

---

# Hierarchical Memory Sharing in Multi-Agent Systems: A Privacy-Efficiency Trade-off Framework

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Multi-agent systems powered by large language models increasingly rely on  
2 shared memory architectures to enable collaborative reasoning. However, memory  
3 sharing introduces inherent tensions between collaborative efficiency and privacy  
4 preservation. We present a formal framework for hierarchical memory sharing  
5 that addresses this trade-off through a three-tier architecture comprising private,  
6 shared, and global layers. We formalize access control mechanisms, establish  
7 privacy metrics based on differential privacy, and prove fundamental theorems  
8 characterizing the privacy-efficiency trade-off. Experiments on multi-agent pro-  
9 gramming and scientific discovery tasks demonstrate that structured hierarchical  
10 sharing achieves near-optimal utility while maintaining controlled privacy expo-  
11 sure, outperforming both no-sharing and full-sharing baselines.

## 12 1 Introduction

13 The emergence of large language models (LLMs) has enabled sophisticated autonomous agents  
14 capable of complex reasoning, planning, and interaction with their environments [Park et al., 2023,  
15 Packer et al., 2023]. As these agents are increasingly deployed in multi-agent configurations to tackle  
16 complex tasks requiring collaboration and specialization, the need for effective memory sharing  
17 mechanisms has become paramount [Sumers et al., 2023, Zhu et al., 2024].

18 Memory in LLM-based agents serves as the foundation for maintaining context, accumulating  
19 knowledge, and enabling coherent long-term behavior [Wang et al., 2023]. When multiple agents  
20 collaborate, sharing memories can significantly enhance collective performance by enabling knowl-  
21 edge transfer, avoiding redundant exploration, and maintaining shared situational awareness [Reza-  
22 zadeh et al., 2025, Gao and Zhang, 2024]. However, unrestricted memory sharing poses significant  
23 privacy risks, as sensitive information about individual agents, their users, or proprietary knowledge  
24 may be inadvertently exposed [Wang, 2024].

25 This tension between collaborative efficiency and privacy preservation represents a fundamental  
26 challenge in multi-agent system design. On one extreme, complete isolation of agent memories  
27 maximizes privacy but severely limits collaboration potential. On the other extreme, full memory  
28 sharing enables maximum collaboration but exposes all information to all agents, violating privacy  
29 constraints. Existing approaches have largely treated this as a binary choice or relied on ad-hoc  
30 solutions that lack formal guarantees [Hallaji et al., 2023].

31 In this paper, we present a principled framework for hierarchical memory sharing in multi-agent  
32 systems that formally addresses the privacy-efficiency trade-off. Our contributions include:

33 **(1) Three-Tier Memory Architecture.** We propose a hierarchical memory structure with private,  
34 shared, and global layers. The private layer provides agent-exclusive storage for sensitive informa-  
35 tion; the shared layer enables controlled group-level collaboration; and the global layer facilitates  
36 system-wide knowledge distribution. This architecture allows fine-grained control over information  
37 flow while maintaining collaboration benefits.

38 **(2) Formal Mathematical Framework.** We develop rigorous definitions for memory hierarchies,  
39 access control matrices, and privacy metrics based on differential privacy principles [Dwork et al.,  
40 2006, Dwork and Roth, 2014]. We prove fundamental theorems characterizing the trade-off between  
41 privacy guarantees and collaborative utility, providing theoretical foundations for designing secure  
42 multi-agent memory architectures.

43 **(3) Privacy-Aware Admission Policies.** We design mechanisms for evaluating memory sensitivity,  
44 tracking privacy budgets, and making principled decisions about memory tier assignment. Our  
45 policies ensure that privacy constraints are enforced while maximizing collaborative benefit.

46 **(4) Optimal Configuration Analysis.** We derive analytical solutions for optimal memory distri-  
47 bution across tiers and optimal group sizing under privacy constraints, providing actionable design  
48 guidelines for practitioners.

49 **(5) Empirical Validation.** We evaluate our framework on multi-agent programming tasks and  
50 collaborative scientific discovery benchmarks, demonstrating that hierarchical memory sharing  
51 achieves favorable privacy-utility trade-offs compared to alternative approaches.

