 2019). The critic has the same architecture and initial weights as π_{RL} .

C.2 π_{BR} TRAINING

Table 5 provides the hyperparameters for the AdamW optimizer. We use a learning rate of 5×10^{-6} for all tasks and models, except for the Countdown-Phi configuration, where we use a higher learning rate of 5×10^{-5} . This choice is based on empirical observations.

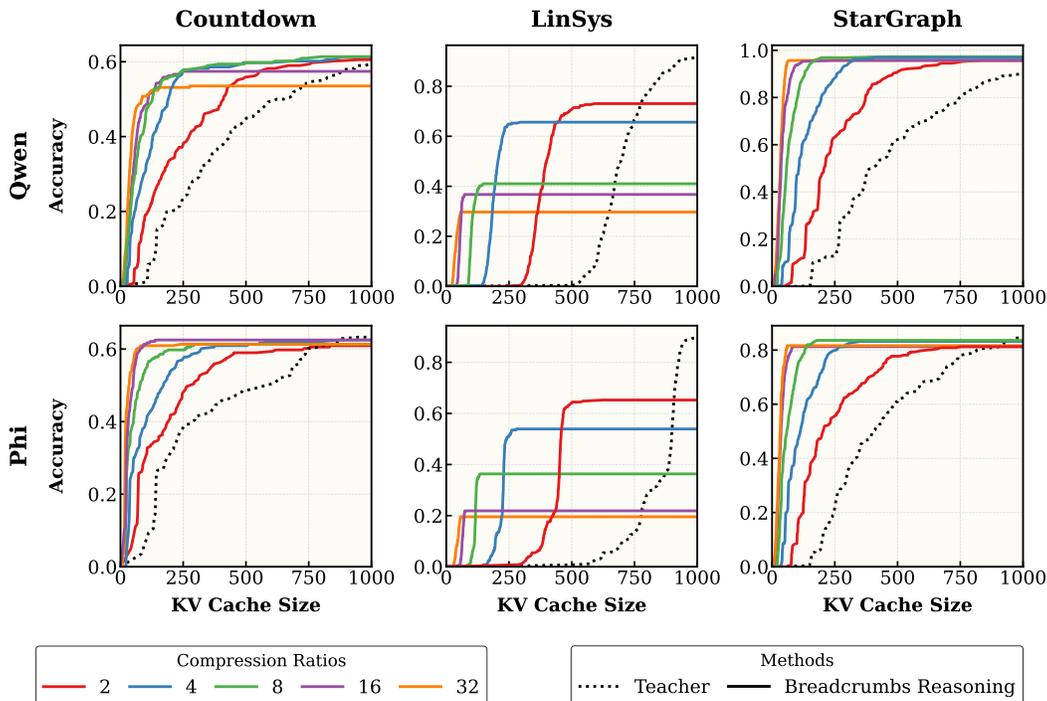


Figure 6: **Accuracy vs. KV Cache Size** Single Ratio Breadcrumbs Reasoning retains most of the teacher’s performance while using significantly fewer KV cache entries, even outperforming the teacher when setting a fixed maximum KV cache size in Countdown and StarGraph.

D ANALYSIS

LLM-as-a-Judge Pipeline. We analyze failures of Qwen on the LinSys task in two stages. First, for every incorrect generation we pair the model’s reasoning trace with the ground truth and prompt a verifier LLM (Qwen3 30B A3B Thinking 2507) to identify the *first* erroneous step; the resulting snippet (e.g., “computed $100 - 3100$ as -2999 ”) is stored. Second, we pass each extracted snippet to a classification judge that sees only the mistake text plus a fixed taxonomy (Arithmetic, Equation Misreading & Fabrication, Algebraic Manipulation, Output, Other), that we identify by manual inspection of the mistakes (and add Other to allow for cases that we do not explicitly categorize). Using best-of- n sampling ($n=3$), the judge assigns the most specific category.

E ADDITIONAL RESULTS

We report in Table 6 the accuracy area under the curve (AUAC) metric with varying the maximum cache size up to 1,000.

Figure 6 shows the Pareto curves for Single Ratio Breadcrumbs Reasoning.

Figure 7 illustrates the increase in accuracy of BR when increasing the generation length from 1,000 to $1,000 \times c$ (i.e., KV cache size of 1,000 entries).

