

---

# AGENTKV: Phase-Aware KV Eviction for Agentic LLMs

---

Taowen (Tony) Liu<sup>1</sup> Jeffrey T. H. Wong<sup>1</sup> Can Xiao<sup>1</sup> Bowen Yang<sup>1</sup> Hao (Mark) Chen<sup>1</sup> Yiren Zhao<sup>1</sup>

## Abstract

Agentic workloads can consume on average  $1,000\times$  more tokens than chatbot workloads (Bai et al., 2026), stressing both KV-cache capacity and decode-time bandwidth. Existing query-centric eviction methods estimate key importance using representative queries drawn from recent tokens, an approach that works when future attention resembles recent attention. We show that agentic generation violates this assumption: its future query distribution is a mixture over `think`, `act`, `tool`, and `others` phases, whose components occupy measurably different query subspaces. As a result, recency representatives over-represent the current phase and under-score keys needed by other phases. We propose AGENTKV, a phase-aware eviction method that maintains a small set of queries per phase and scores cached keys against their union. Across model-task-budget combinations, AGENTKV ranks first or tied first among compressed-cache methods in 78% BFCL tasks and 61%  $\tau^2$ -bench tasks, and improves output-token throughput by up to  $1.80\times$ .

## 1. Introduction

Agentic reasoning is becoming a central use case for large language models (LLMs), but it is substantially more expensive to serve than conventional single-turn chatbot workloads. Agentic workloads can consume  $1,000\times$  more tokens than conventional chatbot workloads (Bai et al., 2026), dramatically increasing sequence length and KV-cache pressure. This high cost arises because, unlike single-turn chatbot interactions, agentic workloads require models to solve tasks across many turns by alternating among recurring phases of reasoning (`think`), action generation (`act`), tool results (`tool`), and other contextual processing (`others`) (Yao et al., 2022). We refer to these recurring stages as **agentic**

**phases**. As a result, KV-cache memory capacity and memory bandwidth have become major bottlenecks for serving throughput and latency.

A common approach to reducing this bottleneck is KV-cache eviction or compression. Many recent methods can be viewed as estimating which cached keys will be useful for future attention queries. Representative approaches retain heavy hitters (Zhang et al., 2023), exploit attention sinks for streaming generation (Xiao et al., 2024), select informative tokens from local observation windows (Adnan et al., 2024), or prune keys using recent attention and key-token saliency (Liu et al., 2023; Li et al., 2024). These methods are effective when recent attention or recent queries are predictive of future attention. However, this assumption is poorly matched to multi-turn agentic inference.

This mismatch arises because multi-turn agent trajectories do not induce a single stationary query distribution. As a result, eviction policies that extrapolate from the most recent query window primarily estimate the current-phase query distribution, and can under-score keys that are weakly relevant to the current phase but important for future phases. We empirically verify this phase structure by analyzing the geometry of phase-specific query representations: principal-angle analysis shows that queries from different agentic phases occupy distinguishable subspaces in the attention-query representation space. This observation motivates AGENTKV, which models the future-query distribution as a mixture over phase-specific query distributions rather than as a continuation of the most recent query window. AGENTKV maintains small query buffers for each agentic phase and estimates key importance by scoring candidate KV entries against the union of these phase representatives.

Policy-level eviction rules alone, however, are insufficient for efficient agentic serving. In persistent multi-turn inference, compressed KV state must be maintained across turns, compacted after eviction, and reused without introducing large memory-copy or allocator overheads. To realize the practical acceleration potential of AGENTKV, we implement it in a persistent serving path with online KV-page compaction and allocator-level reuse.

This paper makes three contributions. **(1)** we identify phase-induced estimator bias in recency-based KV eviction and verify phase-dependent query subspaces using principal-

---

<sup>1</sup>Imperial College London, London, United Kingdom. Correspondence to: Taowen Liu <tony.liu20@imperial.ac.uk>.