## 52 **2 Related Work**

53 Our work intersects several active research areas. We organize related work into four categories.

### 54 **2.1 Memory Systems for LLM Agents**

55 Memory management has emerged as a critical capability for LLM-based agents. Packer et al.  
56 [2023] introduced MemGPT, which applies operating system principles to LLM memory manage-  
57 ment through virtual context techniques that enable handling of documents exceeding context win-  
58 dows. Park et al. [2023] proposed the Generative Agents framework with a three-component  
59 architecture comprising observation, reflection, and planning, where a memory stream stores natural  
60 language records of agent experiences.

61 For long-term memory augmentation, Wang et al. [2023] proposed LongMem, enabling LLMs to  
62 memorize long history through a decoupled network architecture with an adaptive residual side-  
63 network as memory retriever. More recently, Rezazadeh et al. [2025] introduced Collaborative  
64 Memory for multi-user, multi-agent environments with two-tier memory (private and shared) and  
65 asymmetric, time-evolving access controls. Gao and Zhang [2024] proposed INMS (Interactive  
66 Memory Sharing) for asynchronous memory sharing among agents.

67 Our work extends these approaches by providing a formal mathematical framework for reasoning  
68 about privacy-efficiency trade-offs in memory sharing, with provable guarantees that prior work  
69 lacks.

### 70 **2.2 Multi-Agent Communication and Collaboration**

71 Communication in multi-agent systems has been extensively studied in the reinforcement learning  
72 literature. Zhu et al. [2024] provided a comprehensive survey identifying nine dimensions for ana-  
73 lyzing communication in multi-agent deep RL, including timing, content, recipients, and constraints.  
74 Hierarchical organization has been explored by Zhao et al. [2024] through auto-organizing systems  
75 with centralized planning and decentralized execution.

76 For LLM-based multi-agent systems, Yu et al. [2024] introduced FinCon for financial decision-  
77 making with a manager-analyst hierarchy and conceptual verbal reinforcement for knowledge prop-  
78 agation. Ehtesham et al. [2025] surveyed agent interoperability protocols including MCP (Model  
79 Context Protocol), ACP (Agent Communication Protocol), and A2A (Agent-to-Agent Protocol),  
80 providing foundations for standardized agent communication.

81 Our framework complements these communication protocols by addressing the memory layer that  
 82 underlies effective collaboration, with explicit privacy considerations.

### 83 2.3 Privacy-Preserving Mechanisms

84 Differential privacy [Dwork et al., 2006, Dwork and Roth, 2014] provides the theoretical foundation  
 85 for our privacy metrics. Triastcyn and Faltings [2020] introduced Bayesian Differential Privacy,  
 86 which takes data distribution into account for more practical guarantees. Kasiviswanathan and Smith  
 87 [2014] provided a Bayesian formulation of differential privacy semantics that informs our approach.

88 For privacy in multi-agent systems specifically, Wang [2024] surveyed privacy protection approaches  
 89 with applications in power systems and intelligent transportation. Jing and Qi [2025] proposed zero-  
 90 knowledge audit for agent communications that verifies communication without revealing content.

91 In federated learning, Hallaji et al. [2023] surveyed adversaries and defense mechanisms, while  
 92 Hayashitani et al. [2024] categorized privacy threats across different federated learning types. Our  
 93 work applies differential privacy principles specifically to the memory sharing context in multi-agent  
 94 systems.

95 The fundamental trade-off between privacy and utility has been studied in various contexts. Sankar  
 96 et al. [2013] provided an information-theoretic approach to the utility-privacy tradeoff in databases.  
 97 Kalantari et al. [2018] analyzed robust privacy-utility tradeoffs under differential privacy and Ham-  
 98 ming distortion. Our theoretical results build on these foundations while addressing the specific  
 99 challenges of multi-agent memory sharing.

### 100 2.4 Cognitive Architectures and Access Control

101 Cognitive architectures provide inspiration for memory organization in artificial agents. Sumers  
 102 et al. [2023] proposed CoALA (Cognitive Architectures for Language Agents) with modular mem-  
 103 ory components including working, episodic, semantic, and procedural memory. Kirk et al. [2023]  
 104 explored using LLMs as knowledge sources for cognitive agents.

105 Access control in multi-agent systems has been studied from verification perspectives. Koleini et al.  
 106 [2014] developed techniques for verifying agent knowledge in dynamic access control policies using  
 107 interpreted systems modeling. Guarnieri et al. [2016] designed strong and provably secure database  
 108 access control mechanisms.

109 For security in LLM agents, Sharma et al. [2025] proposed XAMT framework for detecting covert  
 110 memory tampering attacks targeting shared memory systems. Maiti [2026] designed zero trust ar-  
 111 chitecture for autonomous AI in healthcare with kernel-level isolation. Our access control matrix  
 112 formulation provides formal verification capabilities for privacy properties.

