

# When KV Cache Reuse Fails in Multi-Agent Systems: Cross-Candidate Interaction is Crucial for LLM Judges

Sichu Liang<sup>1\*</sup>, Zhenglin Wang<sup>1\*</sup>, Jiajia Chu<sup>2</sup>, Pengfei Xia<sup>2</sup>, Hui Zang<sup>2†</sup>, Deyu Zhou<sup>1†</sup>

<sup>1</sup>School of Computer Science and Engineering, Key Laboratory of Computer Network and Information Integration, Ministry of Education, Southeast University, China

<sup>2</sup>Huawei Technologies Co., Ltd.

{coder\_liang, zhenglin, d.zhou}@seu.edu.cn

{chujiajia1, xiapengfei12, zanghui}@huawei.com

## Abstract

Multi-agent LLM systems routinely generate multiple candidate responses that are aggregated by an LLM judge. To reduce the dominant prefill cost in such pipelines, recent work advocates KV cache reuse across partially shared contexts and reports substantial speedups for generation agents. In this work, we show that these efficiency gains do not transfer uniformly to *judge-centric* inference. Across GSM8K, MMLU, and HumanEval, we find that reuse strategies that are effective for execution agents can severely perturb judge behavior: *end-task accuracy may appear stable, yet the judge’s selection becomes highly inconsistent with dense prefill*. We quantify this risk using Judge Consistency Rate (JCR) and provide diagnostics showing that reuse systematically weakens cross-candidate attention, especially for later candidate blocks. Our ablation further demonstrates that explicit cross-candidate interaction is crucial for preserving dense-prefill decisions. Overall, our results identify a previously overlooked failure mode of KV cache reuse and highlight judge-centric inference as a distinct regime that demands dedicated, risk-aware system design. <sup>1</sup>

## 1 Introduction

Modern intelligent systems often reason by considering multiple alternatives rather than committing to a single trajectory—a pattern observed in human decision-making and in AI systems that rely on search, sampling, and planning (Churchland et al., 2008; Browne et al., 2012; Silver et al., 2016). **Large language models (LLMs)** have begun to exhibit similar behavior through techniques such as multi-sample decoding (Leviathan et al., 2023; Gui et al., 2024; Sun et al., 2024), self-consistency (Wang et al., 2023; Li et al., 2025c),

\*Equal Contribution, †Corresponding Author

<sup>1</sup>The code is available in [https://github.com/dbsxfz/kv\\_reuse\\_fails](https://github.com/dbsxfz/kv_reuse_fails)

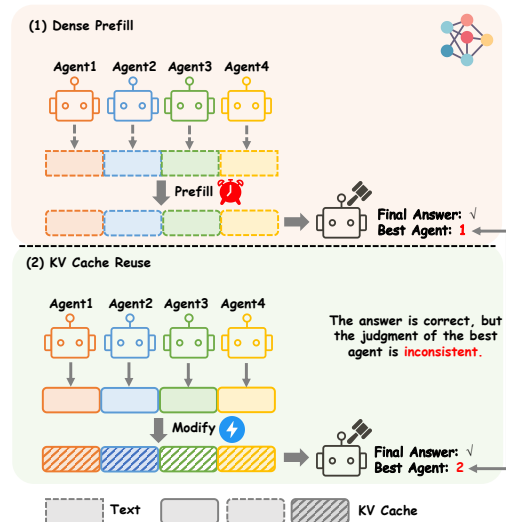


Figure 1: Illustration of decision non-invariance under judge-side KV cache reuse. **Top:** dense prefill recomputes the judge KV cache and selects Agent 1. **Bottom:** KV reuse stitches/modifies cached KV blocks, keeping the final answer correct but changing the selected best agent, despite identical candidate texts.

preference modeling (Rafailov et al., 2023; Ouyang et al., 2022; Sun et al., 2025; Ye et al., 2025b), and search-based inference (Yao et al., 2023; Besta et al., 2024), where multiple candidate solutions are generated and subsequently evaluated.

LLM-based **multi-agent systems (MAS)** (Hong et al., 2024; Wu et al., 2024) represent a concrete instantiation of this paradigm, where multiple execution agents generate diverse candidate outputs that are aggregated by a central meta-agent (Zhang et al., 2025c) responsible for planning, coordination or **judging** (Chen et al., 2024; Fourney et al., 2024; Zhang et al., 2025b). To enable such collaboration, execution agents and the meta-agent must repeatedly exchange intermediate reasoning, candidate outputs, and contextual information, leading to substantial overlap in their effective input contexts (Wang et al., 2025; Zhang et al., 2025a).

While this interaction improves robustness and reasoning quality, it also introduces severe computational redundancy: each agent repeatedly performs a full prefill pass over largely overlapping context, causing inference cost to scale rapidly with the number of agents (Ye et al., 2025a).

A natural solution is **KV cache reuse** (Kwon et al., 2023; Zheng et al., 2024b). During LLM inference, the prefill stage encodes the entire input context and constructs key–value (KV) caches for attention, which can be reused directly when multiple inputs share overlapping prefixes. However, in multi-agent settings, simple prefix-based reuse is often ineffective: although agents typically share substantial contextual overlap (e.g., task instructions or exchanged messages), their prompts frequently diverge due to role specifications, intermediate reasoning, or agent-specific context. Consequently, KV caches across agents are similar but not identical, preventing direct reuse (Zhao et al., 2025; Liu et al., 2025). This has motivated *cross-prefix* (approximate) reuse methods that enable reuse under partially shared prefixes and have reported substantial speedups for generation agents in multi-agent pipelines and serving systems (Yao et al., 2025; Ye et al., 2025a; Yang et al., 2025b).

Most existing studies evaluate KV reuse primarily from the perspective of execution agents, implicitly assuming that efficiency gains transfer uniformly across all components of a multi-agent pipeline. Yet the impact of KV reuse on *judge* agents—which must compare multiple candidates jointly and produce a selection (Li et al., 2024a)—remains largely unexplored.

This gap is non-trivial, as judging constitutes a qualitatively different inference regime from generation (Li et al., 2025b). Prior work has examined the reliability of LLM-based judges (Zheng et al., 2024a), documenting systematic biases such as position preference (Li et al., 2024b) and sensitivity to irrelevant contexts (Shi et al., 2023). However, these failures are typically attributed to model limitations or task-level biases, rather than to *system-level efficiency optimizations* that modify the inference procedure. Consequently, how KV cache reuse alters judge behavior in multi-candidate comparison settings remains an open question. As shown in Figure 1, KV reuse may keep the final answer unchanged while altering which candidate the judge selects, motivating judge-centric diagnostics beyond accuracy.

In this paper, we revisit KV cache reuse from

a *judge-centric* perspective, where an LLM judge must *jointly* compare multiple candidates within a single context and output both a final answer and the selected candidate. We show that judge-side reuse can be *decision-non-invariant*: it may preserve end-task accuracy while substantially changing which candidate the judge selects, especially under order perturbations. These results highlight judge-centric inference as a distinct regime in which preserving cross-candidate interactions is critical.

Our key contributions are:

- **Judge-centric evaluation of KV reuse.** We provide a systematic study of KV reuse in multi-candidate judging across tasks, candidate-generation regimes, and candidate orderings.
- **Decision instability beyond accuracy.** We introduce Judge Consistency Rate (JCR) to quantify decision non-invariance relative to dense prefill, and show that Acc. and JCR can be decoupled—accuracy may appear stable while selection behavior becomes inconsistent under reuse.
- **Mechanistic diagnostics and implications.** Through attention analyses and controlled ablations, we attribute the failure to disrupted cross-candidate interaction; we further develop PAL-KV as a probe to rule out agent identity as the dominant bottleneck and discuss practical directions such as interaction-aware reuse and risk-aware gating.

## 2 Related Work

**KV Cache Reuse.** KV cache reuse is a widely studied approach for reducing the high prefill cost of large language model (LLM) inference. Early systems (Kwon et al., 2023; Zheng et al., 2024b) focus on exact-prefix caching, where KV tensors are reused only when the input context remains identical. Recent work extends KV reuse beyond exact prefix matching. One line studies *cross-model* cache sharing, enabling cache reuse across model variants (Liu et al., 2026; Li et al., 2025a) or even heterogeneous model families (Fu et al., 2026). More closely related to our study is *cross-prefix* KV reuse, which enables approximate reuse when prompts share partial prefixes but are not identical. Representative methods explore reuse through simple position adjustment (Lu et al., 2025), linking

tokens (Yang et al., 2025c), or reuse-time repair/fusion mechanisms (Yao et al., 2025; Hu et al., 2025; Yang et al., 2025a; Cao et al., 2026), and are typically evaluated in settings where cross-segment interactions are limited. *KVCOMM* (Ye et al., 2025a) further brings cross-prefix reuse to *multi-agent systems* by using anchor examples to correct cache deviations across agent-specific prompts.

We study KV reuse in *judge-centric* multi-candidate inference and show that reuse strategies effective for agent-side generation can be decision-non-invariant for judges, revealing a failure mode overlooked by prior work.

**LLM-Based Multi-Agent Systems.** LLM-based multi-agent systems coordinate multiple specialized agents through structured communication graphs (Chen et al., 2024; Zhuge et al., 2024), where *execution agents* produce candidate outputs under different roles or prompts (Wang et al., 2025; Zhang et al., 2025a). A common pattern is a central LLM *meta-agent* that aggregates and verifies multiple agent responses within a single context window to produce the final decision (Li et al., 2023; Hong et al., 2024; Zhang et al., 2025c), making it a key component for both quality and efficiency.

We show that KV reuse, while effective for execution agents, is not always *behavior-preserving* for judge agents, often changing which candidate is selected relative to dense prefill.

**LLM-as-a-Judge.** LLM-based judges are increasingly deployed as *decision modules* in high-stakes workflows including healthcare, finance, and legal analysis (Xie et al., 2024, 2023; Raju et al., 2024), where selection robustness is critical. Most existing formulations treat judging as either *pairwise comparison* or *scalar scoring*, where candidates are evaluated independently or in pairs (Sun et al., 2024, 2025). In contrast, many practical applications require a judge to process *multiple candidates jointly* within a single context and directly select the best output (Tran et al., 2025; Gera et al., 2025). This setting is especially common in multi-agent systems, where a judge must aggregate and attribute outputs from several execution agents.

Such judges require cross-candidate comparison and relative reasoning, unlike standard generation (Li et al., 2025b). We study KV cache reuse under this joint multi-candidate judging setting and show it can change the judge’s selection relative to dense prefill.

### 3 Preliminaries

We study *judge-centric* multi-candidate inference with  $N$  execution agents and a central judge agent. Given an input question  $x$ , each execution agent  $A_i$  is instantiated as a large language model with an agent-specific prompt  $p_i$  (e.g., role description or instruction), and produces a candidate response  $y_i$ . A judge agent  $J$  with a judge-specific prompt  $p_J$  then takes all candidate responses as input and outputs (i) a final answer  $\hat{y}$  and (ii) the index of the selected candidate  $\hat{i}$ :

$$(\hat{y}, \hat{i}) = J(p_J, x, y_{1:N}), \quad \hat{i} \in \{1, \dots, N\}. \quad (1)$$

Unlike pairwise comparison or scalar scoring, the judge must *jointly* reason over multiple candidates within a single context, performing explicit cross-candidate comparison and selection.

**Candidate generation regimes.** We consider two common ways to construct multi-candidate inputs in multi-agent pipelines.