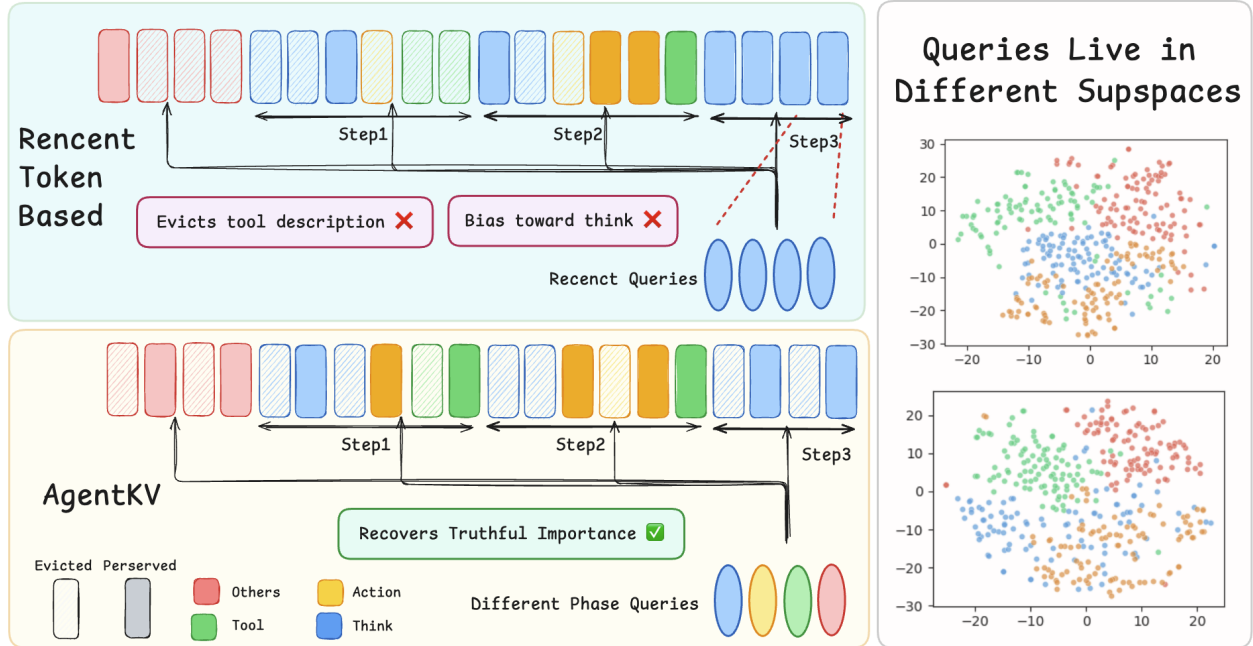


Figure 1. Left: recency-based KV eviction estimates the future-query distribution from the current phase, leaving keys aligned with non-recent phases (e.g., tool, others) undervalued. AGENTKV draws representatives from each phase so that  $\hat{p}(q)$  tracks the four-phase mixture.

Right: t-SNE visualization of query vectors, colored by ground-truth phase: **blue** = think, **yellow** = act, **green** = tool, **red** = others (100 queries per phase). Each phase occupies a distinguishable region. Top graph shows layer 24 head 32, and bottom graph shows layer 48 head 24.

angle analysis. (2) we propose AGENTKV, to our knowledge the first KV-eviction method that explicitly accounts for phase-mixture query distributions in agentic inference. (3) we build a persistent multi-turn serving implementation with online KV-page compaction and allocator-level reuse to reduce compressed-cache serving overheads.

Empirically, across two models, three task domains, and multiple KV budgets, AGENTKV ranks first or tied first among compressed-cache methods in 78% BFCL combinations and 61%  $\tau^2$ -bench combinations. In persistent multi-turn serving, AGENTKV improves output-token throughput by up to  $1.80\times$ .

## 2. Background

### 2.1. Agentic Trajectories

Tool-augmented agents (Yao et al., 2022; Schick et al., 2023; Qin et al., 2024) produce multi-turn trajectories that interleave reasoning and tool use (Figure 2). We define an *agentic phase* as a semantic token category in an agent trajectory, rather than as a turn index. Each emitted token is assigned to one of four phases: `think` (assistant reasoning spans), `act` (assistant tool-call arguments), `tool` (tool-response observations), and `others` (system prompt, tool descrip-

tions, user question). All tokens with the same label belong to the same phase across the trajectory; when an analysis compares steps within a phase, it compares different contiguous spans with that label (e.g., the `think` span before one tool call versus the `think` span before another).

### 2.2. Query-Centric KV Cache Eviction

KV eviction methods reduce memory usage by retaining only the important KV entries under a fixed budget  $B$ . Query-centric methods estimate token importance by comparing stored keys against a small set of selected past queries, which we refer to as *representative queries*, and keeping the top- $B$  tokens according to the importance scores computed against representative queries (Li et al., 2024; Cai et al., 2026; Zhang et al., 2023).

Upon compression, suppose the model has emitted query vectors  $q_t^{(l,h)}$  for positions  $t = 1, 2, \dots, T$ , where  $q_t^{(l,h)} \in \mathbb{R}^{d_h}$  denotes the query vector at layer  $l$  and attention head  $h$ . For brevity, we omit the layer and head indices and write  $q_t$ . Let  $\mathcal{Q} \subseteq q_1, \dots, q_T$  denote the representative query set used to approximate the unknown future queries. Given candidate keys  $K = [k_1, \dots, k_T]$ , query-centric methods compute the



Figure 2. Structure of an agentic trajectory with 4 semantic phases: think, act, tool, and others.

average attention mass assigned to each candidate token  $i$ :

$$A_i = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \text{softmax} \left( \frac{q^\top K}{\sqrt{d_h}} \right)_i,$$

where the softmax is taken along the key axis, so each representative query induces a probability distribution over candidate keys. Each candidate token is then assigned a score  $\text{score}_i = f(A_i)$ , and the top- $B$  tokens are retained while the rest are evicted.

A range of KV eviction methods have been proposed under the query-centric framework, differing primarily along two axes: *which tokens to retain* and *when to apply compression*. H2O (Zhang et al., 2023) retains heavy hitters identified by accumulated attention mass, while StreamingLLM (Xiao et al., 2024) preserves a fixed combination of attention sinks and a sliding recent window. SnapKV (Li et al., 2024) selects tokens based on attention from a short observation window immediately preceding generation. R-KV (Cai et al., 2026) further refines the SnapKV-style formulation by introducing a redundancy-aware objective  $\text{score}_i = \lambda A_i - (1-\lambda)R_i$ , where  $R_i$  is the redundancy penalty.

AGENTKV is complementary to these directions. It changes how the representative query set  $\mathcal{Q}$  is estimated: rather than drawing representatives only from a recent window, it draws them across agentic phases.

### 3. KV Cache Eviction Method for Agents

Let  $p(q)$  be the unknown future query distribution at a compression event, and  $\hat{p}(q)$  an estimator of it built from past queries. We frame KV eviction as ranking keys by their expected attention under  $\hat{p}$ :

$$A_i = \mathbb{E}_{q \sim \hat{p}(q)} \left[ \text{softmax} \left( \frac{q^\top K}{\sqrt{d}} \right)_i \right].$$

Recency representatives instantiate this with  $\hat{p}(q) = p_{\text{recent}}(q)$ , the empirical distribution over the last few queries. This is reliable only if  $p(q)$  is stationary across the trajectory. We ask: **(i)** what structure does  $p(q)$  exhibit in agentic settings (§3.1); and **(ii)** how much bias does  $p_{\text{recent}}$  introduce, and what are the costs and benefits of correcting it (App. C.2).