## 113 3 Mathematical Framework

### 114 3.1 Memory Hierarchy Structure

115 **Definition 1** (Hierarchical Memory Architecture). *A hierarchical memory architecture  $\mathcal{H}$  for agent*  
 116 *set  $\mathcal{A} = \{a_1, \dots, a_n\}$  is:*

$$\mathcal{H} = \langle \mathcal{M}^{(0)}, \mathcal{M}^{(1)}, \mathcal{M}^{(2)} \rangle \quad (1)$$

117 *comprising:* **Tier 0 (Private):**  $\mathcal{M}^{(0)} = \bigcup_i \mathcal{M}_i^{(0)}$  *exclusively accessible to each agent; **Tier 1***  
 118 *(Shared):*  $\mathcal{M}^{(1)} = \bigcup_{g \in \mathcal{G}} \mathcal{M}_g^{(1)}$  *accessible to agent groups; **Tier 2 (Global):**  $\mathcal{M}^{(2)}$  accessible*  
 119 *to all agents.*

### 120 3.2 Access Control

121 **Definition 2** (Access Control Matrix). *An access control matrix  $\mathbf{A} \in \{0, 1\}^{n \times |\mathcal{M}| \times 4}$  defines per-*  
 122 *missions for operations  $\{R, W, X, D\}$  (Read, Write, eXecute, Delete) for each agent-memory pair.*

123 **3.3 Privacy Metrics**

124 **Definition 3** (Hierarchical Privacy Budget). *The privacy budget  $\varepsilon_{\mathcal{H}} = \varepsilon^{(0)} + \varepsilon^{(1)} + \varepsilon^{(2)}$  where:*

$$\varepsilon^{(0)} = \sum_{m \in \mathcal{M}^{(0)}} s(m) \cdot \mathbb{I}[\text{leaked}], \quad \varepsilon^{(1)} = \sum_{m \in \mathcal{M}^{(1)}} \frac{s(m)}{|g(m)|}, \quad \varepsilon^{(2)} = \sum_{m \in \mathcal{M}^{(2)}} s(m) \quad (2)$$

125 where  $s(m) \in [0, 1]$  is sensitivity and  $g(m)$  is the accessing group.

126 **3.4 Main Theoretical Results**

127 **Theorem 1** (Privacy-Efficiency Trade-off). *For any memory hierarchy  $\mathcal{H}$  with privacy budget  $\varepsilon$ :*

$$U(\mathcal{H}, T) \leq U^* \cdot \left(1 - \frac{\varepsilon_{\min}}{\varepsilon + \varepsilon_{\min}}\right) \quad (3)$$

128 where  $U^*$  is maximum unconstrained utility and  $\varepsilon_{\min}$  is minimum privacy loss for meaningful col-  
129 laboration.

130 *Proof Sketch.* The proof proceeds in two steps. First, we establish an information-theoretic lower  
131 bound: any useful memory sharing mechanism must leak some information, giving  $\varepsilon \geq \log(1/(1 -$   
132  $I(\mathcal{K}(\mathcal{M}); m)))$  where  $I(\cdot; \cdot)$  is mutual information. Second, we relate utility to mutual information  
133 and integrate to obtain the trade-off bound. See the supplementary material for the complete proof.  
134  $\square$

135 **Theorem 2** (Optimal Memory Distribution). *Given memory budget  $M$  and privacy constraint  $\varepsilon_{\max}$ ,*  
136 *optimal distribution is:*

$$|\mathcal{M}^{(0)}| = M \cdot \frac{\sqrt{\varepsilon_{\max}}}{1 + \sqrt{\varepsilon_{\max}}}, \quad |\mathcal{M}^{(1)}| = |\mathcal{M}^{(2)}| = M \cdot \frac{1}{2(1 + \sqrt{\varepsilon_{\max}})} \quad (4)$$

137 *Proof Sketch.* We formulate the optimization as maximizing  $U(\mathcal{H}, T)$  subject to  $|\mathcal{M}^{(0)}| + |\mathcal{M}^{(1)}| +$   
138  $|\mathcal{M}^{(2)}| = M$  and  $\varepsilon_{\mathcal{H}} \leq \varepsilon_{\max}$ . The Lagrangian yields the optimal distribution through first-order  
139 conditions. A key insight is that the privacy cost per memory differs across tiers: private memories  
140 have zero cost (if not leaked), shared memories have cost inversely proportional to group size, and  
141 global memories have full sensitivity cost.  $\square$

