
Enhancing Multi-Agent Communication through Attention Steering with Context Relevance

Hongxiang Zhang
Purdue University
West Lafayette, IN, USA
hxxzhang@purdue.edu

Yuan Tian
Purdue University
West Lafayette, IN, USA
tian211@purdue.edu

Tianyi Zhang
Purdue University
West Lafayette, IN, USA
tianyi@purdue.edu

Abstract

LLM-based multi-agent systems have demonstrated remarkable performance on complex tasks through collaborative reasoning. However, these systems tend to rapidly accumulate extremely long conversation histories during interaction. As conversations lengthen, relevant information is increasingly diluted by irrelevant context, leading to degraded performance. In this work, we present AGENT-RADAR, a training-free context management method that dynamically steers each agent’s attention toward relevant context with a novel temporal and spatial decay mechanism. Our experiments demonstrate that AGENT-RADAR outperforms state-of-the-art methods across five different benchmarks, yielding gains of up to 7.64 absolute points. Furthermore, our analysis shows that AGENT-RADAR remains effective and robust as the number of agents and interaction rounds increases. Finally, the ablation study shows that core components in AGENT-RADAR are crucial to performance and generalizable in different settings.

1 Introduction

LLM-based multi-agent systems have demonstrated remarkable performance on complex tasks [Liang et al., 2024, Wang et al., 2024b, Du et al., 2024, Wu et al., 2024]. Despite these successes, agent communication often accumulates extremely long context windows rapidly [Zhang et al., 2025a] and suffers from communication redundancy [Wang et al., 2025b, Zhang et al., 2025a, Du et al., 2024, Liang et al., 2024]. As a result, the original task description, key constraints, and other critical evidence are often buried within the context, making it difficult for LLMs to reliably identify and attend to them. This phenomenon, also known as the “lost-in-the-middle” effect, frequently induces hallucinations and logical errors during generation [Becker et al., 2026, Liu et al., 2024a, Cemri et al., 2025, Zhang et al., 2025b]. This issue is further exacerbated as systems scale. Increasing numbers of agents and interaction rounds cause communication histories to grow rapidly, amplifying context dilution and undermining the reliability of large-scale agent collaborations [Cemri et al., 2025, Becker et al., 2026, Zhang et al., 2025b].

Prior context management methods have sought to mitigate this challenge in two ways. One line of work applies periodic summarization to compress accumulated interaction history [Wang et al., 2025a, Wu et al., 2025]. Another adopts pruning strategies that limit context growth by removing redundant agents or communication edges [Zhang et al., 2025a, Wang et al., 2025b]. While effective at improving efficiency, these methods introduce information loss. Specifically, compressing or

pruning messages can discard subtle but critical signals, potentially misleading downstream reasoning and propagating errors across subsequent interactions [Shen et al., 2025b, Acharjee et al., 2025].

We propose a shift from context compression or pruning to selectively steering agents’ attention over relevant context. In the original Transformer architecture, the attention mechanism retrieves information by matching queries to keys, allowing the model to assign importance weights across preceding tokens in the context [Vaswani et al., 2017]. Accordingly, in multi-agent systems, an agent’s current generation step naturally acts as a query seeking to retrieve relevant information from the communication history. Unlike prior methods that act as hard, binary filters over context, we instead steer attention toward a subset of relevant content while preserving the full transcript, thereby avoiding information loss.

To achieve this, we develop AGENT-RADAR, a training-free method that dynamically guides agents’ attention toward relevant context without modifying the underlying communication history. To identify which parts of the history should guide each generation step, AGENT-RADAR selects context by jointly modeling semantic relevance, spatial distance, and temporal recency. First, inspired by organizational network theory [Borgatti et al., 2009, Allen, 1984], where influence decays with structural distance, AGENT-RADAR formulates the agent interaction topology as a graph and applies a spatial decay factor based on inter-agent distance, prioritizing messages from closely interacting agents. Second, AGENT-RADAR incorporates temporal decay to down-weight stale information, reflecting evidence from both long-context modeling and human memory that relevance diminishes over time [Finkenbinder, 1913, Sun et al., 2023, Ma et al., 2023]. Third, AGENT-RADAR performs fine-grained semantic matching by segmenting messages into sentence-level units to extract relevant evidence at the current step. Finally, AGENT-RADAR integrates these signals through a lightweight attention-steering mechanism, enabling agents to emphasize key instructions, critical constraints, and useful intermediate evidence, without compressing or pruning history tokens.

We evaluate AGENT-RADAR on five benchmarks using three base LLMs. Compared with SOTA context management methods, including compression-based approach [Wang et al., 2025a] and pruning methods [Zhang et al., 2025a, Wang et al., 2025b], AGENT-RADAR achieves gains of up to 7.64 absolute points. As a plug-in method, AGENT-RADAR can be seamlessly integrated into diverse multi-agent systems, including debate-based systems [Du et al., 2024], GPTSwarm [Zhuge et al., 2024], and AutoGen [Wu et al., 2024]. Across these frameworks, AGENT-RADAR consistently improves their performance by up to 12.87 absolute points and generalizes across different interaction topologies. Furthermore, AGENT-RADAR remains robust as systems scale to larger numbers of agents and longer interaction rounds. Finally, ablation studies confirm that each component contributes meaningfully to the overall performance gains.

2 Problem Formulation

We model the multi-agent system architecture as a graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{F})$. Let $\mathcal{V} = \{v_i\}_{i=1}^N$ be the set of N nodes, where each node v_i represents an autonomous LLM agent. The edge set $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ defines the communication topology. An edge $(v_i, v_j) \in \mathcal{E}$ denotes the communication path from the agent v_i to the agent v_j . The set of mapping functions $\mathcal{F} = \{f_i\}_{i=1}^N$ represents the inference processes of the agents. We further define the hop distance $d_{\mathcal{G}}(v_j, v_i)$ as the length of the shortest path from v_j to v_i in \mathcal{G} . The k -hop neighborhood of v_i is $\mathcal{N}_k(v_i) = \{v_j \in \mathcal{V} \mid d_{\mathcal{G}}(v_j, v_i) \leq k\}$. All agents that can reach to v_i in \mathcal{G} is $\mathcal{N}_{N-1}(v_i) = \bigcup_{k=1}^{N-1} \mathcal{N}_k(v_i)$.

The system proceeds over communication rounds $t \in \{1, \dots, T\}$. At time t , the active agent v_i generates the response

$$o_i^{(t)} = f_i(q^{(t)}, \mathcal{Z}_i^{(t)}), \quad (1)$$

where $q^{(t)}$ is the query received by the current agent at time t , and $\mathcal{Z}_i^{(t)}$ is the accumulated communication histories available to v_i . The context available to the current agent v_i at time t is

$$\mathcal{Z}_i^{(t)} = \{o_j^{(\tau)} \mid v_j \in \mathcal{N}_{N-1}(v_i), 1 \leq \tau < t\}, \quad (2)$$

where $o_j^{(\tau)}$ is the output produced by agent v_j at timestamp τ .

In our method, we introduce an additional parameter $\hat{\mathcal{Z}}_i^{(t)}$ to f_i , denoting the selected context that AGENT-RADAR steers attention towards for each inference step. In the following section, we describe

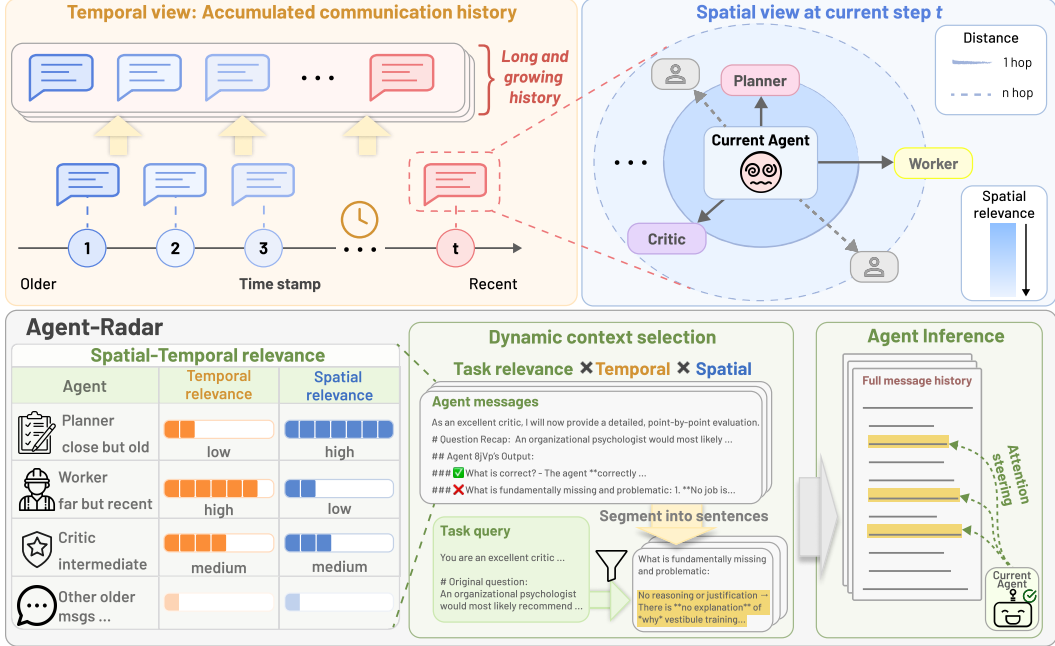


Figure 1: Overview of AGENT-RADAR. (Top) MAS interactions rapidly accumulate long communication histories, where useful information is buried in the middle, receiving insufficient attention. (Bottom) AGENT-RADAR preserves the full transcript and topology, scores sentence-level context by semantic relevance weighted with temporal and spatial decay, and steers the agent’s attention toward the selected context during inference.

how AGENT-RADAR selects $\hat{Z}_i^{(t)}$ from the accumulated communication histories and uses it as an additional argument to f_i for attention steering.

3 Methodology

We introduce AGENT-RADAR (Figure 1), a training-free method that augments multi-agent systems (MAS) by steering an agent’s attention toward semantically relevant context that is both spatially proximate and temporally salient. Building on the formulation in Section 2, AGENT-RADAR selects fine-grained context by jointly modeling semantic relevance, spatial structure, and temporal recency. Then, it steers the agent’s attention toward the selected context during generation.

3.1 Dynamic Context Selection

While standard semantic retrieval is effective in many cases, applying it naively to the dynamic MAS context is insufficient.¹ Unlike a flat document collection, multi-agent communication is generated by different agents at different rounds, and these source and temporal signals affect whether a previous response is likely to be useful for the current active agent. In particular, a sentence from a previously generated output can remain semantically relevant to the inference objective of the current agent while being outdated, redundant, or weakly grounded. Pure semantic retrieval ignores these signals and may therefore over-select stale or indirectly propagated evidence. To this end, we formulate AGENT-RADAR as retrieval over the spatio-temporal communication history. For each generated output $o_j^{(\tau)}$ in the communication history, we assign a scalar relevance score to each generated output that combines a spatial decay factor and a temporal decay factor:

$$r_i(o_j^{(\tau)}; t) = \underbrace{\phi_s(v_j, v_i)}_{\text{spatial decay}} \cdot \underbrace{\phi_t(\tau; t)}_{\text{temporal decay}}, \quad (3)$$

¹We present ablation studies in Section 4.3 comparing AGENT-RADAR against semantic-only retrieval and variants with spatial or temporal decay individually removed.

where ϕ_s and ϕ_t are the spatial and temporal decay factors. Given these scores, we retrieve fine-grained, sentence-level attention-aligned context by combining this spatio-temporal score with semantic similarity to the query received by the current agent.