#### 3.1. Phase Structure Exists in $p(q)$

Drawing on prior work showing phases are linearly detectable from hidden states (Ye et al., 2026), we hypothesize that query vectors from each phase occupy geometrically distinct subspaces. We examine query distributions across attention heads and layers with t-SNE. Right side of figure 1 shows that queries from different agentic phases exhibit intra-phase clustering.

We next quantitatively test whether query directions are more separated across phases than within the same phase. We first split each trajectory by phase label into contiguous spans, e.g., `think`<sub>1</sub>, `act`<sub>1</sub>, `tool`<sub>1</sub>, `think`<sub>2</sub>, `act`<sub>2</sub>, `tool`<sub>2</sub>. We use span’s phase label to decide whether a pair of spans are inter-phase or intra-phase. For each span  $s$ , layer, and head, we collect all query vectors from tokens in that span and stack them into a matrix  $X_s \in \mathbb{R}^{n_s \times d}$ , where  $s$  indexes a span,  $n_s$  is the number of tokens in that span, and  $d$  is the head dimension. We center the rows of  $X_s$  and compute the SVD  $X_s = L_s \Sigma_s V_s^\top$ . The top- $k$  right singular vectors  $V_{s,k} \in \mathbb{R}^{d \times k}$  form an orthonormal basis for the dominant query directions of span  $s$ . Given two spans  $s$  and  $t$ , their principal angles are defined by the singular values of the basis-overlap matrix:

$$\cos \theta_i = \sigma_i(V_{s,k}^\top V_{t,k}), \quad i = 1, \dots, k.$$

We summarize each span pair by the mean principal angle  $\bar{\theta} = \frac{1}{k} \sum_i \theta_i$ . A larger  $\bar{\theta}$  means the two spans occupy less aligned query directions. We then compare inter-phase pairs, whose spans have different phase labels, against intra-phase pairs, whose spans have the same label but occur at different steps (e.g. `think` span in one step versus `think` span in another step). Figures 3 and 4 show that inter-phase angles are consistently larger than intra-phase angles across heads and layers (see Appendix B), confirming that phase labels correspond to separable query subspaces. This structured partitioning has an important implication: recency-based cache representatives cannot faithfully approximate the query distribution of a future phase, as we will show next.

#### 3.2. AGENTKV

The distribution analysis (§3.1) shows that different phases are associated with different query directions. AGENTKV exploits this directly: keep a small per-phase representative set of recent queries for each of the four phases, and at each

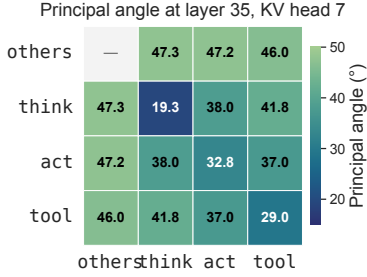


Figure 3. Mean principal angles between phase-specific query subspaces (5 dimensional). Off-diagonal entries show cross-phase angles; diagonal entries show within-phase angles between step- $i$  and step- $j$ .

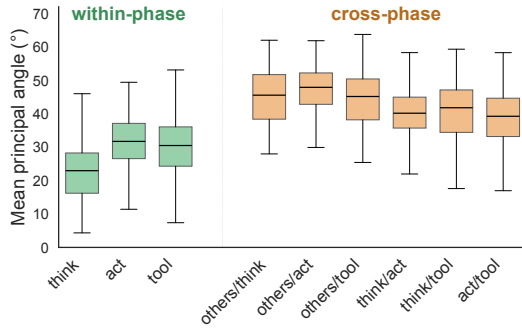


Figure 4. Distribution of mean principal angles for each comparison, pooled across all layer-head-compression-event tuples.

compression event score keys against the union of these representative sets rather than a single recency window.

AGENTKV requires a phase assignment for emitted queries, but it does not use a model-side classifier or any hidden-state phase detector. Instead, the phase assignment is implemented as a low-cost parser over tokenizer-visible text structure. We use a four-state machine to the tokenizer’s chat-template emission: every decoded token carries the current phase label from  $\{\text{think}, \text{act}, \text{tool}, \text{others}\}$ , with state transitions fired by the chat-template role tokens and markers. For each phase, we keep a small FIFO buffer  $Q_\phi$  that always holds the  $\frac{q}{4}$  most-recent query vectors emitted in that phase; older queries from the same phase are overwritten. At each compression event, the four buffers are merged into:

$$Q^{(e)} = Q_{\text{think}} \cup Q_{\text{act}} \cup Q_{\text{tool}} \cup Q_{\text{others}}$$

Thus  $|Q^{(e)}| = q$ , matching the representative count of the recency baseline.

Given candidate keys  $K^{(e)}$  at compression event  $e$ , AGENTKV scores every candidate key against this union. Each candidate token  $i$  is scored by its average attention from the

phase-aware representatives:

$$A_i^{(e)} = \frac{1}{|Q^{(e)}|} \sum_{q \in Q^{(e)}} \text{softmax}\left(\frac{q^\top K^{(e)}}{\sqrt{d_h}}\right)_i.$$

It then keeps the top- $B$  candidates by  $A_i^{(e)}$  and evicts the rest.

## 4. Experiments

### 4.1. Task Quality under Compression

#### 4.1.1. EXPERIMENT SETUP

We use a persistent multi-turn cache setting and compare AGENTKV against recent compressed-cache baselines on BFCL and  $\tau^2$ -bench. Detailed baseline, benchmark, and serving setup are in App. D.1.

**BFCL results.** Table 1 reports accuracy across the BFCL categories. Across the 18 model-task-budget combinations, AGENTKV ranks first or tied first among compressed-cache methods in 14 combinations. The `memory` task also illustrates a complementary interaction with external memory: it relies on mechanisms such as key-value storage, vector storage, or recursive summaries, rather than only on long-range dependencies inside the model’s attention cache. In this setting, eviction can reduce stale or distracting in-cache context while preserving the state needed to query the external memory.