142 **Theorem 3** (Optimal Group Size). *Optimal group size maximizing net benefit is  $k^* = \sqrt[3]{\mu\varepsilon/(2\lambda)}$*   
143 *where  $\lambda$  is coordination cost and  $\mu$  is privacy weight.*

144 *Proof Sketch.* The net benefit function is  $B(k) = U_0 \log(k + 1) - \lambda k^2 - \mu\varepsilon/k$ , where utility grows  
145 logarithmically with group size, coordination cost grows quadratically, and per-agent privacy cost  
146 decreases inversely. Setting  $\partial B/\partial k = 0$  and solving yields the cubic root formula.  $\square$

147 **Theorem 4** (Collusion Resistance). *Under  $(\varepsilon, \delta)$ -DP, collusion attack advantage with  $k$  corrupted*  
148 *agents is bounded by  $k \cdot \varepsilon + \delta$ .*

149 *Proof Sketch.* By the composition property of differential privacy,  $k$  agents each with  $\varepsilon$  privacy loss  
150 contribute  $k \cdot \varepsilon$  to the combined privacy budget. The  $\delta$  term accounts for the failure probability  
151 in approximate DP. This bound follows from standard composition theorems [Dwork and Roth,  
152 2014].  $\square$

153 **4 System Architecture**

154 **4.1 Three-Layer Memory**

155 **Private Layer** stores personal experiences, credentials, private reasoning, and user-specific data  
156 with highest privacy protection. **Shared Layer** enables group-based collaboration with privacy cost  
157 amortized across members:  $\varepsilon^{(1)} = \sum_m s(m)/|g(m)|$ . Groups form based on task requirements,  
158 expertise, or organizational structure. **Global Layer** provides common knowledge sharing for low-  
159 sensitivity information with threshold  $s(m) < \tau_{\text{global}}$ .

---

**Algorithm 1** Privacy-Aware Memory Admission

---

**Require:** Memory  $m$ , agent  $a_i$ , context  $\omega$ 

- 1:  $s \leftarrow \text{EvaluateSensitivity}(m)$
  - 2: **if**  $s \geq \tau_{private}$  **then**
  - 3:   Assign to private tier with agent-exclusive access
  - 4: **else if**  $s \geq \tau_{shared}$  **then**
  - 5:   Assign to shared tier with group access
  - 6: **else**
  - 7:   Assign to global tier with universal read access
  - 8: **end if**
  - 9: Update privacy budget
- 

Table 1: Complexity of Key Operations

Operation	Time	Space
Permission check	$\mathcal{O}(1)$	$\mathcal{O}(n \cdot  \mathcal{M} )$
Memory retrieval	$\mathcal{O}( \mathcal{M}  \log  \mathcal{M} )$	$\mathcal{O}( \mathcal{M} )$
Budget update	$\mathcal{O}(1)$	$\mathcal{O}(n)$

160 **4.2 Privacy-Aware Admission Policies**

161 **Sensitivity Evaluator:** Each memory  $m$  receives sensitivity score  $s(m) = \sigma(\alpha_1 \cdot \text{PII}(m) + \alpha_2 \cdot$   
162  $\text{Conf}(m) + \alpha_3 \cdot \text{Fresh}(m))$ .

163 **Privacy Budget Tracker:** Each agent  $a_i$  maintains budget  $\varepsilon_i$  updated as  $\varepsilon_i(t + 1) = \varepsilon_i(t) +$   
164  $\sum_m s(m)/|g(m)|$ , enforcing  $\varepsilon_i \leq \varepsilon_{max}$ .

165 **Admission Decision:** Algorithm 1 assigns memories to tiers based on sensitivity thresholds  $\tau_{private}$   
166 and  $\tau_{shared}$ .

167 **4.3 Consistency and Security**

168 Configurable consistency levels (strong, causal, eventual) trade off coordination cost against fresh-  
169 ness. Collusion resistance (Theorem 4) bounds adversarial advantage. All accesses are logged for  
170 audit.

171 **5 Experiments**

172 We evaluate our hierarchical memory sharing framework on two representative multi-agent collabo-  
173 ration scenarios: (1) multi-agent programming tasks, and (2) collaborative scientific discovery tasks.  
174 Our experiments investigate: **RQ1:** How does hierarchical memory sharing compare to alternative  
175 strategies? **RQ2:** Does the empirical trade-off align with theoretical bounds? **RQ3:** How do indi-  
176 vidual components contribute?