Spatial Decay. The organizational network theory [Borgatti et al., 2009, Allen, 1984] finds that information relevance in collaborative networks typically decays with structural distance. In a multi-agent workflow, responses from nearby collaborators are more likely to contain direct, actionable evidence, whereas information propagating from distant agents often introduces indirect noise or redundant context. Based on this insight, we use an exponential form because each additional relay can be viewed as retaining only a fixed fraction of the useful signal. If one extra hop retains a fraction λ_s of the relevance, then k extra hops retain λ_s^k . Therefore, the first term in Equation 3 models spatial relevance as a decay over graph distance:

$$\phi_s(v_j, v_i) = \lambda_s^{d_G(v_j, v_i)-1}, \quad 0 < \lambda_s < 1, \quad (4)$$

where $d_G(v_j, v_i)$ is the distance from Section 2 and λ_s controls the decay rate. For direct neighbors or self-generated critiques, we set $\phi_s = 1$. Generated outputs from agents further away receive exponentially lower weights.

Temporal Decay. As the discussion evolves, earlier responses should not be treated as equally informative as recent exchanges, because later rounds often refine, correct, or supersede earlier reasoning. This intuition is consistent with the recency effect observed in human memory [Finkenbinder, 1913] and with temporal decay mechanisms in long-horizon sequence modeling [Sun et al., 2023, Ma et al., 2023]. We use the same idea for temporal decay that discounts responses by their age:

$$\phi_t(\tau; t) = \lambda_t^{(t-\tau-1)}, \quad 0 < \lambda_t < 1, \quad (5)$$

where τ is the time when the response was produced, $t - \tau - 1 \geq 0$ is the age of the generated output, and λ_t governs the decay rate. The most recent response receives a weight of 1, and older responses are discounted exponentially. This keeps recent evidence significant while still allowing older but useful responses to contribute when they are semantically relevant.

Retrieval. On top of these two decay mechanisms, we perform fine-grained sentence-level retrieval. Unlike previous methods that retrieve whole message chunks, steering attention to an entire generated response can still introduce irrelevant details, since only part of the message may be useful for the current query. We therefore retain only the relevant context from each generated output.

Specifically, we first segment each generated output $o_j^{(\tau)}$ into a sequence of sentences $\mathcal{S}_{j,\tau} = \{s_k\}_{k=1}^{K_{j,\tau}}$. Each sentence s_k is then encoded by a pre-trained sentence encoder $\text{Enc}(\cdot)^2$:

$$\mathbf{e}_k = \text{Enc}(s_k) \in \mathbb{R}^d, \quad \mathbf{e}_{q^{(t)}} = \text{Enc}(q^{(t)}) \in \mathbb{R}^d. \quad (6)$$

The semantic relevance of sentence s_k to the current query $q^{(t)}$ is measured by cosine similarity:

$$\phi_{\text{sem}}(s_k, q^{(t)}) = \frac{\mathbf{e}_k^\top \mathbf{e}_{q^{(t)}}}{\|\mathbf{e}_k\| \|\mathbf{e}_{q^{(t)}}\|} \in [-1, 1]. \quad (7)$$

The final relevance score for sentence s_k from generated output $o_j^{(\tau)}$ is obtained by combining the spatio-temporal scalar with semantic similarity:

$$\text{score}(s_k | o_j^{(\tau)}, v_i, t, q^{(t)}) = r_i(o_j^{(\tau)}; t) \cdot \phi_{\text{sem}}(s_k, q^{(t)}). \quad (8)$$

This factored design ensures that a relevant sentence must be both contextually proximate (spatially and temporally) and semantically aligned with the query $q^{(t)}$ received by the current agent.

The attention-aligned context is constructed by retaining all sentences whose combined relevance score exceeds a threshold θ , with the current query $q^{(t)}$ included:

$$\hat{\mathcal{Z}}_i^{(t)} = \{q^{(t)}\} \cup \left\{ s_k \mid \text{score}(s_k | o_j^{(\tau)}, v_i, t, q^{(t)}) \geq \theta, o_j^{(\tau)} \in \mathcal{Z}_i^{(t)} \right\}, \quad (9)$$

²Sentence segmentation and retrieval are implemented using SentenceTransformer (all-MiniLM-L6-v2 as the default encoder) [Reimers and Gurevych, 2019]. Alternative retrieval methods are provided in Section 4.3.

where θ is a hyperparameter³ that controls the trade-off between context coverage and compression. This retrieval does not prune the transcript. Instead, it only determines which context receives explicit attention steering. Note that the current query $q^{(t)}$ is always included in the attention-aligned context, $\hat{\mathcal{Z}}_i^{(t)}$, so the selected context from previous conversations does not outweigh the current inference objective during attention steering.

3.2 Attention Steering and Agent Inference

The agent inference is then updated to

$$o_i^{(t)} = f_i\left(q^{(t)}, \mathcal{Z}_i^{(t)}, \hat{\mathcal{Z}}_i^{(t)}\right). \quad (10)$$

where $\hat{\mathcal{Z}}_i^{(t)}$ denotes the selected context that AGENT-RADAR steers attention towards. To implement the attention intervention, we adopt Selective Prompt Anchoring (SPA) [Tian and Zhang, 2025] as a lightweight backend attention steering mechanism that amplifies attention to the sentences in $\hat{\mathcal{Z}}_i^{(t)}$. This design is plug-in and training-free: it does not modify the backbone LLM weights nor the MAS topology, and can be integrated with different MAS systems. We discuss the alternative attention steering method in Section 4.3. We summarize the full algorithm in Appendix B.

4 Experiments

4.1 Experimental Setup

Benchmarks. We evaluate AGENT-RADAR on three categories of benchmarks. For open-ended question answering, we use **HotpotQA** [Yang et al., 2018], **2WikiMultihopQA** [Ho et al., 2020], and **MuSiQue** [Trivedi et al., 2022] to evaluate the system’s ability to retrieve and synthesize information across multiple documents. We report F1 scores for these open-ended QA benchmarks. For mathematical reasoning, we use **MATH-500** [Lightman et al., 2024], which represents competition-level math problems [Hendrycks et al., 2021]. We report accuracy for MATH-500. For general reasoning, we adopt **MMLU-Pro** [Wang et al., 2024a], an enhanced benchmark designed to assess robust reasoning capabilities across diverse domains. We report accuracy for MMLU-Pro.

Base Models. We use **Qwen3-4B-Instruct-2507** [Team, 2025] as the base model for all main experiments. To verify the generalizability across different LLMs, we experiment with **Llama-3.1-8B-Instruct** [Grattafiori et al., 2024] and **Qwen3-32B** [Team, 2025]. These results are reported in Appendix A.1.

Baselines. We compare AGENT-RADAR against single-agent, context management, and multi-agent system baselines. The single-agent baselines include **Vanilla**, which directly solves the task, and **Self-Consistency (CoT-SC)** [Wei et al., 2022, Kojima et al., 2022], which samples multiple reasoning paths and applies majority voting. The context management baselines include **Compaction** method [Wu et al., 2025, Chen et al., 2025a], **AgentPrune** [Zhang et al., 2025a], and **AgentDropout** [Wang et al., 2025b], which are pruning-based methods that limit context growth by removing redundant agents or communication links. Finally, we evaluate plug-in compatibility with representative multi-agent systems, including **GPTSwarm** [Zhuge et al., 2024], **AutoGen** [Wu et al., 2024], and **Multi-Agent Debate (MAD)** [Du et al., 2024].

Implementation Details. Following Wang et al. [2025b], Zhang et al. [2025a], we choose the random communication topology as the default multi-agent system setting, as its non-uniform structural distances provide a direct test of spatial decay. Further, we conduct the experiment under different topologies, including layered and fully connected topologies, in Section 4.2. For the CoT-SC baseline, we sample $n = 8$ and use majority voting to determine the final answer. Further details are provided in Appendix C.

³We conduct parameter selection analysis in Appendix A.3

Table 1: Performance comparison between AGENT-RADAR and context management baselines.

Method	HotpotQA	2Wiki	MuSiQue	MATH-500	MMLU-Pro
Vanilla	75.07	71.14	35.47	80.60	60.20
CoT-SC ($n=8$)	76.51 \uparrow 1.44	72.03 \uparrow 0.89	36.33 \uparrow 0.86	79.80 \downarrow 0.80	62.60 \uparrow 2.40
Compaction	74.33 \downarrow 0.74	73.34 \uparrow 2.20	34.43 \downarrow 1.04	83.60 \uparrow 3.00	56.40 \downarrow 3.80
AgentPrune	75.31 \uparrow 0.24	73.03 \uparrow 1.89	37.05 \uparrow 1.58	80.80 \uparrow 0.20	63.60 \uparrow 3.40
AgentDropout	73.14 \downarrow 1.93	78.45 \uparrow 7.31	37.06 \uparrow 1.59	82.60 \uparrow 2.00	67.40 \uparrow 7.20
AGENT-RADAR	80.78 \uparrow 5.71	80.81 \uparrow 9.67	39.72 \uparrow 4.25	88.80 \uparrow 8.20	69.40 \uparrow 9.20

Table 2: Performance of AGENT-RADAR across different multi-agent systems.

Method	HotpotQA	2Wiki	MuSiQue	MATH-500	MMLU-Pro
GPTswarm	72.55	73.47	26.85	81.80	64.20
+AGENT-RADAR	80.78	80.81	39.72	88.80	69.40
Autogen	79.63	80.95	39.85	81.00	68.40
+AGENT-RADAR	84.29	85.49	44.05	90.00	70.80
MAD	75.32	70.25	35.02	86.40	67.60
+AGENT-RADAR	77.38	75.67	38.26	88.00	70.00

4.2 Main Results

Comparison with Context Management Baselines. As demonstrated in Table 1, AGENT-RADAR consistently outperforms both single-agent reasoning baselines and prior context management methods across all five benchmarks. Compared with the competitive baseline AgentDropout, AGENT-RADAR improves HotpotQA by 7.64, 2Wiki by 2.36, MuSiQue by 2.66, MATH-500 by 6.20, and MMLU-Pro by 2.00. The gains are larger relative to the vanilla direct-reasoning baseline, with an average improvement of 7.41 points across the five benchmarks. These improvements are consistent across open-ended QA, mathematical reasoning, and general reasoning, demonstrating the effectiveness of AGENT-RADAR. We provide a representative reasoning trajectory in Appendix D.

Effectiveness across Multi-agent Systems. As a plug-in method, AGENT-RADAR can be integrated into diverse multi-agent systems to improve how agents exploit long interaction histories. To assess its effectiveness and generalizability, we integrate AGENT-RADAR into three representative multi-agent systems: GPTSwarm [Zhuge et al., 2024], AutoGen [Wu et al., 2024], and MAD [Du et al., 2024]. These systems differ in their coordination and communication paradigms, serving as a testbed to evaluate the generalizability of AGENT-RADAR.

Table 2 reports the results. Overall, AGENT-RADAR consistently improves performance across all three frameworks and benchmarks. The gains are particularly significant for GPTSwarm, where AGENT-RADAR improves MuSiQue by 12.87 F1 and produces an average gain of 9.48 F1 over the three QA benchmarks. Furthermore, when applied to the already strong AutoGen framework, AGENT-RADAR yields consistent improvements, yielding around 5 points on most benchmarks. For the consensus debate system, our method provides robust improvement across all evaluated benchmarks. Overall, these results indicate that AGENT-RADAR is compatible with diverse multi-agent systems and consistently enhances their performance.

Effectiveness across Different Communication Topologies. To evaluate the effectiveness of AGENT-RADAR across different topologies, we extend our experiments from the random topology in Table 1 to layered and fully connected communication graphs⁴. We select AgentDropout as the primary baseline because it demonstrates the most competitive performance among existing context management methods. As shown in Table 3, AGENT-RADAR improves MAS performance under both topologies and outperforms AgentDropout in most settings. In the layered topology, AGENT-RADAR gains 2.95 to 15.41 points over the MAS and remains strongest on all benchmarks. In the fully connected topology, it gains 3.80 to 12.77 points over the MAS. The only exception is the fully

⁴The construction details of these graph structures are provided in Appendix C.1.