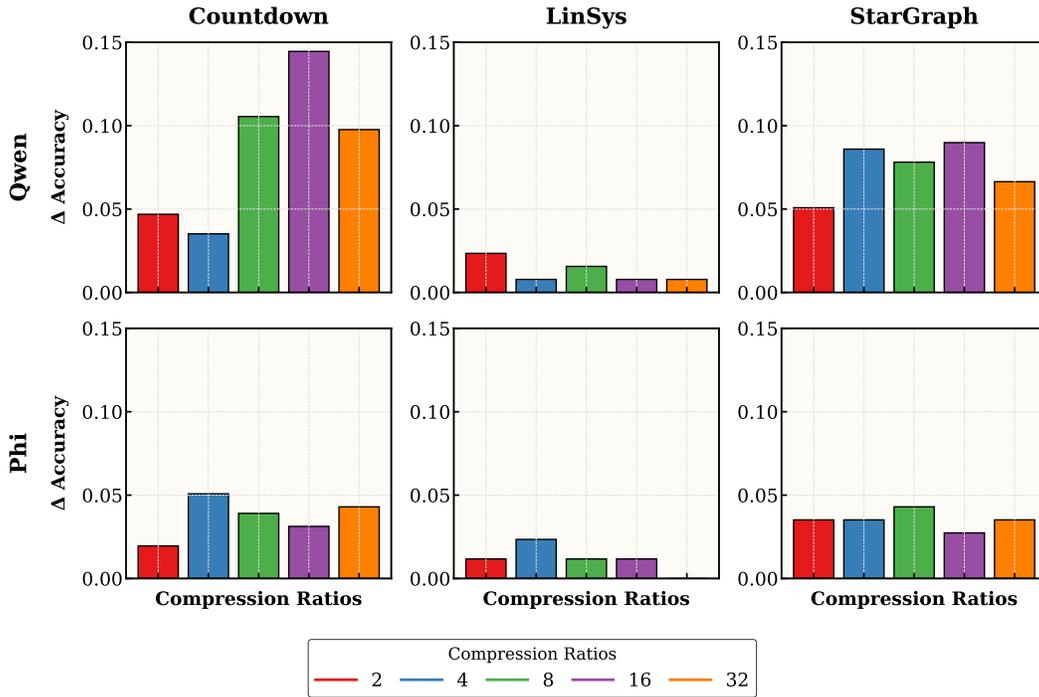
F LLM USAGE

LLMs were used in the process of writing this paper to assist in creating tables and figures.