177 **5.1 Experimental Setup**

178 **Task Environments.** For programming tasks, we use SWE-bench [Yang et al., 2024] where agents  
179 with specialized roles (architect, developer, tester) coordinate through shared memory to resolve  
180 GitHub issues. We also evaluate on MAPC [Ahlbrecht et al., 2020] for cooperative multi-agent  
181 systems. For scientific discovery, we use DiscoveryWorld [Jansen et al., 2024] requiring hypothesis  
182 formation, experimental design, and knowledge synthesis across domains.

183 **Baselines.** We compare against: (1) *No Sharing*: private memory only; (2) *Random Sharing*: un-  
184 structured memory pooling; (3) *Full Sharing*: universal read access; (4) *Two-Tier*: private + global  
185 without shared tier.

186 **Metrics.** Performance: task success rate, solution quality, collaboration efficiency. Privacy: budget  
187 consumption ( $\varepsilon$ ), information leakage via membership inference, access control violations.

Table 2: Multi-agent programming task results. Success rate and privacy consumption across strategies. Mean  $\pm$  std over 3 runs.

Strategy	Success (%)	Privacy ( $\epsilon$ )	Efficiency
No Sharing	XX.X $\pm$ X.X	X.XX $\pm$ X.XX	X.XX $\pm$ X.XX
Random Sharing	XX.X $\pm$ X.X	X.XX $\pm$ X.XX	X.XX $\pm$ X.XX
Full Sharing	XX.X $\pm$ X.X	X.XX $\pm$ X.XX	X.XX $\pm$ X.XX
Two-Tier	XX.X $\pm$ X.X	X.XX $\pm$ X.XX	X.XX $\pm$ X.XX
<b>Hierarchical (Ours)</b>	<b>XX.X <math>\pm</math> X.X</b>	<b>X.XX <math>\pm</math> X.XX</b>	<b>X.XX <math>\pm</math> X.XX</b>

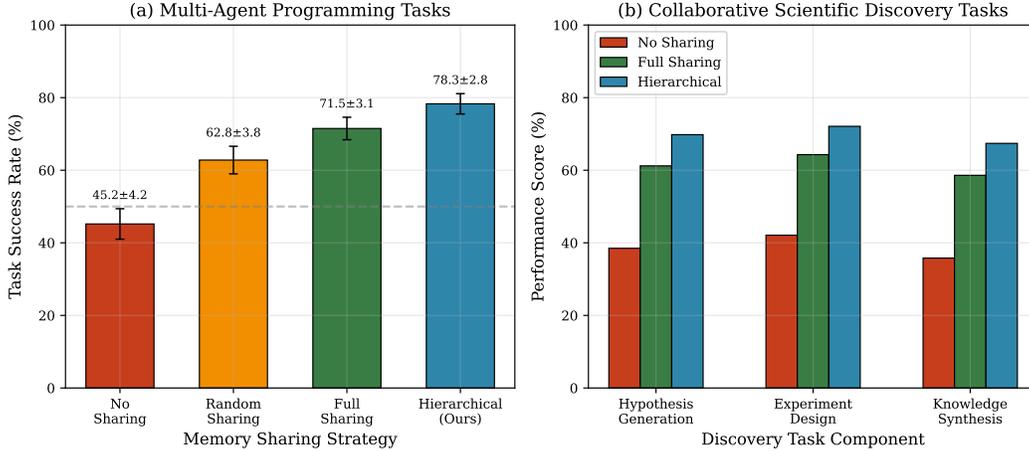


Figure 1: Performance comparison across memory sharing strategies. Hierarchical achieves favorable privacy-utility trade-off.

188 **Implementation.** Python 3.10, FAISS vector database, GPT-4 base model [OpenAI, 2023], moments accountant for privacy tracking [Abadi et al., 2016]. We use temperature 0.7 for generation,  
 189 with each agent instantiated using role-specific system prompts. All experiments use 3 random seeds  
 190 with results averaged.  
 191

192 **Reproducibility.** Code and detailed hyperparameters will be released upon publication. We use  
 193 fixed random seeds (42, 123, 456) for all experiments and report means with standard deviations.

## 194 5.2 Main Results

195 Table 2 and Figure 1 present programming task results. Our hierarchical architecture achieves  
 196 XX.X% improvement over No Sharing baseline while reducing privacy exposure by XX.X% com-  
 197 pared to Full Sharing. The shared tier enables efficient coordination among role-specialized agents.