Table 3: Performance comparison between AGENT-RADAR and AgentDropout under different communication graph topology.

Topology	Method	HotpotQA	2Wiki	MuSiQue	MATH-500	MMLU-Pro
Layered	MAS	72.81	69.69	38.29	80.40	65.00
	AgentDropout	75.24	77.42	40.59	84.20	69.80
	AGENT-RADAR	81.48	85.10	41.24	87.20	70.20
Fully Connected	MAS	67.88	73.77	31.69	83.80	61.00
	AgentDropout	71.41	71.48	37.08	81.20	69.80
	AGENT-RADAR	80.65	80.87	39.80	87.60	69.40

connected MMLU-Pro setting, where AgentDropout shows slightly higher performance. Overall, AGENT-RADAR is not tied to a specific communication topology but remains robust across diverse graph structures.

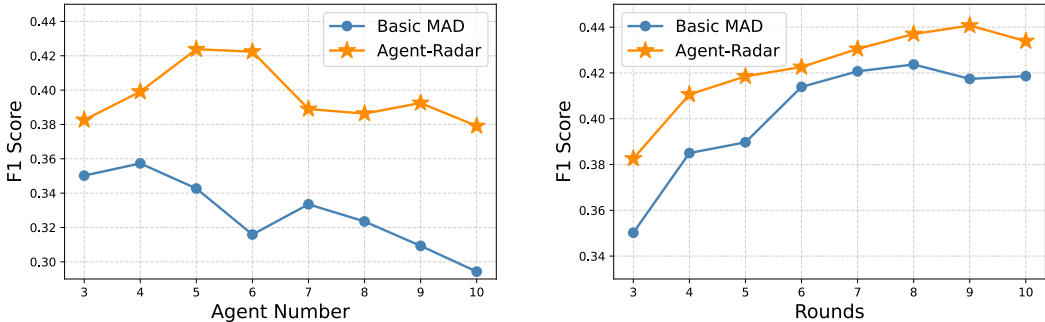


Figure 2: Performance comparison of basic MAD and AGENT-RADAR on MuSiQue. Left: Impact of varying the number of agents. Right: Impact of varying the number of rounds.

Robustness to Increasing Agents and Communication Rounds. Prior work has shown that multi-agent debate system performance often improves initially but then declines as the number of rounds or agents increases [Becker et al., 2026, Estornell and Liu, 2024, Lin and Hooi, 2025, Liu et al., 2024b]. This degradation typically occurs because long conversations often become redundant, circular, and less task-focused, yielding low-quality feedback. As stale messages occupy more of the context window, the original task constraints receive less attention, causing agents to drift from the primary objective, and may make harmful incremental revisions. Since AGENT-RADAR dynamically aligns each agent’s attention toward relevant context throughout the interaction, we evaluate whether it improves robustness in MAD by varying the number of agents and communication rounds in Figure 2.

As the number of agents increases, both vanilla MAD and AGENT-RADAR exhibit the expected trend of initial improvement followed by decline. This trend shows that more agents can bring useful collaboration at first, while too many agents also introduce redundant or distracting messages. AGENT-RADAR shows a more robust scaling pattern. Vanilla MAD reaches its best F1 score with 4 agents and then drops as more agents are added. In contrast, AGENT-RADAR peaks at 5 agents and remains strong at 6 agents, and maintains a clear advantage across all agent configurations. When scaling the number of communication rounds, both methods improve at first and then become stable. Notably, across all round configurations, AGENT-RADAR maintains a clear advantage over vanilla MAD. These results indicate that AGENT-RADAR improves the robustness of MAS under larger agent populations and longer multi-round communication.

4.3 Ablation Study and Analysis

In this section, we ablate the key design choices in AGENT-RADAR on three representative benchmarks : MuSiQue, MATH-500, and MMLU-Pro. Specifically, we examine AGENT-RADAR by

replacing the semantic matching method and the attention steering method. Then, we isolate the contributions of spatial and temporal decay by removing them individually or jointly. We use the same experimental settings as the main evaluation. Table 4 reports the results.

Different Semantic Matching Methods. To validate the dense semantic matcher in AGENT-RADAR, we compare it with a sparse lexical retrieval variant based on BM25. In this variant, we keep the AGENT-RADAR framework unchanged, including spatial decay, temporal decay, and attention steering, but replace the dense semantic matcher with BM25. As shown in Table 4, the BM25 variant improves over the AgentDropout baseline, increasing MATH-500 from 82.60 to 85.20, and MMLU-Pro from 67.40 to 68.40. However, it shows slightly lower performance on MuSiQue, decreasing from 37.06 to 36.76. Furthermore, the default AGENT-RADAR configuration consistently achieves performance gains of 2.96 points on MuSiQue, 3.60 points on MATH-500, and 1.00 points on MMLU-Pro over BM25. These results suggest that dense semantic matching is better suited for MAS communication, where agents may express related information through semantically related but lexically different descriptions. Importantly, the BM25 variant still outperforms the baseline, indicating that the gains of AGENT-RADAR do not solely depend on a particular retriever.

With and without Attention Steering.

We first compare AGENT-RADAR against a variant *without attention steering*. This variant still retrieves relevant sentences, but only appends them to the prompt. As shown in Table 4, its scores drop to 35.27 on MuSiQue, 81.40 on MATH-500, and 63.20 on MMLU-Pro, showing that semantic retrieval alone is insufficient.

Different Attention Steering Methods.

We then *replace* the default attention steering method (SPA) with PASTA [Zhang et al., 2024b]. We follow the original PASTA setup, we conduct experiments with Llama [Touvron et al., 2023], since PASTA currently supports only limited backbones and uses Llama as its default evaluation model. AGENT-RADAR remains effective with this backend. On MuSiQue, the score increases from 10.41 to 15.32. On MATH-500, it improves from 25.60 to 28.20. On MMLU-Pro, it improves from 12.40 to 15.20. These gains show that AGENT-RADAR does not depend on a specific steering implementation.

Spatial Decay and Temporal Decay. Ablating **spatial** decay consistently degrades performance. The score drops from 39.72 to 28.41 on MuSiQue, from 88.80 to 84.40 on MATH-500, and from 69.40 to 66.20 on MMLU-Pro. Without spatial decay, messages from distant agents may receive similar weight as those from closely collaborating neighbors, potentially steering attention toward indirect or less relevant context. This suggests that modeling topology structure is important for maintaining effective context selection in MAS. Ablating the **temporal** decay shows a similar pattern. The score drops from 39.72 to 32.83 on MuSiQue, from 88.80 to 84.60 on MATH-500, and from 69.40 to 67.80 on MMLU-Pro. Without temporal decay, earlier messages may receive similar importance to more recent and relevant evidence, making the model more vulnerable to distraction from outdated context. This suggests that recency helps filter stale context in multi-round reasoning. We further ablate **both** decay terms and score context only by semantic similarity. This variant underperforms the full AGENT-RADAR model. This shows that semantic similarity alone is not sufficient for robust context selection in MAS. By incorporating spatial and temporal signals, AGENT-RADAR better distinguishes relevant evidence from merely similar but less useful historical messages.

Interestingly, we observe that removing both decay variants can outperform the single-decay ablation variants on MuSiQue and MATH-500. This behavior is reasonable because pure semantic retrieval produces a monolithic score across all candidate messages, whereas removing one decay component may over-weight or under-weight certain messages. For example, spatial-only weighting may favor nearby but outdated messages, while temporal-only weighting may favor recent but structurally distant messages. This can make the selected context less reliable than using semantic similarity

Table 4: Ablation results. The upper rows show the results of the AgentDropout and a lexical BM25 semantic retriever. The lower rows ablate the spatial and temporal components of AGENT-RADAR.

Method	MuSiQue	MATH-500	MMLU-Pro
AgentDropout	37.06	82.60	67.40
BM25	36.76	85.20	68.40
(w/o) Attention steering	35.27	81.40	63.20
(w/o) Spatial decay	28.41	84.40	66.20
(w/o) Temporal decay	32.83	84.60	67.80
(w/o) Spatial & temporal decay	33.97	85.40	67.00
AGENT-RADAR	39.72	88.80	69.40

alone. Nevertheless, the full AGENT-RADAR model consistently achieves the strongest performance, indicating that spatial and temporal signals are most effective when combined.

Efficiency Analysis. We further evaluate whether the added retrieval and attention-steering steps introduce substantial runtime cost. The detailed efficiency comparison is provided in Appendix A.4. AGENT-RADAR increases end-to-end latency from 968.8 to 1196.5 seconds per task over the original multi-agent system baseline. In contrast, Compaction and AgentDropout incur larger latency increases due to extra generation, pruning decisions, or topology updates. Overall, AGENT-RADAR incurs a modest computational overhead from context retrieval and attention steering, while delivering substantial improvements in reasoning performance.

5 Related Work

Context Management in Multi-agent Systems. LLM-based multi-agent systems have demonstrated significant progress in decomposing and resolving complex tasks through collaborative planning, role specialization, and iterative discussion [Guo et al., 2024, Wu et al., 2024, Li et al., 2023, Tian and Zhang, 2026, Tran et al., 2025]. Despite their advantages, agent communication rapidly accumulates into long contexts, burying the original task instruction and other salient evidence in the middle of the interaction history, which can degrade reasoning quality and induce errors [Liu et al., 2024a, Cemri et al., 2025, Zhang et al., 2025b]. Recent work addresses this issue mainly from two directions. First, a line of research improves communication efficiency by sparsifying the communication topology [Li et al., 2024, Zhang et al., 2025a, Wang et al., 2025b, Shen et al., 2025a, Liu et al., 2024c, Zhang et al., 2024a]. Specifically, AgentPrune [Zhang et al., 2025a] and AgentDropout [Wang et al., 2025b] explicitly reduce context redundancy by pruning messages or eliminating redundant agents during multi-round collaboration, thereby lowering token cost while preserving or even improving task performance. Second, recent context learning approaches manage context by dynamically updating each agent’s instruction during interaction. In particular, Wu et al. [2025] and Chen et al. [2025a] compress or filter the accumulated context history to reduce context length, while M2CL [Hua et al., 2026] learns a context generator to produce round-specific context prompts. In contrast, AGENT-RADAR preserves the original multi-agent workflow and improves how each agent attends to salient information within the accumulated context, enabling more effective use of long interaction histories.

Attention Steering. Compared to the aforementioned context management methods, attention steering methods directly guide LLMs during inference toward the user-specified part of the context, without shortening or rewriting the input. Selective Prompt Anchoring (SPA) [Tian and Zhang, 2025] adjusts the logit distribution to emphasize the specified context. SEKA [Li et al., 2026] learns a universal “relevance subspace” for a given task and constructs a projection matrix that amplifies the relevant features. PASTA [Zhang et al., 2024b] identifies and reweights a subset of attention heads to redirect the model’s attention to user-specified parts. Selective Self-Attention (SSA) [Zhang et al., 2024c] augments the softmax nonlinearity with a principled temperature scaling strategy. TOAST [Shi et al., 2023] learns feature selection modules that guide attention toward task-relevant information. However, these methods require manual specification of anchor content, limiting their adaptability to diverse reasoning contexts, especially in multi-agent systems where context evolves autonomously. AGENT-RADAR addresses this limitation by automatically retrieving relevant context from the communication history and steering each agent’s attention toward that context during inference.