**$\tau^2$ -bench results.** Table 2 reports the  $\tau^2$ -bench results. Across the 18 model-task-budget combinations, AGENTKV ranks first or tied first among compressed-cache methods in 11 combinations. The smaller advantage on `airline` may reflect its task structure: tool responses expose much of the current environment state, so many turns can be resolved from the current user request and recent tool response.

### 4.2. Serving Throughput

AGENTKV is designed for the target stateful serving path, where the KV cache persists across tool calls and future turns as trajectory state. Thus, after each compression event, retained KV entries must be physically compacted and evicted pages returned to the serving allocator; a smaller logical budget alone would not reduce decode-time memory traffic or improve concurrency. Our implementation performs online in-place compaction, allocator bookkeeping, and state carry-over across tool responses and subsequent user turns, while remaining compatible with chunked prefill, CUDA graph capture, and the normal decode path. Serving setup details are in App. D.1.

Figure 5 shows that AGENTKV improves throughput over

Table 1. BFCL accuracy. **Bold** indicates the best compressed-cache result; Full KV Cache is shown as an uncompressed reference.

Method	web_search			multi_turn			memory		
	512	768	1024	2048	3072	4096	1024	1536	2048
Qwen3-32B									
Full KV Cache	28.0	28.0	28.0	52.0	52.0	52.0	24.8	24.8	24.8
AGENTKV	<b>15.0</b>	<b>22.0</b>	22.0	<b>22.0</b>	<b>42.0</b>	<b>48.0</b>	23.2	<b>21.7</b>	22.8
R-KV	5.0	14.0	<b>25.0</b>	13.0	27.0	37.0	23.2	21.3	23.9
Tri-attention	5.0	21.0	21.0	<b>22.0</b>	41.5	44.0	<b>24.3</b>	21.5	<b>26.9</b>
Qwen3-8B									
Full KV Cache	11.0	11.0	11.0	34.5	34.5	34.5	14.4	14.4	14.4
AGENTKV	<b>12.0</b>	<b>14.0</b>	16.0	<b>17.0</b>	<b>26.5</b>	<b>38.0</b>	<b>18.7</b>	<b>19.1</b>	<b>18.7</b>
R-KV	6.0	6.0	<b>19.0</b>	5.0	14.0	20.0	17.2	16.8	15.9
Tri-attention	3.0	7.0	12.0	7.5	18.0	23.0	12.3	15.7	14.0

Table 2.  $\tau^2$ -bench scores. **Bold** indicates the best compressed-cache result; Full KV Cache is shown as an uncompressed reference.

Method	airline			retail			telecom		
	2048	4096	6144	2048	4096	6144	2048	4096	6144
Qwen3-32B									
Full KV Cache	44.0	44.0	44.0	53.5	53.5	53.5	25.4	25.4	25.4
AGENTKV	<b>36.0</b>	34.0	40.0	<b>8.8</b>	<b>42.1</b>	<b>50.9</b>	6.1	<b>18.4</b>	21.9
R-KV	34.0	36.0	40.0	6.1	26.3	44.7	5.3	7.9	<b>22.8</b>
Tri-attention	30.0	<b>40.0</b>	<b>44.0</b>	7.9	16.7	39.5	<b>7.9</b>	10.5	20.2
Qwen3-8B									
Full KV Cache	40.0	40.0	40.0	42.1	42.1	42.1	15.8	15.8	15.8
AGENTKV	26.0	<b>36.0</b>	<b>34.0</b>	<b>10.5</b>	<b>27.2</b>	<b>30.7</b>	7.9	14.0	<b>18.4</b>
R-KV	<b>40.0</b>	<b>36.0</b>	30.0	5.3	13.2	28.9	10.5	<b>18.4</b>	<b>18.4</b>
Tri-attention	32.0	18.0	<b>34.0</b>	5.3	11.4	21.9	<b>14.6</b>	14.9	12.3

upstream SGLang in two regimes. At 16k/1k, decode is primarily bandwidth-bound because long KV caches are repeatedly read, so reducing KV traffic improves throughput. At 32k/1k, GPU memory capacity becomes the dominant bottleneck: the baseline supports fewer parallel requests, whereas AGENTKV frees pages and enables higher parallelism. On RTX A6000, Qwen3-32B does not fit in memory, so we report Qwen3-8B, whose best output throughput improves by **1.37** $\times$  at 16k/1k and **1.76** $\times$  at 32k/1k. On H100 NVL, Qwen3-8B improves by **1.31** $\times$  at 32k/1k, while Qwen3-32B improves by **1.64** $\times$  at 16k/1k and **1.80** $\times$  at 32k/1k.

### 4.3. Mechanism Ablation

#### 4.3.1. HOW MANY REPRESENTATIVE QUERIES DO WE NEED?

We sweep the online query-representative budget with the same phase-aware policy used in the main experiments.

Table 3 shows that representative budgets above  $q=32$  do not consistently improve end-to-end quality. We therefore use  $q=32$  in the main AGENTKV results.

Table 3. End-to-end query-representative budget ablation for AGENTKV on Qwen3-32B, using two representative settings: BFCL `multi_turn_base` at  $B=3072$  and  $\tau^2$ -bench `retail` at  $B=4096$ .