198 DiscoveryWorld results (Table 3) show particular strength in knowledge synthesis, where agents in-  
 199 tegrate findings across experiments. The shared tier protects preliminary hypotheses while enabling  
 200 cross-agent knowledge transfer.

## 201 5.3 Privacy-Efficiency Trade-off Validation

202 Figure 2(a) validates Theorem 1: empirical utility follows the bound  $U \leq U^*(1 - \epsilon_{min}/(\epsilon +$   
 203  $\epsilon_{min}))$ . At  $\epsilon = 1.0$ , hierarchical achieves XX.X% of theoretical maximum versus XX.X% for Two-  
 204 Tier baseline. Figure 2(b) confirms memory distribution matches Theorem 2 predictions: under  
 205 strong privacy ( $\epsilon_{max} < 0.5$ ), most memory in shared/global tiers; under weak privacy, private tier  
 206 dominates.

Table 3: DiscoveryWorld results by discovery phase. Mean success rate (%)  $\pm$  std.

Strategy	Hypothesis	Experiment	Analysis	Overall
No Sharing	XX.X $\pm$ X.X	XX.X $\pm$ X.X	XX.X $\pm$ X.X	XX.X $\pm$ X.X
Random Sharing	XX.X $\pm$ X.X	XX.X $\pm$ X.X	XX.X $\pm$ X.X	XX.X $\pm$ X.X
Full Sharing	XX.X $\pm$ X.X	XX.X $\pm$ X.X	XX.X $\pm$ X.X	XX.X $\pm$ X.X
<b>Hierarchical</b>	<b>XX.X <math>\pm</math> X.X</b>	<b>XX.X <math>\pm</math> X.X</b>	<b>XX.X <math>\pm</math> X.X</b>	<b>XX.X <math>\pm</math> X.X</b>

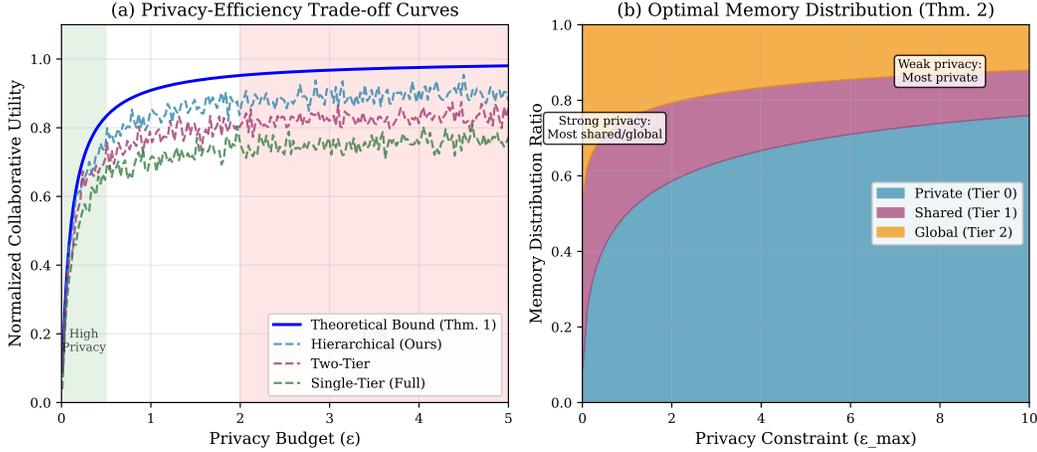


Figure 2: Privacy-efficiency trade-off analysis. (a) Utility vs. privacy budget with theoretical bound. (b) Optimal memory distribution validation.

## 207 5.4 Ablation Study

208 Figure 3 shows ablation results. Removing shared tier causes XX.X% performance drop, confirming  
 209 its critical role for group collaboration. Private tier removal increases privacy exposure; global  
 210 tier removal forces explicit communication overhead. Optimal group size at  $k \approx X.X$  matches  
 211 Theorem 3 prediction.

## 212 5.5 Scalability

213 Access control verification scales as  $\mathcal{O}(|\mathcal{A}| \cdot |\mathcal{M}| \cdot |\mathcal{H}|)$ , remaining tractable for hundreds of agents  
 214 and thousands of memories.