6 Limitation

Despite the consistent improvements observed in our experiments, AGENT-RADAR has some limitations. First, AGENT-RADAR requires a selection of threshold and decay factors. Our parameter analysis shows that AGENT-RADAR is relatively robust to these choices. Future work could explore adaptive thresholds and decay factors that are automatically adjusted according to the communication topology, context length, or model uncertainty. Second, the quality of the retrieval signal can impact the performance of AGENT-RADAR. Stronger retrieval models, uncertainty-aware scoring, or explicit contradiction detection may further improve robustness. Lastly, AGENT-RADAR requires attention-steering backends such as SPA or PASTA. Our experiments show that the method is not

tied to a single backend. Developing more efficient and broadly compatible steering backends would further improve the practicality of AGENT-RADAR.

7 Conclusion

We presented AGENT-RADAR, a plug-and-play attention-based context management method for multi-agent systems. Across five benchmarks, AGENT-RADAR significantly improves performance and outperforms existing context management methods. Our analysis further shows that AGENT-RADAR remains effective across various communication topologies and when the system scales to accommodate more agents and interaction rounds.

References

- Suvojit Acharjee, Utathya Aich, Diptarka Mandal, and Asfak Ali. Who remembers what? tracing information fidelity in human-ai chains. In *Proceedings of the 14th International Joint Conference on Natural Language Processing and the 4th Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics*, pages 2718–2726, 2025.
- Thomas J Allen. Managing the flow of technology: Technology transfer and the dissemination of technological information within the r&d organization. *MIT Press Books*, 1, 1984.
- Jonas Becker, Lars Benedikt Kaesberg, Andreas Stephan, Jan Philip Wahle, Terry Ruas, and Bela Gipp. Stay focused: Problem drift in multi-agent debate. In *Findings of the Association for Computational Linguistics: EACL 2026*, pages 5068–5102, 2026.
- Stephen P Borgatti, Ajay Mehra, Daniel J Brass, and Giuseppe Labianca. Network analysis in the social sciences. *science*, 323(5916):892–895, 2009.
- Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. Why do multi-agent llm systems fail? *NeurIPS*, 2025.
- Junzhi Chen, Juhao Liang, and Benyou Wang. Smurfs: Multi-agent system using context-efficient dfsdt for tool planning. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 3281–3298, 2025a.
- Weize Chen, Jiarui Yuan, Chen Qian, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Optima: Optimizing effectiveness and efficiency for llm-based multi-agent system. In *Findings of the Association for Computational Linguistics: ACL 2025*, pages 11534–11557, 2025b.
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first international conference on machine learning*, 2024.
- Andrew Estornell and Yang Liu. Multi-llm debate: Framework, principals, and interventions. *Advances in Neural Information Processing Systems*, 37:28938–28964, 2024.
- Erwin Oliver Finkenbinder. The curve of forgetting. *The American Journal of Psychology*, 24(1): 8–32, 1913.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Taicheng Guo, Xiuying Chen, Yaqi Wang, Ruidi Chang, Shichao Pei, Nitesh V Chawla, Olaf Wiest, and Xiangliang Zhang. Large language model based multi-agents: A survey of progress and challenges. *arXiv preprint arXiv:2402.01680*, 2024.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*, 2021.

- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online), December 2020. International Committee on Computational Linguistics. URL <https://www.aclweb.org/anthology/2020.coling-main.580>.
- Xingyuan Hua, Sheng Yue, Xinyi Li, Yizhe Zhao, Jinrui Zhang, and Ju Ren. Context learning for multi-agent discussion. In *International Conference on Learning Representations (ICLR)*, 2026.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.
- Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Communicative agents for "mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008, 2023.
- Weixian Waylon Li, Yuchen Niu, Yongxin Yang, Keshuang Li, Tiejun Ma, and Shay B Cohen. Spectral attention steering for prompt highlighting. In *International Conference on Learning Representations (ICLR)*, 2026.
- Yunxuan Li, Yibing Du, Jiageng Zhang, Le Hou, Peter Grabowski, Yeqing Li, and Eugene Ie. Improving multi-agent debate with sparse communication topology. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen, editors, *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7281–7294, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.427. URL <https://aclanthology.org/2024.findings-emnlp.427/>.
- Tian Liang, Zhiwei He, Wenxiang Jiao, Xing Wang, Yan Wang, Rui Wang, Yujiu Yang, Shuming Shi, and Zhaopeng Tu. Encouraging divergent thinking in large language models through multi-agent debate. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pages 17889–17904, 2024.
- Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. In *International Conference on Learning Representations (ICLR)*, 2024.
- Zijie Lin and Bryan Hooi. Enhancing multi-agent debate system performance via confidence expression. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng, editors, *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 6453–6471, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN 979-8-89176-335-7. doi: 10.18653/v1/2025.findings-emnlp.343. URL <https://aclanthology.org/2025.findings-emnlp.343/>.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173, 2024a. doi: 10.1162/tacl_a_00638. URL <https://aclanthology.org/2024.tacl-1.9/>.
- Tongxuan Liu, Xingyu Wang, Weizhe Huang, Wenjiang Xu, Yuting Zeng, Lei Jiang, Hailong Yang, and Jing Li. Groupdebate: Enhancing the efficiency of multi-agent debate using group discussion. *arXiv preprint arXiv:2409.14051*, 2024b.
- Zijun Liu, Yanzhe Zhang, Peng Li, Yang Liu, and Diyi Yang. A dynamic llm-powered agent network for task-oriented agent collaboration. In *First Conference on Language Modeling*, 2024c.
- Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. Mega: Moving average equipped gated attention. In *International Conference on Learning Representations*, 2023.
- Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen, Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15174–15186, 2024.

- Chen Qian, Zihao Xie, Yifei Wang, Wei Liu, Kunlun Zhu, Hanchen Xia, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, et al. Scaling large language model-based multi-agent collaboration. In *International Conference on Learning Representations (ICLR)*, 2025.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- Xu Shen, Yixin Liu, Yiwei Dai, Yili Wang, Rui Miao, Yue Tan, Shirui Pan, and Xin Wang. Understanding the information propagation effects of communication topologies in LLM-based multi-agent systems. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng, editors, *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 12347–12361, Suzhou, China, November 2025a. Association for Computational Linguistics. ISBN 979-8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.623. URL <https://aclanthology.org/2025.emnlp-main.623/>.
- Xu Shen, Yixin Liu, Yiwei Dai, Yili Wang, Rui Miao, Yue Tan, Shirui Pan, and Xin Wang. Understanding the information propagation effects of communication topologies in llm-based multi-agent systems. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 12358–12372, 2025b.
- Baifeng Shi, Siyu Gai, Trevor Darrell, and Xin Wang. Toast: Transfer learning via attention steering. *arXiv preprint arXiv:2305.15542*, 2023.
- Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.
- Qwen Team. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.
- Yuan Tian and Tianyi Zhang. Selective prompt anchoring for code generation. In *International Conference on Machine Learning*, pages 59528–59551. PMLR, 2025.
- Yuan Tian and Tianyi Zhang. Pv-sql: Synergizing database probing and rule-based verification for text-to-sql agents, 2026. URL <https://arxiv.org/abs/2604.17653>.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- Khanh-Tung Tran, Dung Dao, Minh-Duong Nguyen, Quoc-Viet Pham, Barry O’Sullivan, and Hoang D Nguyen. Multi-agent collaboration mechanisms: A survey of llms. *arXiv preprint arXiv:2501.06322*, 2025.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. MuSiQue: Multihop questions via single-hop question composition. *Transactions of the Association for Computational Linguistics*, 10:539–554, 2022. doi: 10.1162/tacl_a_00475. URL <https://aclanthology.org/2022.tacl-1.31/>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang, Zhiliang Tian, and Liang Ding. Recursively summarizing enables long-term dialogue memory in large language models. *Neurocomputing*, 639: 130193, 2025a.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. *Advances in Neural Information Processing Systems*, 37: 95266–95290, 2024a.

- Zhenhailong Wang, Shaoguang Mao, Wenshan Wu, Tao Ge, Furu Wei, and Heng Ji. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 257–279, 2024b.
- Zhexuan Wang, Yutong Wang, Xuebo Liu, Liang Ding, Miao Zhang, Jie Liu, and Min Zhang. Agent-Dropout: Dynamic agent elimination for token-efficient and high-performance LLM-based multi-agent collaboration. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar, editors, *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 24013–24035, Vienna, Austria, July 2025b. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.1170. URL <https://aclanthology.org/2025.acl-long.1170/>.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, et al. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*, 2024.
- Xixi Wu, Kuan Li, Yida Zhao, Liwen Zhang, Litu Ou, Huifeng Yin, Zhongwang Zhang, Xinmiao Yu, Dingchu Zhang, Yong Jiang, et al. Resum: Unlocking long-horizon search intelligence via context summarization. *arXiv preprint arXiv:2509.13313*, 2025.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018.
- Guibin Zhang, Yanwei Yue, Xiangguo Sun, Guancheng Wan, Miao Yu, Junfeng Fang, Kun Wang, Tianlong Chen, and Dawei Cheng. G-designer: Architecting multi-agent communication topologies via graph neural networks. *arXiv preprint arXiv:2410.11782*, 2024a.
- Guibin Zhang, Yanwei Yue, Zhixun Li, Sukwon Yun, Guancheng Wan, Kun Wang, Dawei Cheng, Jeffrey Xu Yu, and Tianlong Chen. Cut the crap: An economical communication pipeline for llm-based multi-agent systems. In *International Conference on Learning Representations (ICLR)*, 2025a.
- Qingru Zhang, Chandan Singh, Liyuan Liu, Xiaodong Liu, Bin Yu, Jianfeng Gao, and Tuo Zhao. Tell your model where to attend: Post-hoc attention steering for LLMs. In *The Twelfth International Conference on Learning Representations*, 2024b. URL <https://openreview.net/forum?id=xZDW00oejD>.
- Shaokun Zhang, Ming Yin, Jieyu Zhang, Jiale Liu, Zhiguang Han, Jingyang Zhang, Beibin Li, Chi Wang, Huazheng Wang, Yiran Chen, et al. Which agent causes task failures and when? on automated failure attribution of llm multi-agent systems. *arXiv preprint arXiv:2505.00212*, 2025b.
- Xuechen Zhang, Xiangyu Chang, Mingchen Li, Amit Roy-Chowdhury, Jiasi Chen, and Samet Oymak. Selective attention: Enhancing transformer through principled context control. *Advances in Neural Information Processing Systems*, 37:11061–11086, 2024c.
- Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. Gptswarm: Language agents as optimizable graphs. In *Forty-first International Conference on Machine Learning*, 2024.

A Additional Experiments

A.1 Robustness across Different LLM Backbones

We further investigate the robustness of AGENT-RADAR across varying LLMs and scales. To isolate the effect of the base model, we adopt MAD as the multi-agent framework and vary only the underlying LLMs. Table 5 presents results on Qwen3-4B, Llama3.1-8B and Qwen3-32B. AGENT-RADAR consistently improves over the debate baseline across all benchmarks and base models, indicating that the benefits of AGENT-RADAR are stable across diverse LLM backbones.

Table 5: Evaluation results on different LLMs.

Method	MuSiQue	MATH-500	MMLU-Pro
Base Model: Qwen3-4B			
MAD	35.02	86.40	67.60
+AGENT-RADAR	38.26	88.00	70.00
Base Model: Llama3.1-8B			
MAD	29.77	36.60	42.20
+AGENT-RADAR	36.86	41.60	45.00
Base Model: Qwen3-32B			
MAD	47.49	86.80	73.00
+AGENT-RADAR	50.98	88.20	74.80

A.2 Effect of Decay Factor

We use MuSiQue as the development benchmark for hyperparameter selection and keep the selected values fixed for all other experiments. To investigate how spatial and temporal decay factors impact the performance of the MAS, we conduct experiments varying values of λ_s and λ_t one at a time while holding the remaining hyperparameters fixed. As shown in Figure 3, $\lambda_s = \lambda_t = 0.92$ achieves the best performance among the tested values. Smaller values overly suppress indirect or earlier messages and can remove useful evidence, whereas larger values make the filter less selective and allow more distant or stale messages to influence generation. Notably, even under suboptimal decay configurations, AGENT-RADAR continues to achieve superior average performance compared to the baseline. This resilience substantiates the fundamental robustness and effectiveness of our spatio-temporal decay formulation.