**(1) Progressive refinement.** Candidates are generated sequentially, where later agents may condition on earlier candidates:

$$y_i = A_i(p_i, x, y_{1:i-1}), \quad i = 1, \dots, N. \quad (2)$$

**(2) Parallel exploration.** Candidates are generated independently without observing each other:

$$y_i = A_i(p_i, x), \quad i = 1, \dots, N. \quad (3)$$

**Judge input ordering.** To disentangle content-based judging from potential ordering effects, we optionally randomize the presentation order of candidates to the judge. Let  $\pi$  be a permutation over  $\{1, \dots, N\}$ . The judge observes candidates in the permuted order  $y_{\pi(1:N)}$ :

$$(\hat{y}, \hat{i}) = J(p_J, x, y_{\pi(1:N)}), \quad (4)$$

where no-shuffle uses the identity permutation and shuffle samples  $\pi$  per example. Under shuffle, the predicted index  $\hat{i}$  is mapped back to the original candidate index via  $\pi$ .

### 4 KV Reuse for Multi-Candidate Judging

Figure 2 summarizes the judge-side cache construction strategies compared in this work. We consider judge-centric inference where a judge processes  $N$  candidate responses within a single context. We

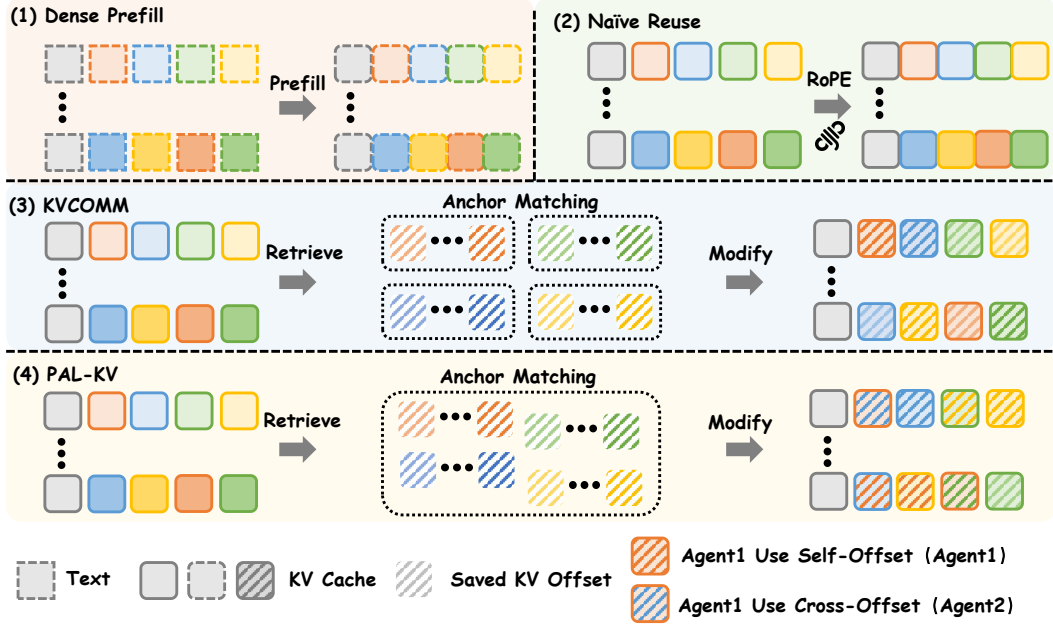


Figure 2: Judge-side KV cache construction for multi-candidate judging. **Dense prefill** recomputes the full judge cache, while **Naïve Reuse** aligns and stitches execution-side candidate KV chunks. **KVCOMM** retrieves anchor-based cache offsets to correct reused chunks, and **PAL-KV** pools anchors across agents for offset retrieval.

use  $\mathcal{K}(\cdot)$  to denote the KV cache obtained by encoding a token sequence, and  $\mathcal{K}(s)[u]$  to denote the KV entries in  $\mathcal{K}(s)$  corresponding to a subsequence  $u \subseteq s$ . For each candidate response  $y_i$ , the execution agent has already computed and cached its KV entries under its execution-time context  $S_i^{\text{exec}}$ :

$$\mathcal{C}_i = \mathcal{K}(S_i^{\text{exec}})[y_i]. \quad (5)$$

We use  $\oplus$  to denote concatenation of KV-cache chunks in the displayed candidate order (i.e., stitching cache segments into a single judge cache). Let  $o_i$  be the starting *token position* (index) of candidate  $y_i$  in the judge input, and let  $\mathcal{C}_i^{\rightarrow o_i}$  denote the position-aligned version of  $\mathcal{C}_i$  (e.g., via RoPE (Su et al., 2024) re-indexing) so that it is valid starting at position  $o_i$  in the judge sequence. For a KV cache reuse method  $m$ , we denote by  $\mathcal{C}_m^J$  the assembled KV cache segment corresponding to the concatenation of the  $N$  candidate responses in the judge prompt.

**Dense Prefill.** Let  $S^J = (p_J, x, y_1, \dots, y_N)$  denote the full judge input sequence. Dense prefill recomputes the judge KV cache from scratch as  $\mathcal{K}(S^J)$ . In particular, the KV chunk for candidate  $y_i$  under dense judging is

$$\mathcal{C}_i^{\text{dense}} = \mathcal{K}(S^J)[y_i], \quad (6)$$

which differs from the execution-time chunk  $\mathcal{C}_i$  due to the different conditioning context.

Then the full judge cache can be written as

$$\mathcal{C}_{\text{dense}}^J = \mathcal{C}_1^{\text{dense}} \oplus \mathcal{C}_2^{\text{dense}} \oplus \dots \oplus \mathcal{C}_N^{\text{dense}}. \quad (7)$$

**Naïve Reuse (Position-Only Reuse).** Naïve reuse stitches cached response chunks into the judge context by position alignment and concatenation:

$$\mathcal{C}_{\text{naive}}^J = \mathcal{C}_1^{\rightarrow o_1} \oplus \mathcal{C}_2^{\rightarrow o_2} \oplus \dots \oplus \mathcal{C}_N^{\rightarrow o_N}. \quad (8)$$

This approach reuses execution-side KV as-is and does not account for prefix-dependent KV deviations induced by the judge-side context. This position-only stitching strategy has been used in prior KV reuse settings, e.g., in RAG-style cache reuse (Lu et al., 2025) and in multi-agent interaction setups that reuse cached internal states (Zou et al., 2025).

**KVCOMM (Anchor-based Offset Correction).** KVCOMM improves reuse by adding an anchor-retrieved correction to each reused chunk (Ye et al., 2025a). For candidate  $i$ , KVCOMM constructs a matching view  $v_i$  under the current judge-side prefix and retrieves a correction

$$\hat{\Delta}_i = \text{RetrieveOffset}(v_i; \mathcal{A}^{(i)}), \quad (9)$$

where  $\mathcal{A}^{(i)}$  is the agent-specific anchor pool used by KVCOMM. It then corrects the aligned chunk as

$$\tilde{\mathcal{C}}_i^J = \mathcal{C}_i^{\rightarrow o_i} + \hat{\Delta}_i, \quad (10)$$

and assembles the judge cache by concatenation:

$$C_{\text{kvcmm}}^J = \tilde{C}_1^J \oplus \tilde{C}_2^J \oplus \dots \oplus \tilde{C}_N^J. \quad (11)$$

We follow KVCMM for anchor construction, matching, and correction estimation, and use its reliability criterion to fall back to dense computation when reuse is deemed unreliable.

**PAL-KV (Pooled-Anchor Lookup).** PAL-KV (Pooled-Anchor Lookup KV reuse) is a minimal modification of KVCMM that changes only the retrieval scope of anchors. While KVCMM retrieves  $\hat{\Delta}_i$  from an agent-specific pool  $\mathcal{A}^{(i)}$ , PAL-KV retrieves from the union of all pools:

$$\hat{\Delta}_i^{\text{pal}} = \text{RetrieveOffset} \left( v_i; \bigcup_{j=1}^N \mathcal{A}^{(j)} \right). \quad (12)$$

and applies

$$\tilde{C}_{i,\text{pal}}^J = C_i^{\rightarrow o_i} + \hat{\Delta}_i^{\text{pal}}. \quad (13)$$

The assembled judge cache is

$$C_{\text{pal}}^J = \tilde{C}_{1,\text{pal}}^J \oplus \tilde{C}_{2,\text{pal}}^J \oplus \dots \oplus \tilde{C}_{N,\text{pal}}^J. \quad (14)$$

Thus, PAL-KV modifies *which* correction is selected, while keeping the reuse and correction mechanism identical to KVCMM.

## 5 Experiments

### 5.1 Experimental Setup

**Multi-candidate generation.** We use the two candidate-generation regimes in Sec. 3, and generate  $N=4$  candidates per example. To isolate judge-side effects, we disable execution-side reuse and *fix the candidate set*: for each example we generate the candidates once with dense prefill and reuse the identical candidate texts for all judge-side methods. We evaluate both no-shuffle and shuffle candidate orders to quantify ordering effects on judge-side reuse.

**Models and implementation.** Our primary experiments use Llama-3.2-3B-Instruct (Meta, 2024) for both execution agents and the judge, while ablations consider models from the Llama-3.1/3.2 (Grattafiori et al., 2024) and Qwen-2.5 (Yang et al., 2024) families across a range of sizes (3B–72B). To induce candidate diversity while keeping the judge deterministic, execution agents use temperature 0.2 and the judge uses temperature 0. For KVCMM and PAL-KV, we use an anchor pool of size  $|\mathcal{A}^{(i)}| = 5$  per agent.

**Hardware and precision.** Our main experiments are run on 2 NVIDIA RTX 4090 GPUs (24GB each). For the larger-scale model extension, we use 8 Ascend 910B NPUs with 64GB memory per card. For both execution agents and the judge, the maximum generation length is set to 512 tokens. All models are run in bfloat16 for numerical stability.

### 5.2 Benchmarks

We evaluate judge-centric KV reuse on three representative benchmarks spanning reasoning, programming, and knowledge-intensive domains: **GSM8K.** A grade-school math reasoning benchmark (Cobbe et al., 2021); **HumanEval.** A code generation benchmark consisting of programming tasks with unit tests (Chen, 2021). **MMLU.** A multi-domain knowledge benchmark covering diverse subjects (Hendrycks et al., 2021).

### 5.3 Evaluation Metrics

**Task Acc.** The primary metric for the judge’s final answer: exact-match accuracy for GSM8K and MMLU, and Pass@1 for HumanEval (also denoted as *Acc.* in tables).

**Judge Consistency Rate (JCR).** The percentage of examples for which KV reuse and dense-prefill judging select the same candidate under the same candidate set and presentation order. JCR measures whether reuse preserves the dense-prefill judge’s selection behavior, rather than merely the correctness of the final answer. For shuffle, we use identical permutations for both runs and map predicted indices back to the original candidate IDs.

**Reuse Rate (Reuse).** The fraction of candidate blocks ( $y_{1:N}$ ) assembled via reuse versus dense recomputation. It excludes the shared judge prefix and output tokens. We disable execution-side reuse to ensure an identical candidate set across methods.

## 6 Results and Analysis

### 6.1 Main Results

**Overview: judge decisions are not invariant under KV reuse.** Table 1 reports task accuracy (Acc.), Judge Consistency Rate (JCR), and the judge-side candidate Reuse Rate (Reuse) across all benchmarks. A central observation is that **KV reuse can substantially change which candidate the judge selects, even when end-task accuracy remains comparable.** This decision non-invariance appears consistently under both multi-

Table 1: Judge-side KV cache reuse results on three benchmarks under different settings. Bold denotes the best Acc., and the JCR column additionally reports the performance drop under shuffle relative to no-shuffle.