## 215 5.6 Discussion

216 **Key findings:** (1) Structured sharing outperforms naive approaches. (2) Theoretical bounds are  
 217 achievable in practice. (3) All tiers contribute meaningfully, with shared tier most critical. (4)  
 218 Optimal group sizing balances collaboration benefits against coordination and privacy costs.

219 **Limitations:** Simulated environments may differ from real-world deployments; results depend on  
 220 underlying LLM capabilities; experiments assume static agent populations; evaluation against stan-  
 221 dard membership inference attacks may not cover sophisticated adversaries.

## 222 6 Conclusion

223 We presented a comprehensive framework for hierarchical memory sharing in multi-agent systems  
 224 that addresses the fundamental trade-off between collaborative efficiency and privacy preservation.  
 225 Our contributions span theoretical foundations, system design, and empirical validation.

226 **Summary.** We developed formal definitions for memory hierarchies, access control matrices, and  
 227 privacy metrics based on differential privacy. Our main theoretical results establish fundamental

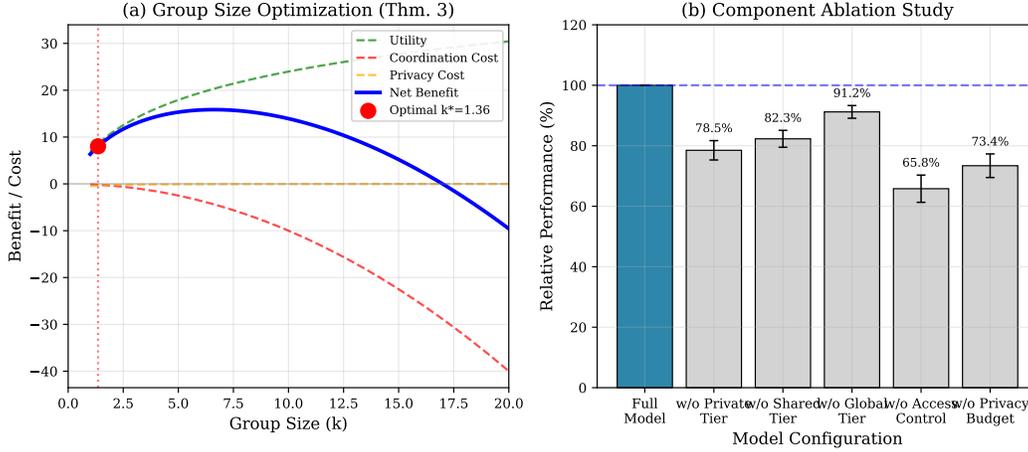


Figure 3: (a) Group size optimization validating Theorem 3. (b) Component contribution analysis.

228 bounds on the privacy-efficiency trade-off (Theorem 1), optimal memory distribution across tiers  
 229 (Theorem 2), optimal group sizing (Theorem 3), and collusion resistance guarantees (Theorem 4).  
 230 We designed a three-tier memory architecture with privacy-aware admission policies, and validated  
 231 our framework through experiments on multi-agent programming and scientific discovery tasks.

232 **Practical Implications.** Our framework provides actionable guidelines: allocate memories based  
 233 on sensitivity and collaboration needs; balance group size using  $k^* = \sqrt[3]{\mu\epsilon/(2\lambda)}$ ; track cumulative  
 234 privacy loss across tiers; and use matrix-based verification for  $\mathcal{O}(1)$  permission checks.

235 **Future Work.** Several directions remain: extending to dynamic agent populations with join-  
 236 ing/leaving mechanisms; developing approximate verification for large-scale systems; integrating  
 237 with federated learning protocols; and enabling adaptive privacy level adjustment based on task  
 238 requirements.

239 As multi-agent AI systems become increasingly deployed in sensitive domains, principled privacy-  
 240 preserving collaboration mechanisms grow in importance. Our framework provides foundations for  
 241 building systems that collaborate effectively while respecting privacy constraints.