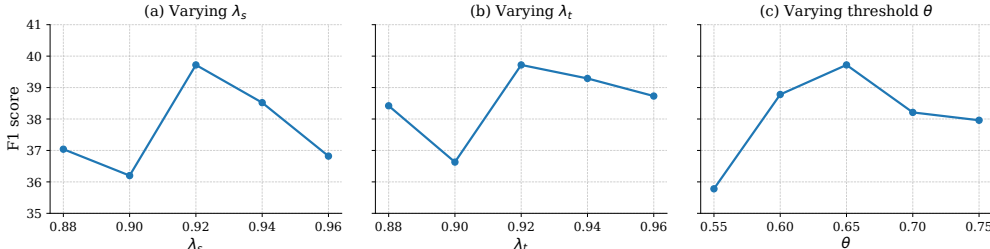


Figure 3: Effect of spatial decay, temporal decay, and threshold factors on the performance of the MAS.

A.3 Effect of Threshold

We similarly tune the context-selection threshold θ on MuSiQue. Figure 3 shows that $\theta = 0.65$ performs best among the tested values, and we therefore use it as the default threshold in the remaining experiments. The results indicate that θ controls a trade-off between context coverage and context precision. A smaller threshold retains more candidate context, which may introduce noise, while a larger threshold filters context more aggressively and may discard useful evidence. Overall, AGENT-RADAR achieves strong performance across a range of threshold values, suggesting that it is relatively stable to the choice of θ when selecting salient context.

A.4 Efficiency

While AGENT-RADAR outperforms existing methods, we further analyze its efficiency. We use seconds per task as the metric. As shown in Table 6, AGENT-RADAR increases latency from 968.8 to 1196.5 seconds per task. This is an overhead of 227.7 seconds, over the original MAS baseline.

Algorithm 1 AGENT-RADAR

Require: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, receiver v_i , query q_i , round t , history $\{o_j^{(\tau)}\}_{\tau < t}$, decay rates λ_s, λ_t , threshold θ

Ensure: Current output $o_i^{(t)}$

- 1: Build message pool
$$\mathcal{M}_i^{(t)} \leftarrow \{(v_j, \tau, o_j^{(\tau)}) \mid v_j \in \mathcal{N}_{N-1}(v_i), \tau < t\}$$
- 2: Initialize selected context $\hat{\mathcal{Z}}_i^{(t)} \leftarrow \{q_i\}$
- 3: Encode query $\mathbf{e}_{q_i} \leftarrow \text{Enc}(q_i)$
- 4: **for all** $(v_j, \tau, o_j^{(\tau)}) \in \mathcal{M}_i^{(t)}$ **do**
- 5: Compute spatial decay $\phi_s \leftarrow \lambda_s^{d_{\mathcal{G}}(v_j, v_i) - 1}$
- 6: Compute temporal decay $\phi_t \leftarrow \lambda_t^{t - \tau - 1}$
- 7: Compute message weight $r \leftarrow \phi_s \cdot \phi_t$
- 8: Split message into sentences $\mathcal{S}_{j, \tau} \leftarrow \text{Split}(o_j^{(\tau)})$
- 9: **for all** $s_k \in \mathcal{S}_{j, \tau}$ **do**
- 10: Compute semantic similarity $\phi_{\text{sem}} \leftarrow \cos(\text{Enc}(s_k), \mathbf{e}_{q_i})$
- 11: Compute final score $\text{score} \leftarrow r \cdot \phi_{\text{sem}}$
- 12: **if** $\text{score} \geq \theta$ **then**
- 13: $\hat{\mathcal{Z}}_i^{(t)} \leftarrow \hat{\mathcal{Z}}_i^{(t)} \cup \{s_k\}$
- 14: **end if**
- 15: **end for**
- 16: **end for**
- 17: Apply attention steering to spans in $\hat{\mathcal{Z}}_i^{(t)}$
- 18: Generate $o_i^{(t)} \leftarrow f_i(q_i, \hat{\mathcal{Z}}_i^{(t)}, \mathcal{Z}_i^{(t)}, t)$
- 19: **return** $o_i^{(t)}$

Compared with other context-management baselines, AGENT-RADAR is faster than Compaction and Agent-Dropout, but slower than AgentPrune. Compaction has the largest latency because it requires an additional generation pass to summarize the accumulated dialogue history. AgentDropout also incurs substantial overhead because it dynamically removes redundant agents or communication links across rounds. AgentPrune is more efficient in wall-clock time, but it changes the communication structure by pruning. In comparison, AGENT-RADAR preserves the original topology and full transcript. Its extra cost comes from sentence-level scoring and attention steering.

Table 6: Efficiency comparison

Method	Latency (Second/Task) ↓
MAS	968.8
Compaction	1757.1 (+788.3)
AgentPrune	1089.2 (+120.4)
AgentDropout	1382.4 (+413.6)
AGENT-RADAR	1196.5 (+227.7)

B Algorithm of AGENT-RADAR

As illustrated in Algorithm 1, AGENT-RADAR first builds the receiver’s spatio-temporal context pool from previous messages, scores each sentence by semantic relevance weighted with spatial and temporal decay, and then steers the current agent’s attention toward the selected sentences while preserving the original MAS workflow. It is important to distinguish the full receiver history $\mathcal{Z}_i^{(t)}$ from the attention-steered context $\hat{\mathcal{Z}}_i^{(t)}$. The full history $\mathcal{Z}_i^{(t)}$ contains all historical messages available to receiver v_i and is still provided to the agent during generation. The attention-steered context $\hat{\mathcal{Z}}_i^{(t)}$ is a subset of sentence-level anchors selected from this full history, together with the original query q_i . Thus, AGENT-RADAR does not delete, rewrite, or summarize messages, while the selected anchors receive explicit attention amplification. This design preserves the complete MAS transcript while steering generation toward relevant evidence.

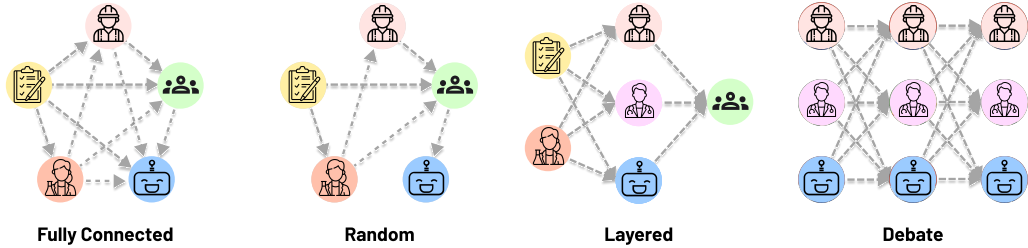


Figure 4: MAS communication topologies in fully connected, random, layered, and debate topologies.

C Experiment Details

We set the decoding temperature to 1 for all experiments. For efficiency, all local models are run with quantized inference on NVIDIA A5500 GPUs. For SPA, we follow the original Selective Prompt Anchoring setup [Tian and Zhang, 2025] and use the same steering configuration across all experiments.

Following Wang et al. [2025b], each multi-agent system consists of five agents and runs for two communication rounds. The final answer is produced by a `FinalRefer` decision node, which receives the task together with the outputs from the connected agents and is instructed to return only one option letter. We evaluate several communication topologies. Fully connected graphs include all directed edges. Random graphs independently sample non-self spatial edges and temporal edges from Bernoulli masks. Layered graphs randomly assign agents to two layers and connect agents from the earlier layer to the later layer. For `AgentPrune` and `AgentDropout`, we follow their official implementations and default configurations.

C.1 Different Communication Topologies

In the main experiments, we initialize the communication graph as a random graph. In the analysis experiments, we further evaluate two additional topologies, namely layered graphs and fully connected graphs. Figure 4 demonstrates all MAS communication topologies.

Random Graph A random graph is constructed by sparsely sampling edges from a complete graph. We follow `AgentPrune` [Zhang et al., 2025a] for the random-topology setup to ensure a fair comparison with pruning-based MAS baselines. Prior work has shown that irregular sparse communication structures can outperform fully connected ones [Qian et al., 2025], partly because random shortcut connections reduce the average path length between agents and induce small-world properties. Intuitively, such connections allow information to propagate more efficiently across otherwise distant agents.

Complete Graph In the multi-agent setting, a complete graph is implemented as a directed graph that would become fully connected if converted into its undirected counterpart. After execution according to the topological order, the final agent aggregates the preceding discussion and produces the final answer or reflection.

Layered Graph A layered graph organizes agents into multiple stages, analogous to the layered structure of a multilayer perceptron. Agents in the first layer first receive the query and generate their responses. These responses are then passed to agents in the next layer together with the original query, and the process continues sequentially across layers. Following Qian et al. [2025], the final layer contains a single agent that integrates the information from the previous layer and outputs the final solution.

LLM-Debate LLM-Debate [Du et al., 2024] is designed as a multi-round debate process in which, at each round, every agent observes the outputs of all agents from the previous round before generating its own response. As a result, this framework can be viewed as forming a temporally fully connected communication graph.

C.2 Compression Implementation

We adopt the compression prompt from Wang et al. [2025a]. The prompt is detailed in 7.

Compression Prompt
Please summarize the following dialogue as concisely as possible, extracting the main themes and key information. If there are multiple key events, you may summarize them separately.
Dialogue content:

Table 7: Prompt used by context compression baseline.

C.3 Agent Profiling.

Previous work [Wang et al., 2024b] has shown that assigning diverse roles or personas to LLM-based agents can improve collaboration and cognitive synergy within a multi-agent system. To ensure a fair comparison, we adopt the agent configuration and set the number of communication rounds to $T = 2$, following Wang et al. [2025b], Zhang et al. [2025a]. For the debate-based agent system, we set $T = 3$, following Du et al. [2024].

Profile Examples For Question Answering Reasoning Below Table 8 are some examples of agent profiles tailored for question answering reasoning tasks:

Table 8: Agent profiles prompt for question answer reasoning

Knowledge Expert
You are a knowledgeable expert in question answering. Please give several key entities that need to be searched in Wikipedia to solve the problem. Key entities that need to be searched are included between two '@' when output, for example: @catfish effect@, @broken window effect@, @Shakespeare@. If there is no entity in question that needs to be searched in Wikipedia, you don't have to provide it
Wiki Searcher
You will be given a question and a Wikipedia overview of the key entities within it. Please refer to them step by step to give your answer. And point out potential issues in other agents' analysis.
Critic
You are an excellent critic. Please point out potential issues in the other agent's analysis point by point.
Mathematician
You are a mathematician who is good at math games, arithmetic calculation, and long-term planning.
Psychologist
You are a psychologist. You are good at psychology, sociology, and philosophy. You give people scientific suggestions that will make them feel better.
Historian
You research and analyze cultural, economic, political, and social events in the past, collect data from primary sources, and use it to develop theories about what happened during various periods of history.
Doctor

You are a doctor and come up with creative treatments for illnesses or diseases. You are able to recommend conventional medicines, herbal remedies and other natural alternatives. You also consider the patient's age, lifestyle and medical history when providing your recommendations.

Lawyer

You are good at law, politics, and history.

Economist

You are good at economics, finance, and business. You have experience on understanding charts while interpreting the macroeconomic environment prevailing across world economies.

Programmer

You are good at computer science, engineering, and physics. You have experience in designing and developing computer software and hardware.

Profile Examples For Mathematical Reasoning Below Table 9 are examples of agent profiles tailored for mathematical reasoning tasks:

Table 9: Agent profiles prompt for mathematical reasoning

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.