Settings	Shuffle	Method	MMLU			GSM8K			HumanEval		
			Acc.	JCR	Reuse	Acc.	JCR	Reuse	Acc.	JCR	Reuse
Progressive Refinement	No	Dense Prefill	45.09	100.00	00.00	<b>71.34</b>	100.00	00.00	33.54	100.00	00.00
		Naïve Reuse	<b>49.02</b>	31.37	100.00	11.83	50.87	100.00	45.96	64.60	100.00
		KVCOMM	39.22	66.67	32.55	64.14	87.91	29.93	<b>51.55</b>	57.14	48.32
		PAL-KV	42.48	69.28	32.55	64.59	89.23	29.93	49.69	66.46	48.32
	Yes	Dense Prefill	44.44	100.00(+00.00)	00.00	<b>70.51</b>	100.00(+00.00)	00.00	31.68	100.00(+00.00)	00.00
		Naïve Reuse	39.86	27.45(-03.92)	100.00	16.98	19.88(-30.99)	100.00	49.06	23.60(-41.00)	100.00
		KVCOMM	39.87	39.87(-26.80)	32.29	63.84	51.75(-36.16)	31.86	<b>49.69</b>	21.74(-35.40)	44.84
		PAL-KV	<b>47.06</b>	31.37(-37.91)	32.29	62.53	51.27(-37.96)	31.86	49.07	24.22(-42.24)	44.84
Parallel Exploration	No	Dense Prefill	49.02	100.00	00.00	<b>67.32</b>	100.00	00.00	28.57	100.00	00.00
		Naïve Reuse	<b>61.44</b>	26.14	100.00	20.34	23.93	100.00	<b>40.99</b>	39.75	100.00
		KVCOMM	47.06	58.16	44.84	39.63	59.19	63.12	31.68	61.49	66.96
		PAL-KV	45.10	62.75	44.84	41.17	59.55	63.12	32.30	64.60	66.96
	Yes	Dense Prefill	<b>51.63</b>	100.00(+00.00)	00.00	<b>70.28</b>	100.00(+00.00)	00.00	25.63	100.00(+00.00)	00.00
		Naïve Reuse	36.60	24.18(-1.96)	100.00	19.86	24.92(+0.99)	100.00	<b>50.31</b>	27.95(-11.80)	100.00
		KVCOMM	43.14	46.41(-11.75)	44.84	38.89	34.43(-24.76)	63.21	37.27	32.92(-28.57)	66.96
		PAL-KV	43.14	45.75(-17.00)	44.84	39.95	35.32(-24.23)	63.21	36.02	36.65(-27.95)	66.96

agent candidate-generation regimes (Progressive Refinement and Parallel Exploration), suggesting that the phenomenon is not tied to a particular collaboration pattern. Overall, these results expose a failure mode that would be missed if one monitors task metrics alone: KV reuse may preserve outcomes while silently perturbing the judge’s underlying selection behavior.

**Accuracy can mask instability: reuse perturbs selection beyond what Acc. reveals.** Under dense prefill, changing the candidate order (shuffle vs. no-shuffle) often affects Acc. only moderately. In contrast, reuse-based methods frequently exhibit **low JCR** even when Acc. remains close to the dense-prefill baseline, indicating that the judge’s cross-candidate comparison is highly sensitive to judge-side reuse. This decoupling is critical for judge-centric pipelines where downstream properties—such as attribution, explanation, and auditability—depend on *which* candidate is chosen, not only whether the final answer is correct. Appendix A.2 further analyzes this phenomenon and shows that Acc. and JCR are weakly correlated.

**Candidate-order perturbation (shuffle) amplifies inconsistency.** A striking pattern is that shuffle sharply reduces JCR for *all* reuse-based methods. This indicates that judge-side reuse is strongly *layout-dependent*: the effective prefix for each candidate block is determined by the entire preceding candidate configuration in the judge prompt. When the order changes, reuse approxima-

Table 2: Effect of masking cross-candidate attention under dense prefill.

Settings	Shuffle	Acc. (%)		JCR (%)	
		Original	Masked	Original	Masked
Progressive Refinement	No	45.09	43.79	100.00	28.76
	Yes	44.44	45.10	100.00	32.03
Parallel Exploration	No	49.02	48.37	100.00	22.22
	Yes	51.63	47.71	100.00	31.37

tions become much less reliable, leading to large drops in JCR even when Acc. changes little.

**PAL-KV as a probe: agent identity is not the dominant bottleneck.** We use PAL-KV as a controlled variant to test whether judge inconsistency is mainly caused by *agent-specific* offset distributions assumed by KVCOMM. By pooling anchors across agents, PAL-KV relaxes the agent-identity constraint in offset retrieval. Empirically, pooled retrieval yields *small but consistent* gains in the fixed-layout setting (no-shuffle), suggesting that some offset patterns transfer across agents when the presentation structure is stable. However, under shuffle, PAL-KV does not mitigate the large JCR drop and behaves similarly to KVCOMM; other probes show the same limitation (Appendix A.4). Together, these results suggest that **the dominant driver of judge-centric failures is the changing cross-candidate context/layout, rather than agent identity per se**. In the next subsection, we provide direct evidence that preserving *cross-candidate interaction* is crucial for maintaining dense-prefill judge decisions.

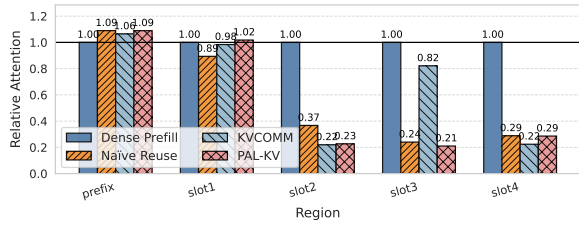


Figure 3: Relative attention mass over regions (prefix and candidate slots) under different KV reuse methods.

## 6.2 Why Does Judge-Centric Reuse Fail? Interaction Diagnostics

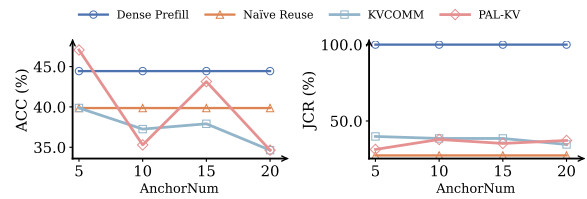
### Attention diagnostics: reuse weakens cross-candidate attention, especially for later slots.

To understand why judge selections change under reuse, we analyze attention patterns during the judge’s first-token generation. Figure 3 shows a consistent trend: attention over the shared prefix and the first candidate block is relatively similar across methods, but **attention to later candidate blocks is substantially weaker and more erratic under reuse-based methods** than under dense prefill. Representative attention maps for individual methods are provided in Appendix A.5. This is consistent with the interpretation that reuse perturbs the fine-grained cross-candidate interactions needed for joint comparison, causing the judge to under-attend to late-arriving evidence.

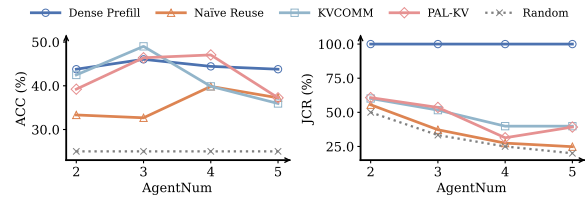
**Masking ablation: cross-candidate interaction is necessary beyond judge-style formatting.** A natural alternative explanation is that reuse fails simply because execution-side KV caches are “incompatible” with the judge prompt. To isolate the role of *interaction* from *format*, we conduct a masking ablation where candidates are placed in the judge prompt but cross-candidate attention is explicitly blocked (candidate  $i$  cannot attend to candidates  $1:i-1$ ). As shown in Table 2, masking has only a limited effect on *Acc* but causes a dramatic collapse in *JCR*, approaching random-choice behavior for  $N=4$ . This provides direct evidence that **judge-centric inference critically relies on cross-candidate interaction**, and that removing such interaction—even under the same judge-style prompt—destroys decision invariance relative to dense prefill.

### 6.3 Ablation Studies and Additional Baselines

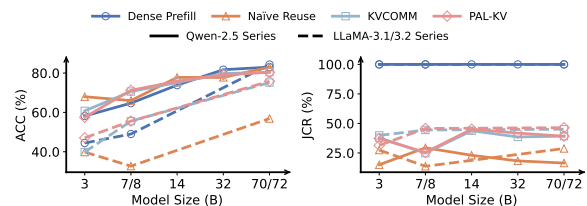
Unless otherwise specified, all ablations focus on the shuffle regime, which is particularly challenging for judge-side reuse due to order-induced lay-



(a) Effect of anchor pool size.



(b) Effect of the number of candidates (agents).



(c) Effect of model family and size.

Figure 4: Ablations under shuffle: judge-side decision non-invariance persists. Varying anchor pool size, candidate count, or model size does not reliably restore JCR for KV reuse methods.

out changes; we use **MMLU under Progressive Refinement with shuffled candidate order** as a representative testbed.

**Anchor pool size is not the bottleneck under shuffle.** A plausible hypothesis is that poor anchor matches cause offset correction to fail. We therefore increase the anchor pool size (Figure 4a); as expected, a larger pool increases the Reuse Rate (Appendix Table 6). Nevertheless, JCR under shuffle stays low and changes little as the pool grows, suggesting that the failure is **not** driven by insufficient anchor coverage. Instead, when the preceding candidate configuration changes, anchor similarity becomes an unreliable proxy for preserving the interaction effects needed for judging.

**More candidates can improve Acc but make decision invariance harder.** We vary the number of candidates  $N$  (Figure 4b). Increasing  $N$  often improves *Acc* by providing stronger candidates, but **JCR decreases monotonically** across reuse-based methods. This highlights a key tension in judge-centric systems: adding candidates improves solution quality but expands the space of cross-candidate interactions that reuse must preserve,

Table 3: Additional related baselines on MMLU under the shuffle judge setting. Despite additional local repair mechanisms, all three baselines remain low JCR across both candidate-generation regimes.

Settings	Method	Acc. (%)	JCR (%)
Progressive Refinement	Dense Prefill	44.44	100.00
	EPIC	51.63	28.10
	CacheClip	48.37	32.68
	SamKV	45.75	30.06
Parallel Exploration	Dense Prefill	51.63	100.00
	EPIC	49.02	33.33
	CacheClip	47.71	34.21
	SamKV	50.98	28.76

making invariant selection increasingly fragile.

**Scaling model size improves Acc but does not reliably restore JCR.** We evaluate different model families and sizes (Figure 4c). While larger models tend to improve *Acc*, **JCR does not show a corresponding improvement** and remains unstable under reuse. This indicates that judge inconsistency is not merely a capacity issue; rather, it reflects sensitivity to interaction patterns that current reuse schemes fail to preserve.

**Additional related baselines also fail to restore judge-side invariance.** We further adapt three representative KV cache reuse baselines originally proposed for Retrieval Augmented Generation (RAG)—EPIC (Hu et al., 2025), CacheClip (Yang et al., 2025a), and SamKV (Cao et al., 2026)—to our judge-side setting. Table 3 reports their results on MMLU under the representative shuffle setting used in this subsection. Although these methods introduce local repair mechanisms, their JCR remains far below dense prefill across both Progressive Refinement and Parallel Exploration, with the best JCR reaching only 34.21. A plausible explanation is that chunk-wise token selection, sparse candidate views, or partial recomputation tends to preserve *problem-solving salience* within each candidate, but may still miss the *fine-grained comparative cues* that determine which candidate wins under cross-candidate judging. Full implementation details and complete results on GSM8K and HumanEval are reported in Appendix A.6.