$q$	BFCL acc. (%)	$\tau^2$ score (%)
8	33.5	32.5
16	38.0	34.2
24	36.5	36.8
32	<b>42.0</b>	<b>42.1</b>
48	40.0	38.6
64	36.5	33.3

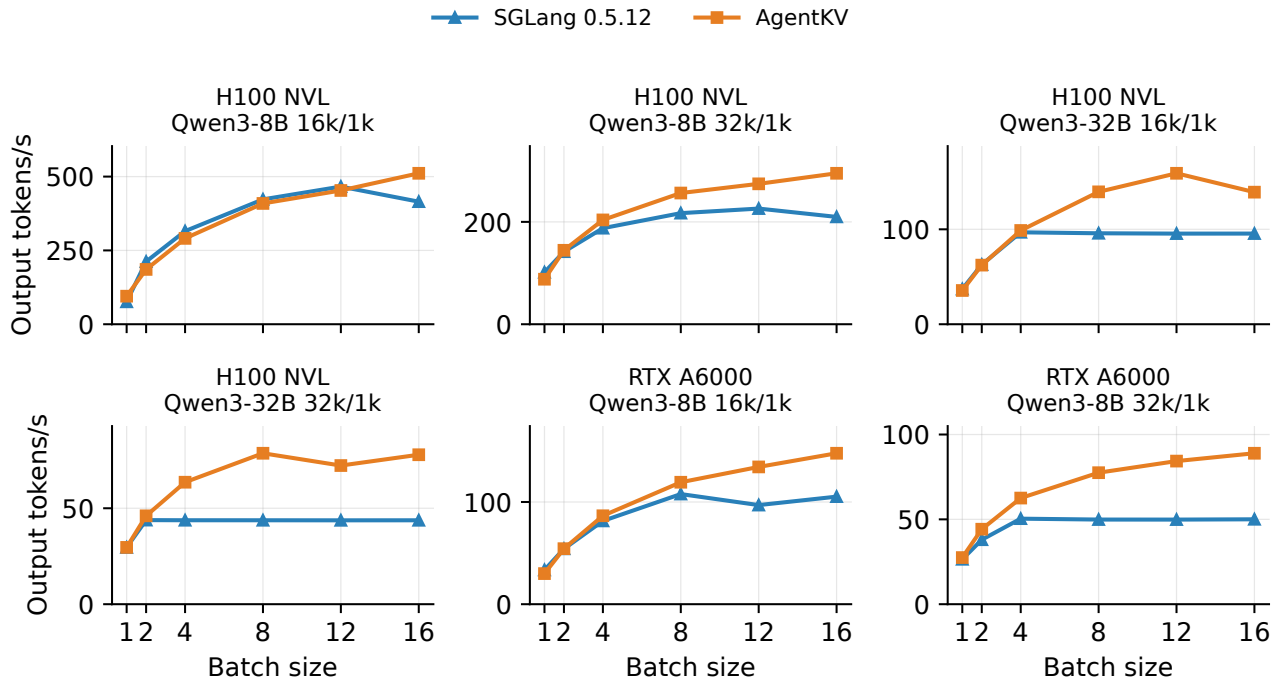


Figure 5. Serving throughput in output tokens per second. Curves compare upstream SGLang 0.5.12 and AGENTKV.

#### 4.3.2. DO WE NEED ALL PHASES?

Table 4 ablates one phase at a time from the representative set. Removing any phase reduces accuracy relative to full AGENTKV, confirming that the four phase buffers provide complementary coverage.

Table 4. Phase-representative ablation on Qwen3-32B BFCL `multi_turn_base` at  $B=3072$ .

Representative strategy	Accuracy (%)
Full AGENTKV	42.0
Drop <code>think</code>	38.5
Drop <code>act</code>	37.5
Drop <code>tool</code>	38.0
Drop <code>others</code>	35.5

## 5. Conclusion

We presented AGENTKV, a phase-aware KV cache eviction method for agentic reasoning. Across model–task–budget combinations, AGENTKV ranks first or tied first among compressed-cache methods in 78% BFCL combinations and 61%  $\tau^2$ -bench combinations, while improving persistent multi-turn output-token throughput by up to  $1.80\times$ . These findings suggest that phase-aware representatives are useful for stateful multi-turn tool-use workloads, where non-recent instructions and observations remain important.

## Limitations

AGENTKV currently uses a manually specified KV budget rather than selecting the budget automatically from a task-level quality or latency target. This can be mitigated by tuning the budget on a small validation set for the target workload, as is common for KV-compression deployment.

Like other KV-compression methods, overly aggressive eviction can degrade task quality, so deployments should validate task-level performance at their chosen budget.