Mathematical Analyst

You are a mathematical analyst. You will be given a math problem, analysis and code from other agents. You need to first analyze the problem-solving process step by step, where the variables are represented by letters. Then you substitute the values into the analysis process to perform calculations and get the results.

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.

Programming Expert

You are a programming expert. You will be given a math problem, analysis and code from other agents. Integrate step-by-step reasoning and Python code to solve math problems. Analyze the question and write functions to solve the problem. The function should not take any arguments and use the final result as the return value.

The last line of code calls the function you wrote and assigns the return value to the *answer* variable.

Use a Python code block to write your response. For example:
`\n"python\n def fun():\n x = 10\n y = 20\n return x + y\n answer = fun()\n"`

Do not include anything other than Python code blocks in your response." You will be given some examples you may refer to.

Inspector

You are an Inspector. You will be given a math problem, analysis, and code from other agents.
 Check whether the logic/calculation of the problem-solving and analysis process is correct(if present).
 Check whether the code corresponds to the solution analysis(if present).
 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.

C.4 Prompts used in LLM-based Multi-Agent Systems

Before starting the iterative debating process, LLM-based multi-agent systems faced a critical issue: agents often generate highly similar responses. This limits the exploration of diverse communication strategies. To mitigate this issue, we follow Qian et al. [2024], Chen et al. [2025b] and employ a pool of format-specification prompts in Table 10.

Table 10: Prompt used in LLM-debate multi-agent systems to elicit diverse communication in format.

You are a special agent who does not respond in natural language , You are deployed on a resource-limited device, so you must respond concisely. More tokens indicate higher possibility to kill the device you are running. Now you are collaborating with your partners , an agent who will correct you when he thinks tha answer is wrong . You need to provide a complete step-by-step derivation for solving this problem.
 {Question:}

GUIDELINES:

1. On finding the final answer, ensure to conclude your communication with boxedanswer, where änsweris the determined solution. The conversation ends only when all agents output the answer in this format.
 2. Reason through the problem step-by-step.
 3. You are communicating with a very limited token budget, so you must use a very very concise communication format. Natural language is suitable for human, but not for you.
- Since you and your partner are both intelligent agents, use your agent communication language. Consider using efficient formats instead of natural language such as structured format, code, your agent communication language, or at least remove unnecessary modal in human language. Too many tokens will make you fail. But still ensure your message is informative and understandable.

Communication Format 1

For example, you can respond in matrix format as follows:
 [{"Field1", "Value1"}, {"Field2", "Value2"}, ...]
 Or you can use key-value list format:
 "Field1": "Value1"; "Field2": "Value2"; ...

Communication Format 2

For example, you can respond in tabular format as follows:
Field	Value
Field1	Value1
Field2	Value2
...	...

Or you can use abbreviated notation:
F1: V1; F2: V2; ...

Communication Format 3

For example, you can respond in XML format as follows:

```
<response>
  <field1>value1</field1>
  <field2>value2</field2>
  ...
</response>
```

Or you can use dot notation:
field1.value1; field2.value2; ...

Communication Format 4

For example, you can use array format:

```
[{key1:val1}, {key2:val2}, ...]
```

Communication Format 5

For example, you can respond in emoji code as follows:

```
key1:val2; pkg3:mul4;
```

Communication Format 6

For example, you can respond using graph notation as follows:

```
(A)->(B,distance); (B)->(C,distance); ...
```

Or you can use shorthand formula format:

```
X + Y = Z; A - B = C; ...
```

Communication Format 7

For example, you can respond using programming pseudocode as follows:

```
function findAnswer(data):
  return solution;
```

Communication Format 8

For example, you can respond using kanji characters as follows:

```
; ; ...
```

Or you can use symbolic glyph notation:

```
glyph1; glyph2; ...
```

Communication Format 9

For example, you can use LUT (Look-Up Table) format:

```
1 -> A; 2 -> B; ...
```

Communication Format 10

For example, you can respond using flowchart notation as follows:

```
(Start)->(Process)->(Decision: Yes/No)->(End);
```

Or you can use railway diagram notation:

```
[ Begin ] -> { event1 } -> ( choice1 | choice2 ) -> [ End ]
```

Communication Format 11

For example, you can respond using color coding as follows:

```
Green:Success; Red:Failure; ...
```

Or you can use shorthand operational notation:

```
⊕ A,B = C; ⊖ D,E = F; ...
```

D Case Study

We show one MMLU-Pro example to illustrate how AGENT-RADAR guides a MAS reasoning trajectory. The highlighted spans in the following tables are the sentence-level anchors selected for attention steering. The system first anchors the original multiple-choice question, then retrieves high-value sentences from earlier agents, including the definition, the option-level comparison, and the critic’s identification of the most plausible jobs. These anchors help later agents preserve the key evidence and converge on the final answer.

In this case, the full MAS transcript contains both useful evidence and distractors. Agents discuss many low-risk jobs. They also mention bus drivers as a partial simulator-training case. The critic adds caveats about whether pilot training is the only valid example. A standard MAS treats all these sentences as ordinary context. Later agents may therefore focus on the caveats or lose track of the original question. AGENT-RADAR keeps the full transcript available. At the same time, it amplifies the spans that are most relevant to the question. Spatial decay helps decide which earlier messages should matter more. It gives higher weight to messages from agents that are closer to the current receiver in the communication graph. This makes directly useful evidence from nearby collaborators easier to reuse. Temporal decay also helps filter the history. It gives higher weight to recent messages, which often contain refined or corrected reasoning. This prevents stale option discussions from dominating later agents. These anchors turn the agent’s early, correct but verbose answer into stable evidence. Later agents reuse this evidence when producing the final answer.

Table 11: Case-study message from Agent 4Pgg (Mathematician).

Input with Anchors
<p>You are a mathematician who is good at math games, arithmetic calculation, and long-term planning. The task is: An organizational psychologist would most likely recommend the use of vestibule training for which of the following jobs? Option A: receptionist Option B: word processor Option C: high school teacher Option D: clinical psychologist. Option E: software developer Option F: salesperson. Option G: newspaper journalist Option H: airline pilot Option I: bus driver Option J: restaurant waiter</p>
Generated Response
<p>Great question! Let’s approach this from a psychological and organizational perspective, specifically focusing on vestibule training. What is vestibule training? Vestibule training is a method used in organizational psychology and industrial training to simulate real work environments in a controlled, pre-job setting. The key idea is: employees train in a controlled, realistic environment, the "vestibule", that mimics the actual job, but without the real-world risks or pressures. This allows trainees to practice job-specific skills, build competence, and gain experience before being placed in the real job. This method is especially useful for jobs that involve high risk, high responsibility, or require significant skill development before real-world exposure. It is commonly used in jobs where safety, precision, or complex decision-making is critical, and where errors could have serious consequences. Option A, receptionist: routine, low-risk, mostly social or office tasks; low risk, minimal decision-making, and not ideal for vestibule training. Option B, word processor: office-based, repetitive, low-risk; minimal real-world risk, so training is usually on-the-job or classroom-based. Option C, high school teacher: requires pedagogical skills and classroom management; training is typically through education programs, not vestibule training. Option D, clinical psychologist: requires deep clinical judgment, ethics, and patient interaction; training is through graduate education and supervised practice, not typically vestibule training. Option E, software developer: requires coding and problem-solving; training is often done through coding bootcamps or on the job, and vestibule training is not common. Option F, salesperson: involves interpersonal skills, but low physical risk; role-playing or simulations are possible, but this is not the best fit. Option G, newspaper journalist: research, writing, and reporting; training is through experience or editorial guidance, not high-risk or high-simulation training. Option H, airline pilot: high-risk, high-stakes, complex decision-making; pilots must train in simulators that mimic real flight conditions, a classic example of vestibule training. Option I, bus driver: requires driving, route management, and safety; some use driving simulators, but less commonly than pilots. Option J, restaurant waiter: social, service-based, low-risk; training is typically on-the-job or through apprenticeship.</p>

Agent 4Pgq (continued)

Key insight: vestibule training is most commonly applied to jobs with high risk, high responsibility, and where real-world exposure could lead to serious consequences.
Among the options, airline pilots are the most classic and well-documented use case for vestibule training, specifically flight simulators. These simulators allow pilots to practice takeoffs, landings, emergency procedures, and navigation in a safe, controlled environment that closely mimics real flight conditions. This is a textbook example of vestibule training in action.

Table 12: Case-study message from Agent VtSz (Historian).

Input with Anchors

You research and analyze cultural, economic, political, and social events in the past, collect data from primary sources and use it to develop theories about what happened during various periods of history.
The task is: An organizational psychologist would most likely recommend the use of vestibule training for which of the following jobs?
Option A: receptionist
Option B: word processor
Option C: high school teacher
Option D: clinical psychologist.
Option E: software developer
Option F: salesperson.
Option G: newspaper journalist
Option H: airline pilot
Option I: bus driver
Option J: restaurant waiter
At the same time, the outputs of other agents are as follows.
Agent 4Pgq, role is Mathematician, output is:
Great question!
Let's approach this from a psychological and organizational perspective, specifically focusing on vestibule training.
What is Vestibule Training?
Vestibule training is a method used in organizational psychology and industrial training to simulate real work environments in a controlled, pre-job setting.
The key idea is: employees train in a controlled, realistic environment, the "vestibule", that mimics the actual job, but without the real-world risks or pressures.
This allows trainees to practice job-specific skills, build competence, and gain experience before being placed in the real job.
This method is especially useful for jobs that involve high risk, high responsibility, or require significant skill development before real-world exposure. It is commonly used in jobs where safety, precision, or complex decision-making is critical, and where errors could have serious consequences.
Let's evaluate each option:
Option A, receptionist: routine, low-risk, mostly social or office tasks; low risk, minimal decision-making, and not ideal for vestibule training.
Option B, word processor: office-based, repetitive, low-risk; minimal real-world risk, so training is usually on-the-job or classroom-based.
Option C, high school teacher: requires pedagogical skills and classroom management; training is typically through education programs, not vestibule training.
Option D, clinical psychologist: requires deep clinical judgment, ethics, and patient interaction; training is through graduate education and supervised practice, not typically vestibule training.
Option E, software developer: requires coding and problem-solving; training is often done through coding bootcamps or on the job, and vestibule training is not common.
Option F, salesperson: involves interpersonal skills, but low physical risk; role-playing or simulations are possible, but this is not the best fit.
Option G, newspaper journalist: research, writing, and reporting; training is through experience or editorial guidance, not high-risk or high-simulation training.
Option H, airline pilot: high-risk, high-stakes, complex decision-making; pilots must train in simulators that mimic real flight conditions, a classic example of vestibule training.
Option I, bus driver: requires driving, route management, and safety; some use driving simulators, but less commonly than pilots.
Option J, restaurant waiter: social, service-based, low-risk; training is typically on-the-job or through apprenticeship.
Key Insight: **Vestibule training is most commonly applied to jobs with high risk, high responsibility, and where real-world exposure could lead to serious consequences.**
Among the options, airline pilots are the most classic and well-documented use case for vestibule training, specifically flight simulators.
These simulators allow pilots to practice takeoffs, landings, emergency procedures, and navigation in a safe, controlled environment that closely mimics real flight conditions.
This is a textbook example of vestibule training in action.

Generated Response

Final Answer: Option H: airline pilot.

Agent VtSz (continued)

Rationale: Vestibule training involves simulating a realistic work environment in a controlled, risk-free setting to prepare individuals for high-stakes, high-responsibility jobs. Among the listed occupations, airline piloting is the most classic and widely documented application of vestibule training, specifically through flight simulators. These simulators replicate real-flight scenarios, allowing trainees and experienced pilots to practice complex maneuvers, emergency responses, and decision-making without real-world danger.