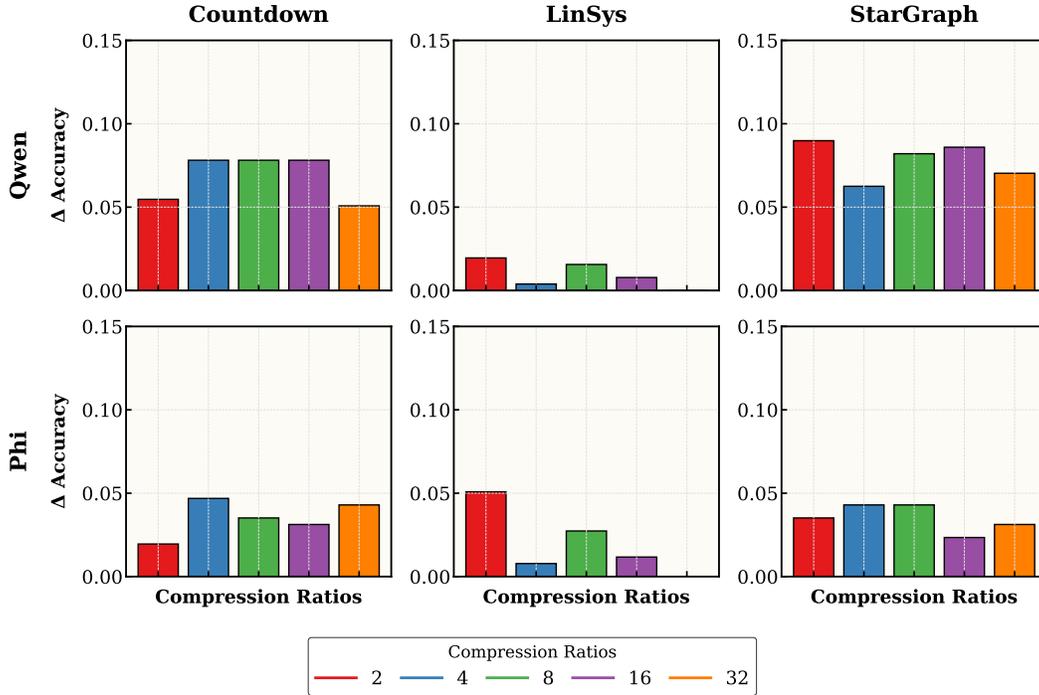
QWEN																		
COMPR. RATIO	COUNTDOWN						LINSYS						STARGRAPH					
	1	2	4	8	16	32	1	2	4	8	16	32	1	2	4	8	16	32
NO COMPRESSION																		
TEACHER	0.377	-	-	-	-	-	0.276	-	-	-	-	-	0.513	-	-	-	-	-
BREADCRUMBS REASONING																		
Sr BR (OURS)	-	0.465	0.533	<u>0.555</u>	<u>0.539</u>	<u>0.513</u>	-	0.451	0.532	<u>0.368</u>	<u>0.347</u>	<u>0.286</u>	-	0.726	0.845	<u>0.906</u>	<u>0.918</u>	<u>0.926</u>
MR BR (OURS)	-	0.446	<u>0.515</u>	0.576	0.576	0.560	-	<u>0.439</u>	<u>0.507</u>	0.424	0.391	0.316	-	<u>0.713</u>	<u>0.839</u>	0.912	0.921	0.935
TRAINING-FREE COMPRESSION																		
PYRAMIDKV	-	0.375	0.201	0.111	0.077	0.011	-	0.000	0.000	0.000	0.000	0.000	-	0.487	0.459	0.491	0.463	0.526
SNAPKV	-	0.376	0.210	0.138	0.100	0.079	-	0.003	0.000	0.000	0.000	0.000	-	0.534	0.478	0.485	0.452	0.489
TOVA	-	<u>0.447</u>	0.272	0.167	0.183	0.199	-	0.000	0.000	0.000	0.000	0.000	-	0.526	0.417	0.424	0.414	0.423
STREAMINGLLM	-	0.007	0.015	0.024	0.049	0.090	-	0.000	0.000	0.000	0.000	0.000	-	0.038	0.032	0.023	0.041	0.105
PHI																		
COMPR. RATIO	COUNTDOWN						LINSYS						STARGRAPH					
	1	2	4	8	16	32	1	2	4	8	16	32	1	2	4	8	16	32
NO COMPRESSION																		
TEACHER	0.433	-	-	-	-	-	0.142	-	-	-	-	-	0.498	-	-	-	-	-
BREADCRUMBS REASONING																		
Sr BR (OURS)	-	0.506	0.557	<u>0.581</u>	0.604	<u>0.597</u>	-	<u>0.373</u>	<u>0.421</u>	<u>0.322</u>	<u>0.205</u>	<u>0.187</u>	-	0.636	0.739	0.786	0.785	0.792
MR BR (OURS)	-	<u>0.480</u>	<u>0.526</u>	0.582	<u>0.568</u>	0.609	-	0.375	0.454	0.408	0.278	0.258	-	<u>0.631</u>	<u>0.710</u>	<u>0.765</u>	<u>0.732</u>	<u>0.738</u>
TRAINING-FREE COMPRESSION																		
PYRAMIDKV	-	0.417	0.225	0.180	0.019	0.000	-	0.011	0.000	0.000	0.000	0.000	-	0.503	0.366	0.292	0.327	0.358
SNAPKV	-	0.406	0.222	0.149	0.062	0.012	-	0.006	0.000	0.000	0.000	0.000	-	0.516	0.424	0.329	0.331	0.304
TOVA	-	0.196	0.064	0.050	0.046	0.094	-	0.000	0.000	0.000	0.000	0.000	-	0.621	0.504	0.433	0.436	0.378
STREAMINGLLM	-	0.014	0.018	0.060	0.150	0.300	-	0.000	0.000	0.000	0.000	0.000	-	0.137	0.017	0.029	0.089	0.097

Table 6: **Model performance (AUAC) on long-context reasoning tasks.** We report AUAC given a maximum cache size of 1,000 entries. Higher is better.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971



(a) Single Ratio Breadcrumb Reasoning



(b) Multi Ratio Breadcrumbs Reasoning

Figure 7: **Performance increase with extended generation.** Breadcrumbs Reasoning improves up to 14.5% with Qwen and 5.1% with Phi beyond the 1,000 token training length.