## References

- Abhyankar, R., He, Z., Srivatsa, V., Zhang, H., and Zhang, Y. Infercept: Efficient intercept support for augmented large language model inference. *arXiv preprint arXiv:2402.01869*, 2024.
- Adnan, M., Arunkumar, A., Jain, G., Nair, P. J., Soloveychik, I., and Kamath, P. Keyformer: Kv cache reduction through key tokens selection for efficient generative inference. *Proceedings of Machine Learning and Systems*, 6:114–127, 2024.
- Bai, L., Huang, Z., Wang, X., Sun, J., Mihalcea, R., Brynjolfsson, E., Pentland, A., and Pei, J. How do coding agents spend your money? analyzing and predicting token consumptions in agentic coding tasks. *Preprint*, 2026.
- Cai, Z., Xiao, W., Sun, H., Zhang, Y., Wan, K., Li, Y., Zhou, Y., Chang, L.-W., Gu, J., Dong, Z., et al. R-kv: Redundancy-aware kv cache compression for reasoning models. *Advances in Neural Information Processing Systems*, 38:60980–61005, 2026.
- Gorilla LLM Team. Berkeley function-calling leaderboard (BFCL), 2025. URL <https://gorilla.cs.berkeley.edu/leaderboard>. Accessed: 2026-03-06.
- Li, H., Mang, Q., He, R., Zhang, Q., Mao, H., Chen, X., Zhou, H., Cheung, A., Gonzalez, J., and Stoica, I. Continuum: Efficient and robust multi-turn llm agent scheduling with kv cache time-to-live. *arXiv preprint arXiv:2511.02230*, 2025.
- Li, Y., Huang, Y., Yang, B., Venkitesh, B., Locatelli, A., Ye, H., Cai, T., Lewis, P., and Chen, D. Snapkv: Llm knows what you are looking for before generation. *Advances in Neural Information Processing Systems*, 37:22947–22970, 2024.
- Liu, Z., Desai, A., Liao, F., Wang, W., Xie, V., Xu, Z., Kyrillidis, A., and Shrivastava, A. Scissorhands: Exploiting the persistence of importance hypothesis for llm kv cache compression at test time. *Advances in Neural Information Processing Systems*, 36:52342–52364, 2023.
- Mao, W., Lin, X., Huang, W., Xie, Y., Fu, T., Zhuang, B., Han, S., and Chen, Y. Triattention: Efficient long reasoning with trigonometric kv compression. *arXiv preprint arXiv:2604.04921*, 2026.
- Patil, S. G., Mao, H., Yan, F., Ji, C. C.-J., Suresh, V., Stoica, I., and Gonzalez, J. E. The berkeley function calling leaderboard (bfcl): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*, 2025.
- Qin, Y., Liang, S., Ye, Y., Zhu, K., Yan, L., Lu, Y., Lin, Y., Cong, X., Tang, X., Qian, B., et al. Toolllm: Facilitating large language models to master 16000+ real-world apis. In *International Conference on Learning Representations*, volume 2024, pp. 9695–9717, 2024.
- Qwen Team. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., and Scialom, T. Toolformer: Language models can teach themselves to use tools. *Advances in neural information processing systems*, 36:68539–68551, 2023.
- Shi, Z., Ruan, C., Qi, P., Huang, G., Wan, X., Lin, M., and Li, J. Tetris: Efficient and predictive kv cache offloading for agentic and reasoning workloads. *SOSP 2025 SAA Workshop*, 2025.
- Srivatsa, V., He, Z., Abhyankar, R., Li, D., and Zhang, Y. Preble: Efficient distributed prompt scheduling for llm serving. In *International conference on learning representations*, volume 2025, pp. 37057–37082, 2025.
- Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. Efficient streaming language models with attention sinks. In *International Conference on Learning Representations*, volume 2024, pp. 21875–21895, 2024.
- Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- Yao, S., Shinn, N., Razavi, P., and Narasimhan, K. *tau-bench*: A benchmark for tool-agent-user interaction in real-world domains. *arXiv preprint arXiv:2406.12045*, 2024.
- Ye, C., Cui, J., and Hadfield-Menell, D. Prompt injection as role confusion, 2026. URL <https://arxiv.org/abs/2603.12277>.
- Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., Song, Z., Tian, Y., Ré, C., Barrett, C., et al. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36:34661–34710, 2023.

## A. Related work

**System-level KV management for agents.** A separate line of work optimizes KV management at the system layer rather than the kv cache eviction algorithms, through interception-aware recomputation, TTL-based retention, predictive offloading, or distributed prompt scheduling (Abhyankar et al., 2024; Srivatsa et al., 2025; Li et al., 2025; Shi et al., 2025). These systems make the cache cheaper, faster, and more shareable, but take the underlying token-importance signal as given. Our contribution sits upstream and concerns the *quality* of that signal under a fixed budget; any improvement we obtain layers cleanly on top of these serving-layer optimizations.

## B. Per-(layer, head) phase-angle matrices

Figure 6 expands Figure 3 into a grid of small  $4 \times 4$  phase matrices, one per (layer, KV head) cell, on the same colour scale. Reading down a column shows how a given KV head’s phase geometry evolves with depth; reading across a row shows head-to-head variation within a layer. The diagonal (within-phase) is consistently darker than the off-diagonal (cross-phase) in every cell, with the gap widening at deeper layers.