## 7 Discussion and Future Directions

Our study highlights a judge-centric failure mode of KV cache reuse: when a judge must *jointly* compare multiple candidate blocks, reuse can sub-

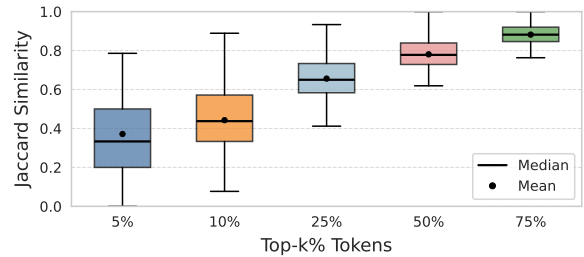


Figure 5: Jaccard similarity of selected Top- $k$ % tokens between the small and large models.

stantially alter the selected candidate (low JCR), especially under unstable layouts such as shuffle. Together with attention diagnostics and the masking ablation, the evidence points to a central mechanism: **judge-centric inference relies on fine-grained cross-candidate interactions that current reuse schemes do not reliably preserve.**

A straightforward workaround is to always run dense prefill for the judge while keeping reuse for other components (e.g., execution agents or retrieval modules). However, this can substantially increase end-to-end latency, undermining the original systems objective of KV cache reuse: balancing behavior preservation and efficiency. We therefore do not view this as a satisfying solution. Instead, we discuss two preliminary directions suggested by the mechanism above. Rather than presenting complete solutions, we use them as initial evidence for how future judge-centric systems might better balance efficiency and decision invariance.

**From the perspective of interaction-preserving reuse.** A promising direction is to enhance judge-side reuse by selectively retaining interaction-relevant context. Removing low-salience tokens (e.g., redundant rationales) may mitigate attention dilution (Shi et al., 2023), helping the judge focus on informative evidence across candidates. Crucially, this aims to concentrate rather than restrict cross-candidate visibility. In a pilot analysis, we find substantial overlap in token selections between Llama-3.2-1B and 3B models, with Jaccard similarity increasing alongside the retention budget (Figure 5). This suggests that small-to-large model cooperation (Zhao et al., 2025; Liu et al., 2025) could help identify judge-relevant tokens. However, fully restoring judge invariance likely requires more sophisticated, interaction-aware objectives beyond simple token-level overlap.

We further conduct a pilot experiment along this direction. We identify interaction-critical to-

Table 4: Pilot results for Direction (i) (interaction-aware selective recomputation) on the representative shuffled MMLU setting.

Method	Original		+Dir. (i)	
	Acc. (%)	JCR (%)	Acc. (%)	JCR (%)
Naïve Reuse	39.87	27.45	37.90	28.76
KVCOMM	39.87	39.87	41.83	47.05
PAL-KV	47.06	31.37	51.63	45.75

kens globally over the full judge input using the same model, rather than selecting tokens separately within each chunk using an auxiliary model (Cao et al., 2026). As shown in Table 4, this Direction (i) intervention improves both Acc. and JCR over the original methods, but the gains remain limited and come with extra latency. These results suggest that even when token selection is performed globally without prioritizing efficiency, restoring judge-side invariance remains difficult: the key challenge lies in recovering the interaction structure that emerges over the full multi-candidate context.

**Future Direction (i):** *Interaction-preserving reuse via selective context retention.*

**From the perspective of meta-reasoning for risk-aware gating.** Given the judge’s sensitivity to KV reuse, a promising direction is to apply *meta-reasoning*—deciding which inference strategy to employ based on the specific situation (Yan et al., 2025). While proactive segment prediction is effective in standard KV systems (Pan et al., 2025), judge-centric reuse requires conservative gating with reliable fallbacks. Rather than predicting the success of a specific strategy ( $AUC \approx 0.6$ ), a more robust meta-goal is to identify instances that are *universally safe* for reuse (i.e., unlikely to alter the judge’s selection across all three reuse methods studied here). Using features derived from the candidate set—such as correctness signals and pairwise similarities—a simple classifier can distinguish “universally safe” vs. “unsafe” instances with an AUC of 0.82 (Figure 6). This suggests that future gating mechanisms should prioritize conservative safety detection over fine-grained strategy selection to ensure decision invariance.

We also test this direction with a pilot risk-aware gating policy using the proposed safe-instance classifier. As shown in Table 5, Direction (ii) yields a much clearer JCR improvement, but only at the cost of a sharply reduced reuse rate, in some cases

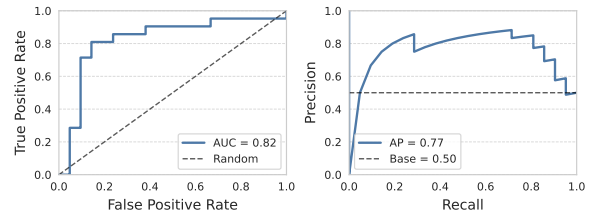


Figure 6: Detecting universally safe vs. unsafe instances for judge-side reuse ( $AUC \approx 0.82$ ,  $AP \approx 0.77$ ).

Table 5: Pilot results for Direction (ii) (risk-aware gating) on the representative shuffled MMLU setting.

Method	Original		+Dir. (ii)	
	JCR (%)	Reuse (%)	JCR (%)	Reuse (%)
Naïve Reuse	27.45	100.00	85.62	16.99
KVCOMM	39.87	32.29	88.88	5.07
PAL-KV	31.37	32.29	87.58	5.07

approaching dense-prefill usage in practice. This trade-off indicates that conservative gating can recover decision invariance more effectively than local repair alone, but does so by sacrificing much of the efficiency benefit that motivates reuse in the first place.

**Future Direction (ii):** *Meta-reasoning and risk-aware fallback gating for judge-side reuse.*

**Outlook.** Overall, our pilot results suggest that mitigation is possible but non-trivial: interaction restoration improves invariance only partially and introduces additional cost, while conservative gating improves invariance by sacrificing reuse. Accelerating judge-centric inference will therefore likely require more careful judge-aware designs that explicitly preserve cross-candidate interactions under variable layouts, potentially combining interaction-aware reuse, selective retention/compression, and principled fallback policies.

## 8 Conclusion

Our results reveal that KV reuse in judge-centric tasks is not inherently behavior-preserving. While task accuracy remains stable, JCR exposes a critical decoupling: judges frequently switch selections due to disrupted cross-candidate interactions, especially under variable layouts. These findings establish JCR as a vital diagnostic and motivate future research into interaction-preserving and risk-aware KV acceleration.

## Limitations

### Judge biases and reliability are not the focus.

LLM-based judges are known to exhibit intrinsic biases and instability (e.g., position preference and sensitivity to formatting). This work does not aim to assess the judge’s *absolute* fairness, calibration, or stability, nor do we analyze whether KV reuse amplifies specific biases of LLM judges. Instead, we focus on a narrower question: whether judge-side KV reuse is *behavior-preserving* relative to dense prefill in multi-candidate *joint* judging, i.e., whether reuse introduces decision non-invariance even when task accuracy appears stable.

**Scope of reuse and system coverage.** Our study focuses on *judge-centric* KV reuse—reusing KV chunks of candidate blocks when constructing the judge’s prefill cache—and uses dense prefill as the reference behavior. We do not claim comprehensive end-to-end optimization for all components of a full multi-agent stack, nor do we report system-wide latency/throughput under different serving configurations. Consequently, our reported Reuse Rate and analyses are intentionally centered on the judge stage rather than end-to-end efficiency.

**Protocol and prompt dependence.** Our findings are obtained under a specific judging protocol (joint multi-candidate selection with structured outputs) and candidate-generation settings (few-shot prompting and chain-of-thought style rationales). Different judge prompts (e.g., pairwise ranking vs. joint selection), output constraints (selection-only vs. judge rewriting), or candidate formatting may change the sensitivity profile. While we expect the core phenomenon—that joint judging relies on delicate cross-candidate interactions—to persist, quantitative results may vary with protocol choices.

**Limited cross-model and cross-architecture reuse.** We do not systematically evaluate judge-centric KV reuse across heterogeneous model families, sizes, or architectures (e.g., a stronger judge supervising weaker generators, or mixed-family agent systems). Recent work has begun to explore cross-model KV communication or sharing beyond same-backbone variants (Fu et al., 2026; Liu et al., 2026; Li et al., 2025a), but these methods target multi-model communication or model-sharing settings rather than judge-centric multi-candidate comparison. Extending judge-side reuse to heterogeneous model mixtures, and understanding whether

reuse remains behavior-preserving in such settings, remains an important direction for future work.

**Model scale and architecture coverage.** Although we extend our evaluation to larger dense models, we do not claim that the magnitude of decision non-invariance will transfer unchanged to larger frontier-scale models or to Mixture-of-Experts (MoE) architectures. In particular, MoE inference introduces additional complexity in routing and cache management. Evaluating judge-centric KV reuse under larger-scale and MoE settings remains important future work.

**Implementation stack and deployment variability.** All experiments are conducted under an open-source inference stack (HuggingFace Transformers<sup>2</sup>) on a fixed GPU setup. Different kernels, serving systems, attention implementations, or cache-management policies could affect absolute accuracy and the magnitude of reuse-induced perturbations. We encourage validation under diverse deployment stacks.

**Downstream safety and compute considerations.** Our study highlights a failure mode where accuracy may mask decision instability; in high-stakes downstream workflows, such non-invariance can complicate accountability and auditing. We do not provide application-specific deployment guidelines, and practitioners should perform domain-specific evaluation and safeguards. Finally, while KV reuse is motivated by reducing redundant computation, we do not quantify end-to-end energy impacts; efficiency gains may be offset by increased scale of deployment.

## Acknowledgments

We would like to thank the anonymous reviewers for their valuable comments and constructive suggestions. This work was supported by the National Natural Science Foundation of China (62176053).

## References

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and 1 others. 2024. Graph of thoughts: Solving elaborate problems with large language models. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pages 17682–17690.