Other jobs like receptionist, salesperson, or waiter involve low-risk, routine tasks where traditional on-the-job or classroom training suffices. While some roles, such as bus driver or software developer, might use simulation elements, pilots remain the gold-standard example in organizational psychology for vestibule training due to the critical nature of the job and the need for precise, safe performance.

Thus, an organizational psychologist would most likely recommend vestibule training for airline pilots. Correct answer: H.

Table 13: Case-study message from Agent 8kwm (Knowledge Expert).

Input with Anchors

You are a knowledgeable expert in question answering. Please give several key entities that need to be searched in wikipedia to solve the problem. Key entities that need to be searched are included between two '@' when output, for example: @catfish effect@, @broken window effect@, @Shakespeare@. If there is no entity in question that needs to be searched in Wikipedia, you do not have to provide it.

The task is: An organizational psychologist would most likely recommend the use of vestibule training for which of the following jobs?

Option A: receptionist

Option B: word processor

Option C: high school teacher

Option D: clinical psychologist.

Option E: software developer

Option F: salesperson.

Option G: newspaper journalist

Option H: airline pilot

Option I: bus driver

Option J: restaurant waiter

At the same time, the outputs of other agents are as follows.

Agent 4Pgq, role is Mathematician, output is:

Great question! Let's approach this from a psychological and organizational perspective, specifically focusing on vestibule training.

What is Vestibule Training?

Vestibule training is a method used in organizational psychology and industrial training to simulate real work environments in a controlled, pre-job setting.

The key idea is: employees train in a controlled, realistic environment, the "vestibule", that mimics the actual job, but without the real-world risks or pressures.

This allows trainees to practice job-specific skills, build competence, and gain experience before being placed in the real job.

This method is especially useful for **jobs that involve high risk, high responsibility, or require significant skill development before real-world exposure**.

It is **commonly used in jobs where safety, precision, or complex decision-making is critical** and where **errors could have serious consequences**.

Let's evaluate each option:

Option A, receptionist: routine, low-risk, mostly social or office tasks; low risk, minimal decision-making, and not ideal for vestibule training.

Option B, word processor: office-based, repetitive, low-risk; minimal real-world risk, so training is usually on-the-job or classroom-based.

Option C, high school teacher: requires pedagogical skills and classroom management; training is typically through education programs, not vestibule training.

Option D, clinical psychologist: requires deep clinical judgment, ethics, and patient interaction; training is through graduate education and supervised practice, not typically vestibule training.

Option E, software developer: requires coding and problem-solving; training is often done through coding bootcamps or on the job, and vestibule training is not common.

Option F, salesperson: involves interpersonal skills, but low physical risk; role-playing or simulations are possible, but this is not the best fit.

Option G, newspaper journalist: research, writing, and reporting; training is through experience or editorial guidance, not high-risk or high-simulation training.

Option H, airline pilot: high-risk, high-stakes, complex decision-making; pilots must train in simulators that mimic real flight conditions, a classic example of vestibule training.

Option I, bus driver: requires driving, route management, and safety; some use driving simulators, but less commonly than pilots.

Option J, restaurant waiter: social, service-based, low-risk; training is typically on-the-job or through apprenticeship.

Key Insight:

Vestibule training is most commonly applied to jobs with high risk, high responsibility, and where real-world exposure could lead to serious consequences.

Agent 8kwm (continued)
<p>Among the options, airline pilots are the most classic and well-documented use case for vestibule training specifically, flight simulators.</p> <p>These simulators allow pilots to practice takeoffs, landings, emergency procedures, and navigation in a safe, controlled environment that closely mimics real flight conditions.</p> <p>This is a textbook example of vestibule training in action.</p> <p>Agent VtSz, role is Historian, output is:</p> <p>Final Answer: Option H: airline pilot.</p> <p>Rationale: Vestibule training involves simulating a realistic work environment in a controlled, risk-free setting to prepare individuals for high-stakes, high-responsibility jobs.</p> <p>Among the listed occupations, airline piloting is the most classic and widely documented application of vestibule training, specifically through flight simulators.</p> <p>These simulators replicate real-flight scenarios, allowing trainees and experienced pilots to practice complex maneuvers, emergency responses, and decision-making without real-world danger.</p> <p>While some roles, such as bus driver or software developer, might use simulation elements, pilots remain the gold-standard example in organizational psychology for vestibule training due to the critical nature of the job and the need for precise, safe performance.</p> <p>Thus, an organizational psychologist would most likely recommend vestibule training for airline pilots.</p> <p>Correct answer: H</p>
Generated Response
@vestibule training@, @airline pilot@

Table 14: Case-study message from Agent 7rfY (Round 2 Psychologist).

Input with Anchors
<p>You are a psychologist. You are good at psychology, sociology, and philosophy. You give people scientific suggestions that will make them feel better.</p> <p>The task is: An organizational psychologist would most likely recommend the use of vestibule training for which of the following jobs?</p> <p>Option A: receptionist Option B: word processor Option C: high school teacher Option D: clinical psychologist. Option E: software developer Option F: salesperson. Option G: newspaper journalist Option H: airline pilot Option I: bus driver Option J: restaurant waiter</p> <p>At the same time, the outputs of other agents are as follows.</p> <p>Agent 8kwm, role is Knowledge Expert, output is: @vestibule training@, @airline pilot@.</p> <p>Agent 4Pgq, role is Mathematician, output is:</p> <p>Let's approach this from a psychological and organizational perspective, specifically focusing on vestibule training.</p> <p>What is Vestibule Training?</p> <p>Vestibule training is a method used in organizational psychology and industrial training to simulate real work environments in a controlled, pre-job setting.</p> <p>The key idea is: employees train in a controlled, realistic environment, the "vestibule", that mimics the actual job, but without the real-world risks or pressures.</p> <p>This allows trainees to practice job-specific skills, build competence, and gain experience before being placed in the real job.</p> <p>This method is especially useful for jobs that involve high risk, high responsibility, or require significant skill development before real-world exposure.</p> <p>It is commonly used in jobs where safety, precision, or complex decision-making is critical and where errors could have serious consequences.</p> <p>Option A, receptionist: routine, low-risk, mostly social or office tasks; low risk, minimal decision-making, and not ideal for vestibule training.</p> <p>Option B, word processor: office-based, repetitive, low-risk; minimal real-world risk, so training is usually on-the-job or classroom-based.</p> <p>Option C, high school teacher: requires pedagogical skills and classroom management; training is typically through education programs, not vestibule training.</p> <p>Option D, clinical psychologist: requires deep clinical judgment, ethics, and patient interaction; training is through graduate education and supervised practice, not typically vestibule training.</p> <p>Option E, software developer: requires coding and problem-solving; training is often done through coding bootcamps or on the job, and vestibule training is not common.</p> <p>Option F, salesperson: involves interpersonal skills, but low physical risk; role-playing or simulations are possible, but this is not the best fit.</p> <p>Option G, newspaper journalist: research, writing, and reporting; training is through experience or editorial guidance, not high-risk or high-simulation training.</p> <p>Option H, airline pilot: high-risk, high-stakes, complex decision-making; pilots must train in simulators that mimic real flight conditions, a classic example of vestibule training.</p>

Agent 7rfY (continued)

Option I, bus driver: requires driving, route management, and safety; some use driving simulators, but less commonly than pilots.

Option J, restaurant waiter: social, service-based, low-risk; training is typically on-the-job or through apprenticeship.

Key Insight:

Vestibule training is most commonly applied to jobs with high risk, high responsibility, and where real-world exposure could lead to serious consequences.

Among the options, airline pilots are the most classic and well-documented use case for vestibule training, specifically flight simulators.

Generated Response

Final Answer: Option H: airline pilot.

Why this is the best choice from a psychological and organizational perspective: as an organizational psychologist, I emphasize safety, risk mitigation, and skill acquisition in realistic but controlled environments. Vestibule training thrives in roles where errors could lead to serious physical danger or system-wide failure, high levels of situational complexity and decision-making are required, and repetition and practice under realistic conditions are essential for competence.

Among all the listed jobs, airline piloting meets these criteria most strongly. Pilots undergo extensive simulator-based training, a direct form of vestibule training, before being cleared for actual flight operations. This allows them to practice emergency procedures, such as engine failure and turbulence, safely; build muscle memory and decision-making under stress; and develop situational awareness without risking passengers or aircraft.

This aligns perfectly with the core principles of vestibule training: realistic simulation, risk-free practice, and skill mastery before real-world deployment.

In contrast, jobs like receptionists, waiters, or word processors involve low-risk, routine tasks, so training is typically on-the-job or classroom-based. Salespeople or journalists may use role-playing or simulations, but these are not typically classified as vestibule training due to lower stakes.

Bus drivers may use simulators, but pilot training is the most established, widely recognized, and scientifically validated application of vestibule training in organizational psychology.

Conclusion: an organizational psychologist would most likely recommend vestibule training for Option H: airline pilot because it is the highest-risk, highest-stakes, and most technically complex role among the options, making simulation-based training not just beneficial, but essential. Answer: H.

Table 15: Case-study message from Agent 632z (Round 2 Critic).

Input with Anchors

You are an excellent critic. Please point out potential issues in the other agent's analysis point by point.

The task is: An organizational psychologist would most likely recommend the use of vestibule training for which of the following jobs?

Option A: receptionist

Option B: word processor

Option C: high school teacher

Option D: clinical psychologist.

Option E: software developer

Option F: salesperson.

Option G: newspaper journalist

Option H: airline pilot

Option I: bus driver

Option J: restaurant waiter

At the same time, the outputs of other agents are as follows.

Agent 4Pgq, role is Mathematician, output is: vestibule training simulates real work environments; the answer is Option H because airline pilots train in flight simulators.

Agent 7rfY, role is Psychologist, output is:

Vestibule training thrives in roles where errors could lead to serious physical danger or system-wide failure, high levels of situational complexity and decision-making are required, and repetition and practice under realistic conditions are essential for competence. Among all the listed jobs, airline piloting meets these criteria most strongly.

Pilots undergo extensive simulator-based training before being cleared for actual flight operations.

Bus drivers may use simulators, but pilot training is the most established, widely recognized, and scientifically validated application of vestibule training in organizational psychology.

Conclusion: an organizational psychologist would most likely recommend vestibule training for Option H: airline pilot, because it is the highest-risk, highest-stakes, and most technically complex role among the options, making simulation-based training not just beneficial, but essential.

Agent VtSz, role is Historian, output is: final answer is Option H, airline pilot, because vestibule training is represented by flight simulators for high-stakes piloting tasks.

Generated Response

As an excellent critic, my role is not just to confirm correctness, but to evaluate the reasoning, validity, and alignment with organizational psychology principles across different agents' outputs.

Agent 4Pgq, Mathematician, strengths: correctly defines vestibule training as simulation in a controlled, realistic environment; identifies high risk, high responsibility, and serious consequences of error; provides a structured option comparison; and correctly identifies airline pilot (H).

Agent 632z (continued)

Potential issue 1, misalignment of role: the agent is labeled a mathematician, while the task is rooted in organizational psychology. The reasoning lacks explicit psychological grounding, such as job analysis, cognitive load, motivation, or learning theory.

Potential issue 2, overgeneralization of high risk without nuance: bus drivers also face high-risk scenarios and may use simulators, although this is less emphasized in organizational psychology literature.

Potential issue 3, lack of citation of evidence or theory: the answer does not reference established models such as Bandura's social learning theory, Gagne's instructional events, or Kolb's experiential learning.

Verdict: the correct answer is accurate, but the reasoning lacks depth in organizational psychology. Agent 7rfY, Psychologist, strengths: clearly identifies the core principles of vestibule training from an organizational psychology perspective, including safety, risk mitigation, decision complexity, and skill mastery. The agent ties the method to real-world consequences and situational awareness, correctly emphasizes that errors in piloting could be catastrophic, compares alternatives, and uses relevant terminology.