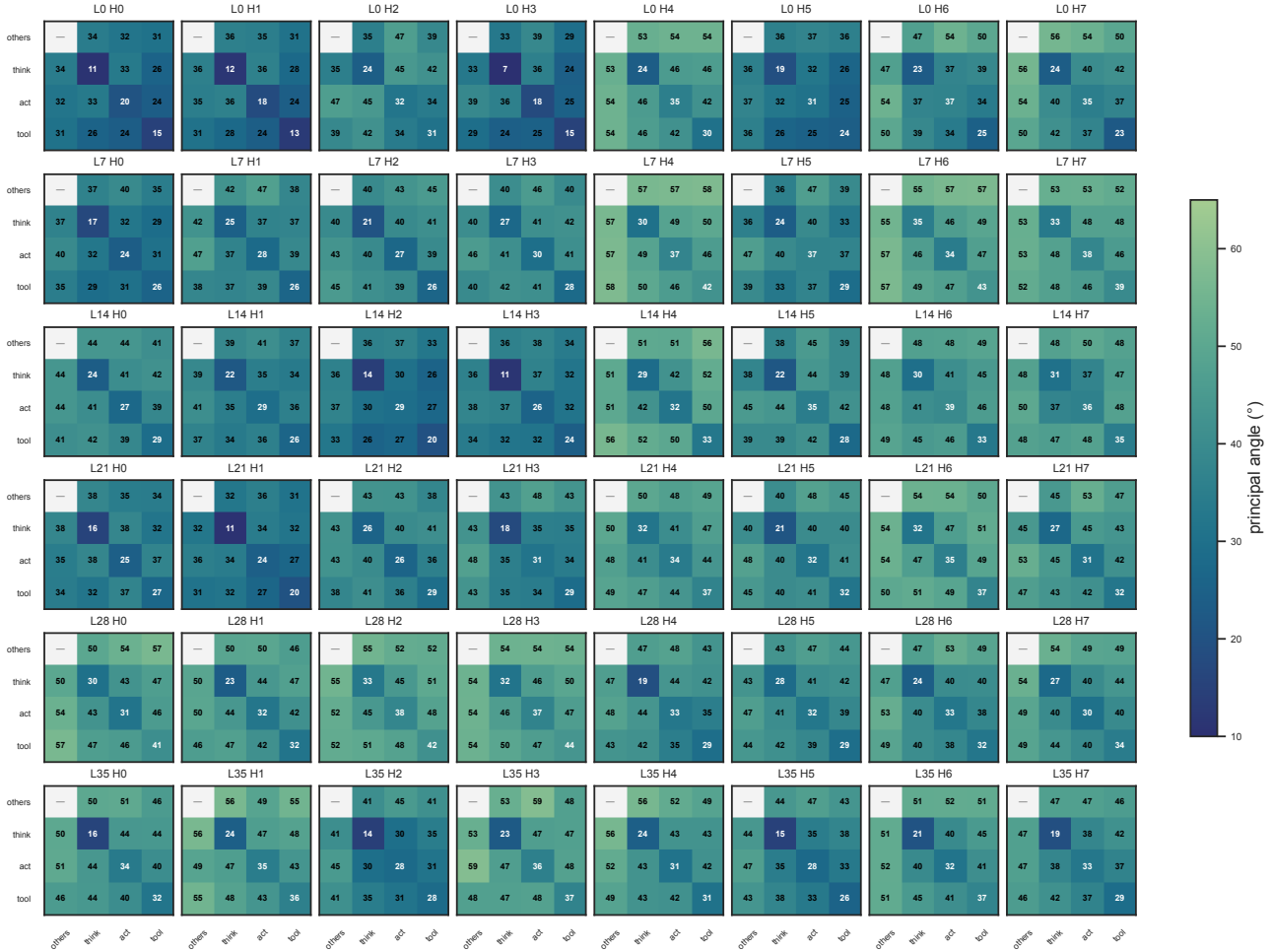


Figure 6. Per-(layer, KV head)  $4 \times 4$  principal-angle matrices, one matrix per cell. Same data and colour scale as Figure 3. Off-diagonal entries are cross-phase angles, diagonal entries are within-phase step- $i$ -vs-step- $j$  angles, and “—” denotes a within-phase pair we cannot measure (single segment, e.g. others).

## C. Utility-Coverage Diagnostics

The main text shows that query directions are more separated across agentic phases than within a phase. This appendix asks whether that geometric separation has a measurable effect on KV eviction. We use an offline utility-coverage diagnostic that compares how well different representative-query choices preserve tokens that will receive future attention.

### C.1. Setup and Metric

**Setup.** We use Qwen3-32B trajectories from BFCL `web_search_base`. All statistics are aggregated over 8 evenly spaced layers and all 8 KV heads. A compression event is triggered every 128 decode tokens. For each trajectory, we first run the full-attention model end-to-end and record all query, key, and value vectors. We then replay each compression event offline: each representative-query strategy scores the same candidate cache and selects a keep set under the same budget.

**Future-attention utility.** At a compression event, let  $i_{\text{comp}}$  denote the current position. The existing cache consists of candidate keys  $k_i$  for  $i \leq i_{\text{comp}}$ , and the future horizon consists of queries  $q_j$  for  $j > i_{\text{comp}}$ . For each candidate token  $i$ , we define its future-attention utility as

$$U_i = \sum_{j > i_{\text{comp}}} a_{j,i}, \quad a_{j,i} = \text{softmax} \left( \frac{q_j^\top K_{\leq j}}{\sqrt{d}} \right)_i.$$

This quantity measures how much attention token  $i$  receives from future queries in the full-attention trajectory. It is available only offline and is used here as an oracle diagnostic, not as a deployable eviction score.

**Utility coverage.** A representative-query strategy selects a subset of past queries, uses them to score the candidate cache, and retains the top- $B$  tokens as the keep set  $\mathcal{S}$ . For phase  $\phi$ , let  $\mathcal{K}_\phi$  be the set of candidate tokens belonging to that phase and let  $\mathcal{R}_\phi = \mathcal{S} \cap \mathcal{K}_\phi$  be the retained tokens from that phase. The phase-wise utility coverage is

$$\text{UC}_\phi = \frac{\sum_{i \in \mathcal{R}_\phi} U_i}{\sum_{i \in \mathcal{K}_\phi} U_i} \in [0, 1].$$

Higher  $\text{UC}_\phi$  means that the retained cache preserves more of the future-attended tokens from phase  $\phi$ .

### C.2. Across-Phase Coverage

Figure 7(a) compares recent-query representatives with phase-balanced representatives. The figure reports utility coverage normalized by the max-overall-utility baseline. The reported numbers are means over approximately 190 compression events.

Recent-query representatives cover adjacent `think` and `act` spans well, because the recent window is usually concentrated in those phases. However, they preserve very little `tool`-phase utility. Increasing the number of recent queries from  $q=8$  to  $q=24$  does not fix this failure mode: `tool` coverage remains nearly unchanged, moving from 0.043 to 0.042 at  $B=512$  and from 0.116 to 0.117 at  $B=1536$ . This indicates that the issue is not simply too few representatives; rather, the recent window is sampling from the wrong part of the phase mixture.

Phase-balanced representatives substantially improve coverage of non-recent phases. For `tool`, utility coverage rises from 0.042 to 0.175 at  $B=512$  and from 0.117 to 0.416 at  $B=1536$ . This supports the claim that recency-only representatives are biased toward adjacent phases, while multi-phase representatives better cover the future-query mixture.

**Interpreting the max-overall-utility baseline.** The max-overall-utility baseline ranks cached tokens directly by their realized future utility  $U_i$ . It maximizes total retained utility across all phases, but it is not a per-phase upper bound. For example, recent-query and phase-balanced strategies can exceed the max-overall-utility baseline on `think` or `act`, because the max-overall-utility reallocates more of its fixed budget to phases such as `others` and `tool`, where the absolute future utility is larger. Thus, the max-overall-utility should be interpreted as a non-deployable diagnostic reference for total utility, not as a direction-wise or phase-wise optimum.