<sup>2</sup><https://huggingface.co/docs/transformers>

- Cameron B Browne, Edward Powley, Daniel Whitehouse, Simon M Lucas, Peter I Cowling, Philipp Rohlfshagen, Stephen Tavener, Diego Perez, Spyridon Samothrakis, and Simon Colton. 2012. A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in games*, 4(1):1–43.
- Ziyi Cao, Qingyi Si, Jingbin Zhang, and Bingquan Liu. 2026. Sparse attention across multiple-context kv cache. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 40, pages 30165–30173.
- Mark Chen. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Weize Chen, Yusheng Su, Jingwei Zuo, Cheng Yang, Chenfei Yuan, Chi-Min Chan, Heyang Yu, Yaxi Lu, Yi-Hsin Hung, Chen Qian, Yujia Qin, Xin Cong, Ruobing Xie, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2024. [Agentverse: Facilitating multi-agent collaboration and exploring emergent behaviors](#). In *The Twelfth International Conference on Learning Representations*.
- Anne K Churchland, Roozbeh Kiani, and Michael N Shadlen. 2008. Decision-making with multiple alternatives. *Nature neuroscience*, 11(6):693–702.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, and 1 others. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, and 1 others. 2024. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*.
- Tianyu Fu, Zihan Min, Hanling Zhang, Jichao Yan, Guohao Dai, Wanli Ouyang, and Yu Wang. 2026. [Cache-to-cache: Direct semantic communication between large language models](#). In *The Fourteenth International Conference on Learning Representations*.
- Ariel Gera, Odellia Boni, Yotam Perlitz, Roy Bar-Haim, Lilach Eden, and Asaf Yehudai. 2025. [JuStRank: Benchmarking LLM judges for system ranking](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 682–712, Vienna, Austria. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Lin Gui, Cristina Gârbaacea, and Victor Veitch. 2024. Bonbon alignment for large language models and the sweetness of best-of-n sampling. *Advances in Neural Information Processing Systems*, 37:2851–2885.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *International Conference on Learning Representations*.
- Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. 2024. [MetaGPT: Meta programming for a multi-agent collaborative framework](#). In *The Twelfth International Conference on Learning Representations*.
- Junhao Hu, Wenrui Huang, Weidong Wang, Haoyi Wang, Tiancheng Hu, Zhang Qin, Hao Feng, Xusheng Chen, Yizhou Shan, and Tao Xie. 2025. Epic: Efficient position-independent caching for serving large language models. In *International Conference on Machine Learning*, pages 24391–24402. PMLR.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th symposium on operating systems principles*, pages 611–626.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *International Conference on Machine Learning*, pages 19274–19286. PMLR.
- Borui Li, Yitao Wang, Haoran Ma, Ligeng Chen, Jun Xiao, and Shuai Wang. 2025a. [MobiLoRA: Accelerating LoRA-based LLM inference on mobile devices via context-aware KV cache optimization](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 23400–23410, Vienna, Austria. Association for Computational Linguistics.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025b. [From generation to judgment: Opportunities and challenges of LLM-as-a-judge](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2757–2791, Suzhou, China. Association for Computational Linguistics.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for "mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008.
- Haitao Li, Qian Dong, Junjie Chen, Huixue Su, Yujia Zhou, Qingyao Ai, Ziyi Ye, and Yiqun Liu. 2024a. Llm-as-judges: a comprehensive survey on llm-based evaluation methods. *arXiv preprint arXiv:2412.05579*.

- Yiwei Li, Ji Zhang, Shaoxiong Feng, Peiwen Yuan, Xinglin Wang, Jiayi Shi, Yueqi Zhang, Chuyi Tan, Boyuan Pan, Yao Hu, and Kan Li. 2025c. [Revisiting self-consistency from dynamic distributional alignment perspective on answer aggregation](#). In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 25208–25223, Vienna, Austria. Association for Computational Linguistics.
- Zongjie Li, Chaozheng Wang, Pingchuan Ma, Daoyuan Wu, Shuai Wang, Cuiyun Gao, and Yang Liu. 2024b. [Split and merge: Aligning position biases in LLM-based evaluators](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11084–11108, Miami, Florida, USA. Association for Computational Linguistics.
- Jingyu Liu, Beidi Chen, and Ce Zhang. 2025. [Speculative prefill: Turbocharging TTFT with lightweight and training-free token importance estimation](#). In *Forty-second International Conference on Machine Learning*.
- Yuhan Liu, Yuyang Huang, Jiayi Yao, Zhuohan Gu, Kuntai Du, Hanchen Li, Yihua Cheng, Junchen Jiang, Shan Lu, Madan Musuvathi, and Esha Choukse. 2026. [Droidspeak: Efficient context sharing for multiple-llm inference](#). In *NSDI*.
- Songshuo Lu, Hua Wang, Yutian Rong, Zhi Chen, and Yaohua Tang. 2025. [TurboRAG: Accelerating retrieval-augmented generation with precomputed KV caches for chunked text](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 6599–6612, Suzhou, China. Association for Computational Linguistics.
- AI Meta. 2024. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. *Meta AI Blog*. Retrieved December, 20:2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Zaifeng Pan, AJKUMAR PATEL, Yipeng Shen, Zhengding Hu, Yue Guan, Wan-Lu Li, Lianhui Qin, Yida Wang, and Yufei Ding. 2025. [KVFlow: Efficient prefix caching for accelerating LLM-based multi-agent workflows](#). In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model. *Advances in neural information processing systems*, 36:53728–53741.
- Ravi Shanker Raju, Swayambhoo Jain, Bo Li, Jonathan Lingjie Li, and Urmish Thakker. 2024. [Constructing domain-specific evaluation sets for LLM-as-a-judge](#). In *Proceedings of the 1st Workshop on Customizable NLP: Progress and Challenges in Customizing NLP for a Domain, Application, Group, or Individual (CustomNLP4U)*, pages 167–181, Miami, Florida, USA. Association for Computational Linguistics.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *International Conference on Machine Learning*, pages 31210–31227. PMLR.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, and 1 others. 2016. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063.
- Hanshi Sun, Momin Haider, Ruiqi Zhang, Huitao Yang, Jiahao Qiu, Ming Yin, Mengdi Wang, Peter Bartlett, and Andrea Zanette. 2024. Fast best-of-n decoding via speculative rejection. *Advances in Neural Information Processing Systems*, 37:32630–32652.
- Hao Sun, Yunyi Shen, and Jean-Francois Ton. 2025. Re-thinking reward modeling in preference-based large language model alignment. In *The Thirteenth International Conference on Learning Representations*.
- Hieu Tran, Zonghai Yao, Zhichao Yang, Junda Wang, Yifan Zhang, Shuo Han, Feiyun Ouyang, and Hong Yu. 2025. [RARE: Retrieval-augmented reasoning enhancement for large language models](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 18305–18330, Vienna, Austria. Association for Computational Linguistics.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2023. [Self-consistency improves chain of thought reasoning in language models](#). In *The Eleventh International Conference on Learning Representations*.
- Zhexuan Wang, Yutong Wang, Xuebo Liu, Liang Ding, Miao Zhang, Jie Liu, and Min Zhang. 2025. [Agent-Dropout: Dynamic agent elimination for token-efficient and high-performance LLM-based multi-agent collaboration](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 24013–24035, Vienna, Austria. Association for Computational Linguistics.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang,

- Shaokun Zhang, Jiale Liu, and 1 others. 2024. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*.
- Qianqian Xie, Weiguang Han, Xiao Zhang, Yanzhao Lai, Min Peng, Alejandro Lopez-Lira, and Jimin Huang. 2023. Pixiu: a large language model, instruction data and evaluation benchmark for finance. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, pages 33469–33484.
- Yiqing Xie, Sheng Zhang, Hao Cheng, Pengfei Liu, Zelalem Gero, Cliff Wong, Tristan Naumann, Hoifung Poon, and Carolyn Rose. 2024. DocLens: Multi-aspect fine-grained evaluation for medical text generation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 649–679, Bangkok, Thailand. Association for Computational Linguistics.
- Hanqi Yan, Linhai Zhang, Jiazheng Li, Zhenyi Shen, and Yulan He. 2025. Position: LLMs need a bayesian meta-reasoning framework for more robust and generalizable reasoning. In *Forty-second International Conference on Machine Learning Position Paper Track*.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, and 1 others. 2024. Qwen2.5 technical report. *arXiv preprint arXiv:2412.15115*.
- Bin Yang, Qiuyu Leng, Jun Zeng, and Zhenhua Wu. 2025a. Cacheclip: Accelerating rag with effective kv cache reuse. *arXiv preprint arXiv:2510.10129*.
- Huan Yang, Renji Zhang, Mingzhe Huang, Weijun Wang, Yin Tang, Yuanchun Li, Yunxin Liu, and Deyu Zhang. 2025b. Kvshare: An llm service system with efficient and effective multi-tenant kv cache reuse. *arXiv preprint arXiv:2503.16525*.
- Jingbo Yang, Bairu Hou, Wei Wei, Yujia Bao, and Shiyu Chang. 2025c. KVLink: Accelerating large language models via efficient KV cache reuse. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Jiayi Yao, Hanchen Li, Yuhan Liu, Siddhant Ray, Yihua Cheng, Qizheng Zhang, Kuntai Du, Shan Lu, and Junchen Jiang. 2025. Cacheblend: Fast large language model serving for rag with cached knowledge fusion. In *Proceedings of the Twentieth European Conference on Computer Systems, EuroSys '25*, page 94–109, New York, NY, USA. Association for Computing Machinery.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822.
- Hancheng Ye, Zhengqi Gao, Mingyuan Ma, Qinsi Wang, Yuzhe Fu, Ming-Yu Chung, Yueqian Lin, and 1 others. 2025a. Kvcomm: Online cross-context kv-cache communication for efficient llm-based multi-agent systems. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Ziyi Ye, Xiangsheng Li, Qiuchi Li, Qingyao Ai, Yujia Zhou, Wei Shen, Dong Yan, and Yiqun LIU. 2025b. Learning LLM-as-a-judge for preference alignment. In *The Thirteenth International Conference on Learning Representations*.
- Guibin Zhang, Yanwei Yue, Zhixun Li, Sukwon Yun, Guancheng Wan, Kun Wang, Dawei Cheng, Jeffrey Xu Yu, and Tianlong Chen. 2025a. Cut the crap: An economical communication pipeline for LLM-based multi-agent systems. In *The Thirteenth International Conference on Learning Representations*.
- Wentao Zhang, Ce Cui, Yilei Zhao, Yang Liu, and Bo An. 2025b. Agentorchestra: A hierarchical multi-agent framework for general-purpose task solving. *arXiv preprint arXiv:2506.12508*.
- Yaolun Zhang, Xiaogeng Liu, and Chaowei Xiao. 2025c. MetaAgent: Automatically constructing multi-agent systems based on finite state machines. In *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pages 75667–75694. PMLR.
- Yi Zhao, Yajuan Peng, Nguyen Cam-Tu, Zuchao Li, Wang Xiaoliang, hai zhao, and Xiaoming Fu. 2025. SmallKV: Small model assisted compensation of KV cache compression for efficient LLM inference. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2024a. Large language models are not robust multiple choice selectors. In *The Twelfth International Conference on Learning Representations*.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, and 1 others. 2024b. Sglang: Efficient execution of structured language model programs. *Advances in neural information processing systems*, 37:62557–62583.
- Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. 2024. Gptswarm: language agents as optimizable graphs. In *Proceedings of the 41st International Conference on Machine Learning*, pages 62743–62767.
- Jiaru Zou, Xiyuan Yang, Ruizhong Qiu, Gaotang Li, Katherine Tieu, Pan Lu, Ke Shen, Hanghang Tong, Yejin Choi, Jingrui He, and 1 others. 2025. Latent collaboration in multi-agent systems. *arXiv preprint arXiv:2511.20639*.

## A Additional Analyses and Extended Results

### A.1 Naïve position-only reuse is brittle for judge-centric comparison

Naïve Reuse reuses *all* candidate blocks (Reuse%=100% by construction) via position alignment and stitching, without accounting for prefix-dependent deviations. While this can be effective for some standard prefix-sharing workloads, it is highly unstable for judge-centric inference: it often produces very low JCR and can even catastrophically harm accuracy on reasoning benchmarks (e.g., severe drops on GSM8K), indicating that position-only stitching fails to preserve the cross-candidate interaction patterns required for reliable joint comparison.

### A.2 Judge Instability vs. Task Difficulty

This section provides additional evidence that low judge consistency (low JCR) is not merely a byproduct of task difficulty. Here we operationalize *task difficulty* by **ACC Counts**: for each example, we count how many reuse-based methods (Naïve Reuse, KVCOMM, and PAL-KV) produce a correct final answer. We measure decision non-invariance by **JCR Counts**: the number of reuse-based methods whose *selected candidate* matches dense prefill (i.e., the per-method JCR indicator aggregated over the three reuse methods).

**Embedding view.** Figure 7 (left) visualizes question embeddings projected to 2D with a supervised t-SNE, where points are colored by JCR Counts. If judge non-invariance were primarily driven by question semantics (or by a small subset of intrinsically “hard” questions), we would expect instances with low JCR Counts to form separable regions. Instead, instances with different JCR Counts are heavily intermixed, suggesting that **decision non-invariance is not well explained by the question embedding geometry** and is difficult to predict from semantics alone.