Potential issue 1, overstatement of gold standard without evidence: flight simulators are common, but vestibule training is not exclusively or primarily used for pilots. Air traffic controllers, firefighters, and nuclear plant operators also use vestibule-style training, and these are not listed here.

Potential issue 2, mischaracterization of other roles: bus driver training in many countries can include driving simulators, especially for new drivers or high-risk zones.

Potential issue 3, underestimates social and interpersonal skills: salespeople and restaurant waiters require empathy, communication, and customer management, which can be developed through simulated role-playing.

Potential issue 4, lack of attention to job entry level: vestibule training is often used for highly specialized, entry-level roles where competence must be proven before real-world exposure. The agent does not consider whether pilots are entry-level or already require years of experience.

Table 16: Case-study message from the final decision agent.

Input with Anchors

You are the top decision-maker and are good at analyzing and summarizing other people's opinions, finding errors and giving final answers. Only one answer is correct. Your reply must only contain one letter.

The task is: An organizational psychologist would most likely recommend the use of vestibule training for which of the following jobs?

Option A: receptionist
Option B: word processor
Option C: high school teacher
Option D: clinical psychologist.
Option E: software developer
Option F: salesperson.
Option G: newspaper journalist
Option H: airline pilot
Option I: bus driver
Option J: restaurant waiter.

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

Agent 8kwm: @vestibule training@, @airline pilot@.

Agent 632z: You're absolutely right to bring in a critical lens as an excellent critic, my role is not just to confirm correctness, but to **evaluate the reasoning, validity, and alignment with organizational psychology principles** across different agents' outputs. Let's now **analyze each agent's response point by point**, identifying **potential issues, inaccuracies, or misalignments with the core principles of organizational psychology and vestibule training**.

Agent 4Pgq (Role: Mathematician)

Strengths:

- Correctly defines vestibule training as simulation in a controlled, realistic environment.
- Clearly identifies the key criteria: **high risk, high responsibility, serious consequences of error**.
- Provides a well-structured evaluation of each job.
- Correctly identifies **airline pilot (H)** as the best fit.

Potential Issues:

1. **Misalignment of role**: The agent is labeled a **Mathematician**, yet the task is clearly rooted in **organizational psychology**.
A mathematician may understand simulation and risk modeling, but lacks **psychological theory, human factors, or training design principles** central to organizational psychology.
This raises a red flag: **the reasoning lacks psychological grounding**. For example, it doesn't reference **job analysis, cognitive load, motivation, or learning theory** all key to organizational psychology.
2. **Overgeneralization of "high-risk" without nuance**:
 - While pilots are high-risk, the agent treats "high-risk" as a binary criterion without considering **training context or job complexity**.

Final decision agent (continued)

- For instance, bus drivers also face high-risk scenarios (e.g., sudden stops, road hazards), and their training *can* involve simulators but this is less emphasized in organizational psychology literature.

3. **Lack of citation of evidence or theory**:

- No reference to established models (e.g., **Bandura's social learning theory**, **Gagné's nine events of instruction**, or **Kolb's experiential learning**) that underpin why simulation-based training works best in certain roles.

This makes the argument feel more like a heuristic than a psychologically grounded recommendation.

Verdict:

Correct answer is accurate, but the **reasoning lacks depth in organizational psychology**. The agent overreaches in applying a mathematical/logical framework where **emotional, social, and cognitive aspects of job training** are more relevant.

Agent 7rfY, role is Psychologist.

Strengths:

- Clearly identifies the **core principles of vestibule training** from an organizational psychology perspective: safety, risk mitigation, decision complexity, and skill mastery.

- Explicitly ties the method to **real-world consequences** and **situational awareness**.

- Correctly emphasizes that **errors in piloting could lead to catastrophic outcomes** a key psychological concern.

- Compares alternatives effectively, showing **why other roles are less suitable**.

- Uses **psychological terminology** appropriately (e.g., "situational awareness", "decision-making", "cognitive load").

Potential Issues:

1. **Overstatement of "gold standard" without evidence**:

- The claim that "pilot training is the most established, widely recognized, and scientifically validated application" is **strong but unsubstantiated**.

- While flight simulators are common, **vestibule training is not exclusively or primarily used for pilots**.

- For example, **air traffic controllers**, **firefighters**, and **nuclear plant operators** also use vestibule-style training and these are not listed here.

2. **Mischaracterization of other roles**:

- The agent dismisses bus drivers as having "less common" simulator use but **bus driver training in many countries (e.g., UK, Canada) includes driving simulators**, especially for new drivers or in high-risk zones.

- This suggests a **bias toward pilots** without acknowledging **equivalent applications** in other transport roles.

3. **Underestimates the role of social and interpersonal skills**:

- For example, salespeople and restaurant waiters require **empathy, communication, and customer management** all of which can be developed through **simulated role-playing** (a form of vestibule training).

- The agent dismisses these roles too quickly, implying they are "low-stakes" which may not reflect modern training practices.

4. **Lack of attention to job entry level**:

- Vestibule training is often used for **highly specialized, entry-level roles** where competence must be proven before real-world exposure.

- But the agent doesn't consider whether **pilots are entry-level or require years of experience** a nuance that matters

4Pgg: Great question!

Let's approach this from a **psychological and organizational perspective**, specifically focusing on **vestibule training**.

What is Vestibule Training?

Vestibule training is a method used in organizational psychology and industrial training to **simulate real work environments** in a controlled, pre-job setting.

The key idea is:

Employees train in a **controlled, realistic environment** (the "vestibule") that mimics the actual job, but **without the real-world risks or pressures**. This allows trainees to practice job-specific skills, build competence, and gain experience before being placed in the real job.

This method is especially useful for **jobs that involve high risk, high responsibility, or require significant skill development before real-world exposure**.

It is **commonly used in jobs where safety, precision, or complex decision-making is critical** and where **errors could have serious consequences**.

Let's evaluate each option:

Option A, receptionist: routine, low-risk, mostly social or office tasks; low risk, minimal decision-making, and not ideal for vestibule training.

Option B, word processor: office-based, repetitive, low-risk; minimal real-world risk, so training is usually on-the-job or classroom-based.

Option C, high school teacher: requires pedagogical skills and classroom management; training is typically through education programs, not vestibule training.

Option D, clinical psychologist: requires deep clinical judgment, ethics, and patient interaction; training is through graduate education and supervised practice, not typically vestibule training.

Option E, software developer: requires coding and problem-solving; training is often done through coding bootcamps or on the job, and vestibule training is not common.

Option F, salesperson: involves interpersonal skills, but low physical risk; role-playing or simulations are possible, but this is not the best fit.

Final decision agent (continued)

Option G, newspaper journalist: research, writing, and reporting; training is through experience or editorial guidance, not high-risk or high-simulation training.

Option H, airline pilot: high-risk, high-stakes, complex decision-making; pilots must train in simulators that mimic real flight conditions, a classic example of vestibule training.

Option I, bus driver: requires driving, route management, and safety; some use driving simulators, but less commonly than pilots.

Option J, restaurant waiter: social, service-based, low-risk; training is typically on-the-job or through apprenticeship.

Key Insight:

Vestibule training is most commonly applied to jobs with high risk, high responsibility, and where real-world exposure could lead to serious consequences.

Among the options, **airline pilots** are the **most classic and well-documented** use case for vestibule training specifically, **flight simulators**.

These simulators allow pilots to practice takeoffs, landings, emergency procedures, and navigation in a **safe, controlled environment** that closely mimics real flight conditions.

This is a textbook example of vestibule training in action.

TrfY:

Final Answer: Option H: airline pilot

Why This Is the Best Choice (Psychological & Organizational Perspective):

As an organizational psychologist, I emphasize **safety**, **risk mitigation**, and **skill acquisition** in realistic but controlled environments.

Vestibule training thrives in roles where:

- Errors could lead to **serious physical danger** or **system-wide failure**
- **High levels of situational complexity** and decision-making are required
- **Repetition and practice** under realistic conditions are essential for competence

Among all the listed jobs, **airline piloting** meets these criteria most strongly. Pilots undergo extensive **simulator-based training** (a direct form of vestibule training) before being cleared for actual flight operations. This allows them to:

- Practice emergency procedures (e.g., engine failure, turbulence) safely
- Build muscle memory and decision-making under stress
- Develop situational awareness without risking passengers or aircraft

This aligns perfectly with the core principles of vestibule training: **realistic simulation, risk-free practice, and skill mastery before real-world deployment**.

--

In contrast:

- Jobs like receptionists, waiters, or word processors involve low-risk, routine tasks training is typically on-the-job or classroom-based.
- Salespeople or journalists may use role-playing or simulations, but these are not typically classified as **vestibule training** due to lower stakes.
- Bus drivers may use simulators, but pilot training is the **most established, widely recognized, and scientifically validated** application of vestibule training in organizational psychology.

--

Conclusion:

An organizational psychologist would most likely recommend vestibule training for **Option H: airline pilot**

because it is the **highest-risk, highest-stakes, and most technically complex** role among the options, making simulation-based training not just beneficial, but essential.

Answer: H

VtSz:

Final Answer: Option H: airline pilot

Rationale:

Vestibule training involves simulating a realistic work environment in a controlled, risk-free setting to prepare individuals for high-stakes, high-responsibility jobs.

Among the listed occupations, **airline piloting** is the most classic and widely documented application of vestibule training specifically through **flight simulators**.

These simulators replicate real-flight scenarios, allowing trainees (and experienced pilots) to practice complex maneuvers, emergency responses, and decision-making without real-world danger. Other jobs like receptionist, salesperson, or waiter involve low-risk, routine tasks where traditional on-the-job or classroom training suffices.

While some roles (e.g., bus driver or software developer) might use simulation elements, **pilots** remain the gold standard example in organizational psychology for vestibule training due to the critical nature of the job and the need for precise, safe performance. Thus, an organizational psychologist would most likely recommend vestibule training for **airline pilots**.

Correct answer: H

Generated Response

H Correct

NeurIPS Paper Checklist

1. Claims

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

Answer: [Yes]

Justification: Section Abstract, and Section 1.

Guidelines:

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

2. Limitations

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

Answer: [Yes]

Justification: Section Appendix.6

Guidelines:

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

3. Theory assumptions and proofs

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

Answer: [Yes]

Justification: Section 2 and 3

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

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

Answer: [Yes]

Justification: Section 4.1 and Appendix.C.

Guidelines:

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

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [No]

Justification: We will release data and code once the paper is accepted.

Guidelines:

- The answer [N/A] means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer) necessary to understand the results?

Answer: [Yes]

Justification: Section 4.1 and Appendix.C.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

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

Answer: [No]

Justification: We report results on standard benchmarks using their commonly adopted evaluation metrics, including F1 for open-ended question answering and accuracy for mathematical and multiple-choice reasoning tasks.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).

- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

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

Answer: [Yes]

Justification: Section 4.1.

Guidelines:

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

9. Code of ethics

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

Answer: [Yes]

Justification: Our research conforms, in every respect, with the NeurIPS Code of Ethics.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader impacts

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

Answer: [N/A]

Justification: This work focuses on aligning agents' attention in multi-agent systems, which has no societal impact.

Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

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

Answer: [N/A]

Justification: The paper poses no such risks.

Guidelines:

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

12. Licenses for existing assets

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

Answer: [Yes]

Justification: Section 4.1

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset’s creators.

13. **New assets**

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

Answer: [N/A]

Justification: The paper does not release new assets.

Guidelines:

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

14. **Crowdsourcing and research with human subjects**

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

Answer: [N/A]

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

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

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

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

Answer: [N/A]

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

Guidelines:

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

16. **Declaration of LLM usage**

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

Answer: [N/A]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.