### C.3. Within-Phase Query Selection

The across-phase results show that covering multiple phases is important. We next ask whether the choice of sampler within each phase matters once the phase budget is fixed. All methods in this ablation use  $q=24$  total representatives, with six queries selected from each phase. We use the same trajectories, compression events, and utility-coverage metric as in §C.1.

We compare four deployable within-phase samplers and one non-deployable oracle reference:

- **Last-6 per phase.** Select the last six query positions from each phase before the compression event.
- **Multi-scale.** Select six queries from each phase using a hierarchical positional sampler. The slots are distributed roughly 3:2:1 across three temporal windows: the last 4, last 32, and last 128 queries of the phase. Within each window, selected positions are evenly spaced.
- **Weighted coverage.** Assign each phase- $p$  key  $k_i$  an importance weight

$$w_i \propto \sum_{q \in \text{last-8 queries of } p} q^\top k_i.$$

This weight estimates how much the most recent queries from the same phase attend to that key. We then greedily select six anchors from the last 64 queries of the phase to maximize

$$\sum_{i \in \mathcal{K}_p} w_i \cdot \max_{a \in \text{anchors}} q_a^\top k_i.$$

This rule is Q/K-aware but does not use future utility  $U_i$ .

- **Diverse queries.** Select phase queries by greedy farthest-point sampling in cosine distance. The sampler starts from the most recent query in the phase and repeatedly adds the candidate whose maximum similarity to the selected set is smallest.
- **max-overall-utility.** Rank cached tokens directly by realized future utility  $U_i$  and keep the top- $B$  tokens. This method does not use anchor queries or phase quotas and is not deployable.

Figure 7(b) shows that the four deployable samplers have similar utility-coverage profiles. Their aggregate utility coverage differs only modestly across budgets. This suggests that, for this diagnostic, the main gain comes from allocating representatives across phases rather than from a more elaborate sampling rule within each phase. AGENTKV therefore uses the simple last- $m$  per-phase sampler as the default deployable choice.

## D. AgentKV: implementation details

This appendix gives the engineering details needed to reproduce experiments. AGENTKV is implemented as a drop-in replacement for the attention layer’s compression, the model itself is unmodified.

### D.1. Experiment setup

All end-to-end evaluations use a persistent multi-turn cache setting, so each turn continues from the KV state produced by previous turns. To implement this setting for Qwen3, we retain prior-step `<think>` blocks that the default chat template would otherwise strip when constructing the next prompt (Qwen Team, 2025). Online compression is triggered every 128 decode tokens.

We report R-KV (Cai et al., 2026) and Tri-attention (Mao et al., 2026) as recent end-to-end compressed-cache baselines. R-KV provides the SnapKV-family comparison: it retains the recency-attention term and adds a redundancy penalty. Tri-attention is a recent KV-compression method for efficient long reasoning.

We evaluate Qwen3-32B and Qwen3-8B (Qwen Team, 2025) on BFCL (Patil et al., 2025; Gorilla LLM Team, 2025) agentic tool-calling tasks, `web_search_base` (web retrieval and summarization), `multi_turn_base` (multi-step function-calling), and `memory`. These tasks cover long-context retrieval, multi-step tool execution, and external-memory use. `memory` contains multiple sub-categories; we report its overall `Memory Summary` score. We also evaluate both models on the  $\tau^2$ -bench (Yao et al., 2024) `airline`, `retail`, and `telecom` domains. The tables report the corresponding benchmark scores for each model-task-budget setting.

For serving, we compare the AGENTKV implementation against upstream SGLang 0.5.12 with FlashInfer attention. We report output-token throughput because both systems process the same prompts with the same prefill configuration, so the comparison isolates the steady-state decode path where KV-cache size directly affects serving cost. Each case issues requests with prompt/generation lengths of 16k/1k or 32k/1k on H100 NVL or RTX A6000. Both systems use `memory_ratio=0.90`, `max_prefill_length=8192`, and CUDA graph capture. We use token budget `B = 4096`.

### D.2. Phase detection

We attach a four-state machine to the tokenizer’s chat-template emission. State transitions fire on the corresponding chat-template / tool-call markers and the bracket tokens of `<think>...</think>` and `<tool_call>...</tool_call>`:

- `others`: tokens outside the `think`, `act`, and `tool` phases, including the system prompt, tool descriptions, and user-turn prefix.
- `think`: tokens inside the assistant’s `<think>...</think>` block of the current step.
- `act`: tokens inside the assistant’s `<tool_call>...</tool_call>` JSON payload of the current step.
- `tool`: tokens inside the tool’s `<|im_start|>tool...<|im_end|>` response message.

Phase labels are tracked at single-token granularity, so a single prefill segment that crosses a marker (e.g. `</think>` immediately followed by `<tool_call>`) is split correctly. The bracket markers themselves are assigned to the phase they introduce. No model-side classifier is used; the phase signal is purely tokenizer-visible.

### D.3. Hyperparameters

Following tables record the configuration used for experiments. For R-KV and Tri-attention, we follow the hyperparameters in their setup, except for R-KV we use the same observation window size  $q = 32$ .

Table 5. Default AGENTKV hyperparameters.

Hyperparameter	Value
Total anchor budget $q$	32
Per-phase ring capacity	8 ( $q/4$ )
Compression interval $\Delta$	128 tokens
Phases tracked	think / act / tool / others
KV cache budget $B$	task-dependent (see §4)

Table 6. Default R-KV hyperparameters.

Hyperparameter	Value
Total anchor budget $q$	32
Redundancy score $\lambda$	0.07
Compression interval $\Delta$	128 tokens
KV cache budget $B$	task-dependent (see §4)

Table 7. Default Tri-attention hyperparameters.

Hyperparameter	Value
Compression interval $\Delta$	128 tokens
KV cache budget $B$	task-dependent (see §4)