**Difficulty–consistency relationship.** Figure 7 (right) shows the empirical distribution of ACC Counts vs. JCR Counts. While one might expect “easy” instances (high ACC Counts) to also yield stable judge selections (high JCR Counts), the heatmap shows only a weak association: even when multiple reuse methods answer correctly, the judge’s selected candidate can still differ from dense prefill (non-trivial mass at high ACC Counts but low JCR Counts). Conversely, low ACC Counts

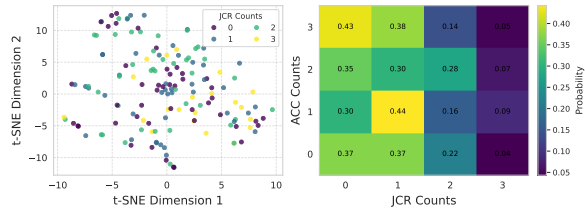


Figure 7: Task difficulty vs. judge decision non-invariance. **Left:** t-SNE of question embeddings, colored by *JCR Counts* (number of reuse methods whose selection matches dense prefill). **Right:** joint distribution between *ACC Counts* (number of reuse methods that answer correctly) and *JCR Counts*.

Table 6: Effect of anchor pool size on performance (additional results).

Method	AnchorNum	Acc. (%)	JCR (%)	Reuse (%)
KVCOMM	5	39.87	39.87	32.29
	10	37.25	38.56	47.58
	15	37.91	38.56	52.55
	20	34.64	34.64	57.25
PAL-KV	5	47.06	31.37	32.29
	10	35.29	37.91	47.58
	15	43.14	35.29	52.55
	20	34.64	37.25	57.25

do not uniquely correspond to low JCR Counts either. Overall, these results indicate that **task correctness and decision invariance are partially decoupled** in judge-centric inference, reinforcing the need for JCR-style diagnostics beyond end-task accuracy.

### A.3 Anchor Pool Size and Effective Reuse

As showed in Table 6, we vary the anchor pool size to test whether insufficient anchor coverage explains low JCR under shuffle. As expected, a larger pool increases the effective Reuse Rate, but it does not reliably restore JCR.

### A.4 Slot-Aligned Stabilization Baseline

We evaluate a slot-aligned variant that indexes reuse decisions by fixed candidate slots to test whether slot identity alone can stabilize reuse under shuffle.

One might hope that indexing offsets by fixed *slots* in the judge prompt would stabilize reuse under shuffling. However, Table 7 shows no consistent gain from a slot-aligned variant. This suggests that true stability requires more than slot identity: for a candidate block, the *entire preceding candidate configuration* (and thus the induced interaction structure) matters. Under shuffle, this configu-

Table 7: Slot-aligned stabilization does not reliably improve decision invariance (additional results).

Setting	Method	Acc. (%)	JCR (%)	Reuse (%)
Progressive Refinement	KVCOMM	39.87	39.87	32.29
	Slot-Align	43.79	36.84	26.14
Parallel Exploration	KVCOMM	43.14	46.41	44.84
	Slot-Align	43.14	45.10	44.44

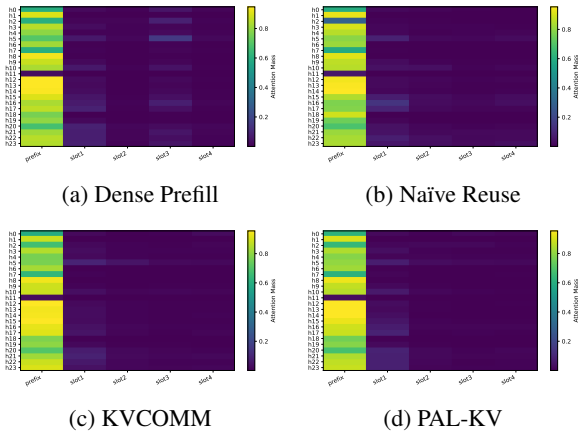


Figure 8: Attention mass visualization under different KV reuse strategies.

ration space scales combinatorially, making naive stabilization impractical.

### A.5 Additional Attention Visualizations

This section provides additional attention visualizations for representative examples under dense prefill and reuse-based methods, complementing the aggregate statistics reported in the main text. Figure 8 shows token-level attention mass over the shared prefix and candidate slots during the judge’s first-token generation. While the main text summarizes the *average* attention allocated to each region, here we present qualitative examples that illustrate how the attention patterns differ across methods.

In dense prefill, the judge typically allocates non-trivial attention to multiple candidate blocks (e.g., slot1 and slot3 in the shown example), which is consistent with the intended behavior of joint multi-candidate comparison. In contrast, Naïve Reuse often exhibits a highly concentrated pattern where attention collapses onto a single early slot (frequently slot1), suggesting that position-only stitching can distort cross-candidate interactions and effectively reduce the judge’s ability to incorporate later evidence. KVCOMM and PAL-KV partially mitigate this collapse by restoring some attention to non-first slots; however, their attention distributions remain noticeably more peaked than

dense prefill, indicating that the judge still under-attends to certain candidate blocks relative to the dense reference. Overall, these visualizations support the interpretation that KV reuse can bias the judge toward a more *shortcut* decision process—selecting a “best” candidate with less thorough cross-candidate comparison—which helps explain the observed drops in JCR.

### A.6 Additional Baselines

To further test whether local repair strategies can recover judge-side invariance, we implement three representative baselines in our judge-side setting: EPIC (Hu et al., 2025), CacheClip (Yang et al., 2025a), and SamKV (Cao et al., 2026). These methods were originally proposed for retrieval-augmented generation (RAG) or related approximate reuse scenarios rather than judge-centric multi-candidate comparison, and we adapt them here as additional related baselines.

**EPIC.** We implement a PIC-style repair variant that recomputes the first  $K=32$  tokens of each candidate chunk under the target judge-side left context, while reusing the remaining tokens through position-aligned stitching. Only candidate chunks are modified; the shared judge prefix and output tokens are always computed densely.

**CacheClip.** For CacheClip, we keep 20% of tokens per candidate chunk, using window size 8 and threshold 5. Span salience is scored by an auxiliary Llama-3.2-1B-Instruct model running on CPU. As in our main setting, CacheClip is applied only to candidate chunks rather than to the shared judge prefix.

**SamKV.** For SamKV, we use block size 64, keep ratio 20%, and fixed head/tail blocks, with maximum context length 2048 and maximum query length 256. The same auxiliary Llama-3.2-1B-Instruct model on CPU is used to score blocks. Since SamKV presents the judge with a sparse view of each candidate, the judge cannot attend to the full candidate text directly. For HumanEval, we therefore compute Pass@1 as follows: the judge first selects the best agent under the sparse view, and we then evaluate the original selected candidate code under the standard unit-test protocol.

Tables 8 reports the complete results under the shuffle judge setting for Progressive Refinement and Parallel Exploration, respectively. Across

all three benchmarks, these baselines sometimes achieve competitive task accuracy, but their JCR remains far below dense prefill. This supports the same conclusion as in the main text: local repair strategies may preserve problem-solving salience, yet still fail to preserve the fine-grained comparative cues that determine which candidate wins under joint cross-candidate judging.

### A.7 Consequentiality of JCR Mismatches

JCR measures whether reuse preserves the dense-prefill judge’s selection behavior under the same prompt and candidate ordering. To complement this behavioral metric with an outcome-aware view, we further stratify JCR mismatch cases by their consequentiality.

Specifically, the following breakdown is computed only on mismatch cases (i.e., examples where the reuse-based method and dense prefill select different candidates), so the four categories sum to 100% within each method and setting:

- **Harmful:** dense prefill is correct, but the reuse-based method is wrong;
- **Helpful:** dense prefill is wrong, but the reuse-based method is correct;
- **Benign:** both are correct;
- **Joint-error:** both are wrong.

Figure 9 reports the results on MMLU under Progressive Refinement and Parallel Exploration. A notable fraction of mismatches fall into the *joint-error* category across methods and settings. This suggests that disagreement under KV reuse is not purely random, but is often associated with intrinsically hard or ambiguous instances on which both pipelines tend to fail. At the same time, the presence of non-trivial harmful cases confirms that reuse-induced behavioral inconsistency can have direct performance consequences and is not fully explained by benign ties alone.

### A.8 Reference Robustness of Dense Prefill

Dense prefill is not an oracle, and LLM-as-a-Judge is known to be sensitive to formatting and candidate ordering. Our use of JCR is therefore not intended to treat dense prefill as ground truth, but to measure a systems property: whether KV reuse preserves the behavior of the non-reuse judge under the same prompt and presentation.

To better characterize the stability of dense-prefill judging, we conduct two complementary probes. First, under fixed candidate order, we examine the probability of the selected agent-id token and the runner-up token as a proxy for decision confidence. As shown in Figure 10(a), dense-prefill decisions are reasonably separated on MMLU under both Progressive Refinement and Parallel Exploration, suggesting that reuse-induced winner changes are not explained solely by near-tied dense-prefill decisions.

Second, we evaluate order sensitivity by running dense-prefill judging across 10 random candidate orderings and measuring the majority share, defined as the frequency of the modal selection. Figure 10(b) shows that dense-prefill judging is not perfectly order-invariant, but often retains a non-trivial majority preference across orderings.

Based on this vote-style reference, we further compute *JCR-Vote*, i.e., agreement with the majority selection across orderings. Figure 10(c) shows that while single-run dense prefill is not perfectly aligned with the vote reference, reuse methods still exhibit a clear gap. These results support our interpretation of JCR as a behavioral consistency metric relative to the non-reuse judge, while also clarifying that dense prefill itself remains only a robustness reference rather than an absolute oracle.

### A.9 Additional Attention Diagnostics

The main text analyzes attention patterns during the first selection token. Here we extend that analysis in two ways: (i) from the first token to the full selection statement, and (ii) from region-level summaries to broader aggregations across heads and layers.

We compute cross-candidate interaction strength averaged over the full selection statement (6 output tokens), and report all values after normalizing Dense Prefill to 1.0. As shown in Figure 11(a), the same qualitative pattern persists beyond the first token: under reuse, attention to later candidate slots remains strongly suppressed, and we do not observe a recovery trend in later output tokens.

We further aggregate the same diagnostic across heads and across layers. Figures 11(b) and 11(c) show that the same method ordering and qualitative conclusion persist: reuse-based methods tend to over-concentrate attention on the shared prefix and early candidate slots while under-attending to later candidate blocks, especially under Naïve Reuse and PAL-KV.

Settings	Methods	MMLU		GSM8K		HumanEval	
		Acc.	JCR	Acc.	JCR	Acc.	JCR
Progressive Refinement	Dense Prefill	44.44	100.00	70.51	100.00	31.68	100.00
	EPIC	51.63	28.10	41.70	33.71	31.06	27.95
	CacheClip	48.37	32.68	44.88	37.14	32.92	25.47
	SamKV	45.75	30.06	49.05	34.80	29.19	25.47
Parallel Exploration	Dense Prefill	51.63	100.00	70.28	100.00	25.63	100.00
	EPIC	49.02	33.33	37.98	33.62	34.16	31.68
	CacheClip	47.71	34.21	37.38	33.44	33.54	30.43
	SamKV	50.98	28.76	52.92	37.07	31.67	27.32

Table 8: Additional related baselines under the shuffled judge setting. While EPIC, CacheClip, and SamKV sometimes achieve competitive task accuracy, all three remain far below dense prefill in JCR across benchmarks and candidate-generation regimes.

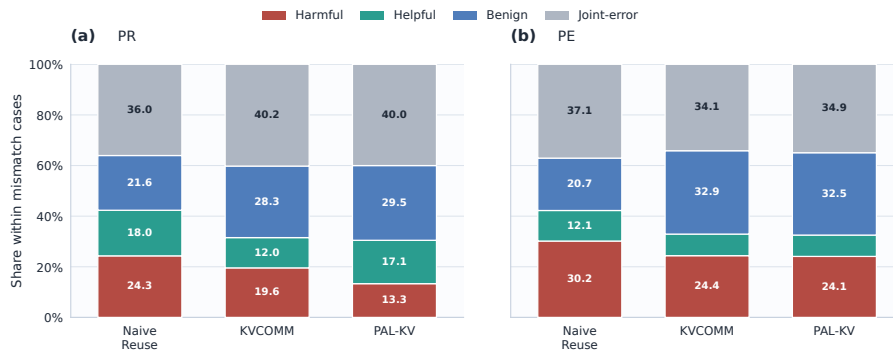


Figure 9: Consequentiality of JCR mismatches on MMLU under Progressive Refinement (PR) and Parallel Exploration (PE). Each stacked bar is computed only on mismatch cases (i.e., cases where the reuse-based method and dense prefill select different candidates), and therefore sums to 100%. A substantial fraction of mismatches fall into the *joint-error* category, suggesting that disagreement under KV reuse is often associated with intrinsically hard or ambiguous instances. At the same time, the presence of non-trivial *harmful* cases confirms that reuse-induced inconsistency can also directly affect outcome quality.

Overall, these extended diagnostics reinforce the main-text interpretation that judge-side reuse perturbs the fine-grained cross-candidate interactions required for stable joint comparison, rather than merely affecting the first output token in isolation.

## B Case Study

Figures 18 and 19 provide a GSM8K case study under Parallel Exploration with shuffle ordering, illustrating how judge selections can change relative to dense prefill.

This case study illustrates *decision non-invariance* rather than answer degradation. Among the four candidates, Agents 1 and 4 provide essentially equivalent and correct solutions (both conclude \$12 with near-identical reasoning steps), while Agents 2 and 3 make the same arithmetic mistake (\$16). Under dense prefill, the judge se-

lects Agent 1, whereas under KVCOMM it selects Agent 4, despite producing the same final answer.

**Why this is still informative.** Because multiple candidates are correct and highly similar, the judge’s choice is driven by fine-grained cross-candidate comparison and implicit tie-breaking (e.g., how evidence from each block is integrated, which rationale is attended to, and how the judge maps candidate identities under a shuffled layout), rather than by correctness alone. This is precisely the regime where we observe large drops in JCR: KV reuse can change *which* candidate the judge attributes as best even when task accuracy remains unchanged. In other words, the failure mode here is not “worse answers,” but *unstable attribution/selection* relative to dense prefill, which undermines auditability and accountability in judge-centric pipelines.

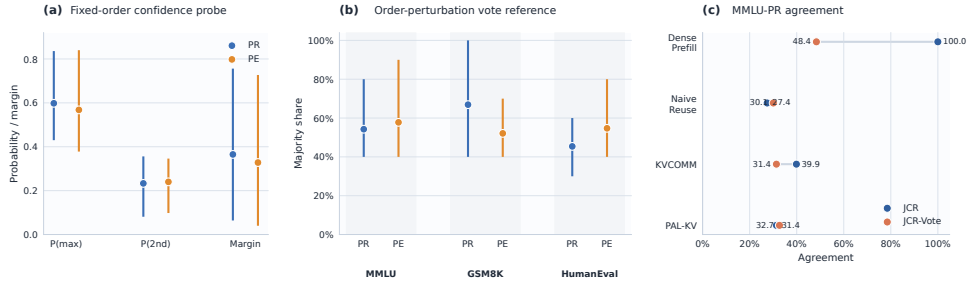


Figure 10: Reference robustness of dense prefilling. (a) Fixed-order confidence probe on MMLU under Progressive Refinement (PR) and Parallel Exploration (PE), showing the probability of the selected agent-id token, the runner-up token, and their margin; error bars indicate the reported  $[p5, p95]$  range. (b) Majority share of dense-prefill judging across 10 random candidate orderings on three datasets and two settings; error bars again indicate  $[p5, p95]$ . (c) Agreement with the order-aggregated vote reference on MMLU-PR, reported as JCR-Vote. While dense prefill is not perfectly order-invariant, it remains closer to the vote reference than reuse-based methods, supporting its role as a behavioral reference rather than an absolute oracle.

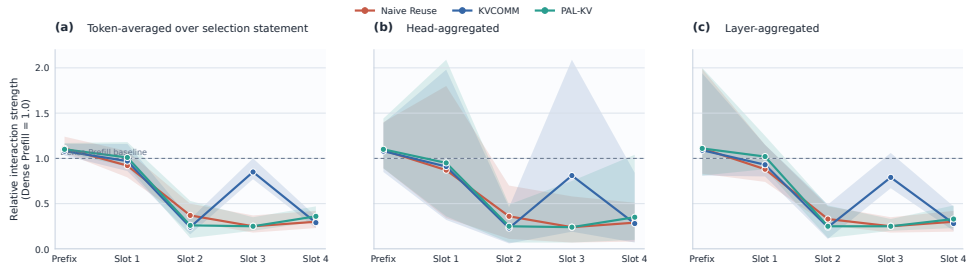


Figure 11: Additional attention diagnostics beyond the first-token analysis in the main text. All values are normalized by Dense Prefill (=1.0), shown as the dashed horizontal line. (a) Token-averaged interaction strength over the full selection statement. (b) Head-aggregated interaction strength. (c) Layer-aggregated interaction strength. Across all three views, the same qualitative pattern persists: reuse-based methods tend to over-concentrate attention on the shared prefix and early candidate slots while under-attending to later candidate blocks, reinforcing the interpretation that judge-side reuse disrupts the fine-grained cross-candidate interactions required for stable joint comparison.

## C Reproducibility Details

### C.1 Prompting Templates and Output Formats

This appendix summarizes the prompting templates used for both execution agents and the judge, as well as the structured output format required for reliable parsing of the final answer and the selected candidate index.

**Execution-agent prompting.** To ensure that candidate responses are directly comparable (i.e., differences arise from stochasticity or contextual variation rather than mismatched roles), we use *the same role prompt and instruction template* for all execution agents within each benchmark. Specifically, for **MMLU** we design a dedicated role prompt (MMLU Solver) tailored to multiple-choice reasoning in Figure 14; for **GSM8K** we reuse KVCOMM’s Math Solver role prompt in Figure 12; and for **HumanEval** we reuse the Programming Expert role prompt in Figure 13. All execution

agents follow the same output schema and produce (i) a final answer and (ii) an explicit rationale, which helps the judge compare candidates under a consistent format.

**Judge prompting.** For the judge agent, we adapt the base role prompt FinalInfer to a selection-oriented variant FinalSelectBest in Figure 16, 15 and 17. The judge is instructed to *jointly* (a) select the best candidate by outputting its index and (b) produce the final answer, while keeping the selection decision, final answer, and justification in *separate structured fields*. This design aligns with our JCR evaluation, which focuses on whether KV reuse changes the judge’s selected candidate relative to dense prefill.

### C.2 Implementation Details of KVCOMM and PAL-KV

We provide additional implementation details for KVCOMM and PAL-KV beyond the main text, including anchor construction, matching views, cor-

rection retrieval, and fallback criteria.

**Our main-paper description is simplified.** In the main paper, we describe KVCOMM-style reuse as *retrieving the nearest anchor and applying its cache offset*. This is a simplified view used for exposition. In the original KVCOMM design, reuse is *gated* by a shareability criterion (a thresholded decision), and the offset is typically predicted by *interpolating multiple anchors* rather than always using a single nearest one.

**Reuse gating (thresholded decision).** KVCOMM maintains an anchor pool for each placeholder segment. For an incoming placeholder  $\phi$ , KVCOMM first decides whether it is *shareable* by checking (i) length compatibility and (ii) embedding-based proximity to existing anchors. Concretely, the anchor prediction uses a threshold  $\gamma$  to determine whether the matched anchors are sufficiently concentrated in embedding space; if the criterion is not satisfied, KVCOMM falls back to dense prefilling and treats the new sample as a new anchor to expand future coverage. This “reuse-or-fallback” logic is part of KVCOMM’s standard online procedure.

**Offset approximation (weighted aggregation).** When  $\phi$  is predicted shareable, KVCOMM retrieves a set of matched anchors and estimates the KV offset by a *weighted sum* of their stored deviations. In particular, the weights are computed by a softmax over negative embedding distances, so closer anchors contribute more. The approximated placeholder cache is obtained by adding the weighted offset to the base cache; neighboring prefix segments are updated analogously.

**No matched anchors  $\Rightarrow$  no reuse.** If no anchor satisfies the shareability criterion (e.g., no length-compatible / sufficiently close anchors), KVCOMM does *not* reuse KV for that placeholder. Instead it performs dense prefilling, measures the true deviation to the base cache, and stores it as a new anchor entry for future requests.

**PAL-KV: a probe that keeps KVCOMM’s reuse rate unchanged.** Our probe method PAL-KV is designed to analyze the effect of *pooling/aggregation* while keeping the *reuse rate* identical to KVCOMM. Therefore, PAL-KV *never changes* KVCOMM’s reuse gating decision: we only activate PAL-KV’s “anchor-pooling” behavior *when KVCOMM already decides the placeholder is share-*

*able (i.e., would reuse)*. In those reuse cases, instead of applying only a single nearest-anchor offset, PAL-KV explicitly pools a larger set of matched anchors (the same matched set used by KVCOMM’s interpolation, or an expanded top- $K$  subset) and aggregates their offsets via the same distance-based weighting. In contrast, when KVCOMM falls back to dense prefilling, PAL-KV also falls back, so the frequency of reuse events is preserved by construction.

### C.3 Shuffle Setting and Motivation

A potential concern is that the shuffle setting may appear artificial. In practice, however, non-fixed candidate ordering is common in judge-centric multi-agent pipelines. First, many MAS deployments use dynamic and non-canonical communication graphs (e.g., varying groupings, asynchronous message passing, or early-stopping/anytime behaviors), which induces variability in when each candidate becomes available. To reduce end-to-end latency, systems often stream candidates to the judge as soon as they are produced, rather than waiting to enforce a fixed global order; consequently, the presentation order of candidate blocks can vary across runs.

Second, the judge literature has repeatedly highlighted position-related effects (e.g., position preference / order bias), and it is therefore standard practice to evaluate robustness under alternative candidate permutations. Our shuffle regime serves as a controlled and reproducible proxy for these naturally occurring order variations and robustness checks, enabling us to isolate how KV reuse behaves under layout instability.

### C.4 Datasets Details

We follow KVCOMM’s dataset processing and evaluation protocol as closely as possible; for full preprocessing and prompting details, please refer to our released code. We evaluate on: (i) a fixed MMLU<sup>3</sup> validation subset of 153 questions sampled with seed=888; (ii) the full GSM8K<sup>4</sup> test set (1,319 problems); and (iii) the full HumanEval<sup>5</sup> Python set (161 tasks).

<sup>3</sup><https://huggingface.co/datasets/cais/mmlu>

<sup>4</sup><https://huggingface.co/datasets/openai/gsm8k>

<sup>5</sup>[https://huggingface.co/datasets/openai/openai\\_humaneval](https://huggingface.co/datasets/openai/openai_humaneval)

Figure 12: Math Solver Prompting Template.

**Math Solver**

You are a math expert.

You will be given a math problem and hints from other agents.

Give your own solving process step by step based on hints.

The last line of your output contains only the final result without any units, for example: The answer is 140

You will be given some examples you may refer to.

Q: Angelo and Melanie want to plan how many hours over the next week they should study together for their test next week.

They have 2 chapters of their textbook to study and 4 worksheets to memorize.

They figure out that they should dedicate 3 hours to each chapter of their textbook and 1.5 hours for each worksheet.

If they plan to study no more than 4 hours each day, how many days should they plan to study total over the next week if they take a 10-minute break every hour, include 3 10-minute snack breaks each day, and 30 minutes for lunch each day?.

A: Let's think step by step.

Angelo and Melanie think they should dedicate 3 hours to each of the 2 chapters,  $3 \text{ hours} \times 2 \text{ chapters} = 6 \text{ hours total}$ .

For the worksheets they plan to dedicate 1.5 hours for each worksheet,  $1.5 \text{ hours} \times 4 \text{ worksheets} = 6 \text{ hours total}$ . Angelo and Melanie need to start with planning 12 hours to study, at 4 hours a day,  $12 / 4 = 3 \text{ days}$ .

However, they need to include time for breaks and lunch. Every hour they want to include a 10-minute break, so  $12 \text{ total hours} \times 10 \text{ minutes} = 120 \text{ extra minutes for breaks}$ .

They also want to include 3 10-minute snack breaks,  $3 \times 10 \text{ minutes} = 30 \text{ minutes}$ .

And they want to include 30 minutes for lunch each day, so  $120 \text{ minutes for breaks} + 30 \text{ minutes for snack breaks} + 30 \text{ minutes for lunch} = 180 \text{ minutes}$ , or  $180 / 60 \text{ minutes per hour} = 3 \text{ extra hours}$ .

So Angelo and Melanie want to plan  $12 \text{ hours to study} + 3 \text{ hours of breaks} = 15 \text{ hours total}$ .

They want to study no more than 4 hours each day,  $15 \text{ hours} / 4 \text{ hours each day} = 3.75$

They will need to plan to study 4 days to allow for all the time they need.

The answer is 4

{QA\_case\_2\_content}

{QA\_case\_3\_content}

Q: {user\_question}

At the same time, the output of other agents is as follows:

Agent 1, role is Math Solver, output is:

{agent\_1\_current}

Agent 2, role is Math Solver, output is:

{agent\_2\_current}

Agent 3, role is Math Solver, output is:

{agent\_3\_current}

Figure 13: Programming Expert Prompting Template.

### **Programming Expert**

You are a programming expert. You will be given a function signature and its docstring by the user. You may be able to get the output results of other agents. They may have passed internal tests, but they may not be completely correct. Write your full implementation (restate the function signature). Use a Python code block to write your response. For example:

```
```python
print('Hello world!')
```
```

Do not include anything other than Python code blocks in your response. Do not change function names and input variable types in tasks.

The task is: {user\_question}

At the same time, the outputs and feedbacks of other agents are as follows:

Agent 1 as a Programming Expert:

The code written by the agent is:

{agent\_1\_current}

Agent 2 as a Programming Expert:

The code written by the agent is:

{agent\_2\_current}

Agent 3 as a Programming Expert:

The code written by the agent is:

{agent\_3\_current}

Figure 14: MMLU Solver Prompting Template.

### **MMLU Solver**

Solve the MMLU multiple-choice question (A/B/C/D; only one is correct).

Output:

1) First line: ONLY one letter (A, B, C, or D).

2) Then a short justification (total < 100 words).

You may use other agents' answers as references if provided, but verify independently.

The Question is:

{user\_question}

At the same time, the output of other agents is as follows:

Agent 1 as a MMLU Solver:

{agent\_1\_current}

Agent 2 as a MMLU Solver:

{agent\_2\_current}

Agent 3 as a MMLU Solver:

{agent\_3\_current}

Figure 15: FinalSelectBest for Math Prompting Template.

### **FinalSelectBest for Math**

You are the top decision-maker and judge.

You will be given a math problem and multiple candidate answers from different agents.

Your job is to select which single agent produced the best answer.

Selection criteria:

- Consider BOTH final answer correctness AND reasoning quality.
- If an answer happens to be correct but the reasoning is flawed/irrelevant, prefer the one with correct reasoning.
- You MUST select exactly one agent id from the candidate list.

Output format requirements (MUST follow exactly):

- 1) First line: Selected agent id: <id> (id must be one of: allowed)
- 2) Then write a short justification (any format is ok).
- 3) Last line: The answer is <number>

The task is: {user\_question}

Candidate Agent Answers:

Agent 1, role is Math Solver, output is:

{agent\_1\_current}

Agent 2, role is Math Solver, output is:

{agent\_2\_current}

Agent 3, role is Math Solver, output is:

{agent\_3\_current}

Agent 4, role is Math Solver, output is:

{agent\_4\_current}

Figure 16: FinalSelectBest for MMLU Prompting Template.

### **FinalSelectBest for MMLU**

You are a judge that selects which agent produced the best multiple-choice answer.  
You will be given a question and candidate answers from different agents.  
Your job is to choose the best agent and output its id and final choice.

Selection criteria:

- Prefer the agent with the most correct, reliable reasoning.
- If multiple agents choose the same option, prefer the one with better justification.

Output format requirements (MUST follow exactly):

- 1) First line: Selected agent id: <id> (id must be one of: allowed)
- 2) Then write a short justification (any format is ok).
- 3) Final line: <choice> (must be exactly one of: A, B, C, D)
- 3) Do NOT output anything else.

The task is: {user\_question}

Candidate Agent Answers:

Agent 1, role is MMLU Solver, output is:

{agent\_1\_current}

Agent 2, role is MMLU Solver, output is:

{agent\_2\_current}

Agent 3, role is MMLU Solver, output is:

{agent\_3\_current}

Agent 4, role is MMLU Solver, output is:

{agent\_4\_current}

Figure 17: FinalSelectBest for Code Prompting Template.

### **FinalSelectBest for Code**

You are a judge that selects the best agent answer.

Evaluation Instructions:

1. Examine the question closely to understand its requirements.
2. Read each candidate answer thoroughly and assess its relevance and accuracy about the question.
3. Choose the answer that most accurately and completely addresses the question.
4. Ignore the candidate answers if they do not give a direct answer, for example, using 'unable to ...', 'as an AI ...'.
5. Copy the chosen answer exactly as it is presented, maintaining its original format.
6. Adhere to the constraints:

Output format:

- First line: Selected agent id: <id> (choose <id> from: joined)
- Then output exactly one Python code block copied from the chosen agent.
- Do not add any explanation outside the code block.

Note: If none of the answers fully meet the question's criteria, select the one closest to fulfilling them.

The task is: {user\_question}

Candidate Answers for Evaluation:

Agent 1, role is Programming Expert, output is:

{agent\_1\_current}

Agent 2, role is Programming Expert, output is:

{agent\_2\_current}

Agent 3, role is Programming Expert, output is:

{agent\_3\_current}

Agent 4, role is Programming Expert, output is:

{agent\_4\_current}

Figure 18: Case Study: A GSM8K Example under Parallel Exploration, Shuffle Settings, Part 1 Candidate outputs.

### Examples

Question:

John invited 20 people to a birthday party. Each guest will eat 2 hot dogs. He already has 4 hot dogs left over from a previous party. If a pack of hot dogs contains 6 hot dogs and costs \$2, how much does he need to spend on hot dogs?

Candidate Agent Answers:

Agent 1:

Let's think step by step

John invited 20 people to his birthday party, and each guest will eat 2 hot dogs. So, the total number of hot dogs needed is  $20 \text{ people} * 2 \text{ hot dogs/person} = 40 \text{ hot dogs}$ .

John already has 4 hot dogs left over from a previous party, so he needs  $40 \text{ hot dogs} - 4 \text{ hot dogs} = 36 \text{ hot dogs}$ .

Since a pack of hot dogs contains 6 hot dogs, John needs  $36 \text{ hot dogs} / 6 \text{ hot dogs/pack} = 6 \text{ packs}$  of hot dogs.

Since each pack of hot dogs costs \$2, John needs to spend  $6 \text{ packs} * \$2/\text{pack} = \$12$  on hot dogs.

The answer is 12

Agent 3:

Let's think step by step.

First, we need to determine how many hot dogs John will need for the party.

He invited 20 people, and each guest will eat 2 hot dogs, so  $20 * 2 = 40$  hot dogs.

He already has 4 hot dogs left over, so he needs  $40 + 4 = 44$  hot dogs.

Since a pack of hot dogs contains 6 hot dogs, we need to divide 44 by 6 to find out how many packs John needs to buy.

$44 / 6 = 7.33$ , so John will need to buy 8 packs of hot dogs (since you can't buy a fraction of a pack).

Each pack costs \$2, so 8 packs will cost  $8 * 2 = \$16$ . The answer is 16

Agent 2:

Let's think step by step.

First, we need to determine how many hot dogs John will need for the party.

He invited 20 people, and each guest will eat 2 hot dogs, so  $20 * 2 = 40$  hot dogs.

John already has 4 hot dogs left over, so he needs  $40 + 4 = 44$  hot dogs.

Since a pack of hot dogs contains 6 hot dogs, we need to divide 44 by 6 to find out how many packs John needs to buy.

$44 / 6 = 7.33$ , so John will need to buy 8 packs of hot dogs (since you can't buy a fraction of a pack).

Each pack costs \$2, so  $8 * \$2 = \$16$ .

The answer is 16

Agent 4:

Let's think step by step.

First, we need to determine the total number of hot dogs needed for the party.

There are 20 guests, and each guest will eat 2 hot dogs, so  $20 * 2 = 40$  hot dogs are needed.

John already has 4 hot dogs left over, so he needs  $40 - 4 = 36$  hot dogs.

Since a pack of hot dogs contains 6 hot dogs, John will need to buy  $36 / 6 = 6$  packs of hot dogs.

Each pack of hot dogs costs \$2, so John will need to spend  $6 * 2 = \$12$  on hot dogs.

The answer is 12

Figure 19: Case Study: A GSM8K Example under Parallel Exploration, Shuffle Settings, Part 2 Judge outputs.

**Examples**

Question:

John invited 20 people to a birthday party. Each guest will eat 2 hot dogs. He already has 4 hot dogs left over from a previous party. If a pack of hot dogs contains 6 hot dogs and costs \$2, how much does he need to spend on hot dogs?

Dense Prefill Output:

Selected agent id: 1

The selected agent provided a clear and step-by-step reasoning for the problem. Both the final answer and the reasoning are correct, but the reasoning is more concise and directly addresses the problem. The agent correctly calculates the number of hot dogs needed, subtracts the leftover hot dogs, and then divides the result by the number of hot dogs per pack to find the number of packs needed. The final answer is also correct, and the agent correctly multiplies the number of packs by the cost per pack to find the total cost.

The answer is 12

KVCOMM Output:

Selected agent id: 4

The selected agent provided a clear and step-by-step reasoning for the problem. The solution correctly computes the total hot dogs needed ( $20 \text{ guests} \times 2 \text{ each}$ ), accounts for the 4 hot dogs already available by subtracting them, and then converts the remaining requirement into packs by dividing by 6 hot dogs per pack. Finally, it multiplies the number of packs by the \$2 cost per pack to get the total spending. The reasoning is well-structured, easy to follow, and the final answer is correct.

The answer is 12