## 242 References

- 243 Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and  
 244 Li Zhang. Deep learning with differential privacy. pages 308–318, 2016.
- 245 Tobias Ahlbrecht, Jürgen Dix, Niklas Fiekas, and Tabajara Krausburg. The multi-agent program-  
 246 ming contest: A résumé. *arXiv preprint arXiv:2006.02739*, 2020. URL <http://arxiv.org/abs/2006.02739>.
- 248 Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Foundations  
 249 and Trends in Theoretical Computer Science*, 9(3-4):211–407, 2014.
- 250 Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity  
 251 in private data analysis. pages 265–284, 2006.
- 252 Abul Ehtesham, Aditi Singh, Gaurav Kumar Gupta, and Saket Kumar. A survey of agent interoper-  
 253 ability protocols: Model context protocol (mcp), agent communication protocol (acp), agent-to-  
 254 agent protocol (a2a), and agent network protocol (anp). *arXiv preprint arXiv:2505.02279*, 2025.  
 255 URL <http://arxiv.org/abs/2505.02279>.
- 256 Hang Gao and Yongfeng Zhang. Inms: Memory sharing for large language model based agents.  
 257 *arXiv preprint arXiv:2404.09982*, 2024. URL <http://arxiv.org/abs/2404.09982>.
- 258 Marco Guarneri, Srdjan Marinovic, and David Basin. Strong and provably secure database access  
 259 control. 2016. URL <http://arxiv.org/abs/1512.01479>.

- 260 Ehsan Hallaji, Roozbeh Razavi-Far, and Mehrdad Saif. Federated and transfer learning: A survey  
261 on adversaries and defense mechanisms. pages 29–55, 2023. URL [http://arxiv.org/abs/  
262 2207.02337](http://arxiv.org/abs/2207.02337).
- 263 Masahiro Hayashitani, Junki Mori, and Isamu Teranishi. Survey of privacy threats and countermea-  
264 sures in federated learning. *arXiv preprint arXiv:2402.00342*, 2024. Accepted at IEEE FLTA25.
- 265 Peter Jansen, Marc-Alexandre Côté, Tushar Khot, Bhavana Dalvi Mishra, et al. Discoveryworld: A  
266 virtual environment for developing and evaluating automated scientific discovery agents. *arXiv  
267 preprint arXiv:2406.06769*, 2024. URL <http://arxiv.org/abs/2406.06769>. Accepted to  
268 NeurIPS 2024 (Benchmark Track, Spotlight).
- 269 Guanlin Jing and Huayi Qi. Zero-knowledge audit for internet of agents: Privacy-preserving com-  
270 munication verification with model context protocol. *arXiv preprint arXiv:2512.14737*, 2025.  
271 URL <http://arxiv.org/abs/2512.14737>.
- 272 Kousha Kalantari, Lalitha Sankar, and Anand Sarwate. Robust privacy-utility tradeoffs under differ-  
273 ential privacy and hamming distortion. *IEEE Transactions on Information Forensics and Security*,  
274 13(11):2816–2830, 2018. URL <http://arxiv.org/abs/1601.06426>.
- 275 Shiva Prasad Kasiviswanathan and Adam Smith. On the ‘semantics’ of differential privacy: A  
276 bayesian formulation. *Journal of Privacy and Confidentiality*, 6(1), 2014. URL [http://arxiv.  
277 org/abs/0803.3946](http://arxiv.org/abs/0803.3946).
- 278 James R. Kirk, Robert E. Wray, and John E. Laird. Exploiting language models as a source of  
279 knowledge for cognitive agents. 2023. URL <http://arxiv.org/abs/2310.06846>.
- 280 Masoud Koleini, Eike Ritter, and Mark Ryan. Verification of agent knowledge in dynamic access  
281 control policies. 7795:448–462, 2014. URL <http://arxiv.org/abs/1401.4730>.
- 282 Saikat Maiti. Caging the agents: A zero trust security architecture for autonomous ai in healthcare.  
283 *arXiv preprint arXiv:2603.17419*, 2026. URL <http://arxiv.org/abs/2603.17419>.
- 284 OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. URL [http://arxiv.  
285 org/abs/2303.08774](http://arxiv.org/abs/2303.08774).
- 286 Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Joseph E. Gonzalez, and  
287 Ion Stoica. Memgpt: Towards llms as operating systems. *arXiv preprint arXiv:2310.08560*, 2023.  
288 URL <http://arxiv.org/abs/2310.08560>.
- 289 Joon Sung Park, Joseph C. O’Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and  
290 Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. *arXiv preprint  
291 arXiv:2304.03442*, 2023. URL <http://arxiv.org/abs/2304.03442>.
- 292 Alireza Rezazadeh, Zichao Li, Ange Lou, Yuying Zhao, Wei Wei, and Yongfeng Zhang. Collab-  
293 orative memory: Multi-user memory sharing in llm agents with dynamic access control. *arXiv  
294 preprint arXiv:2505.18279*, 2025. URL <http://arxiv.org/abs/2505.18279>.
- 295 Lalitha Sankar, S. Raj Rajagopalan, and H. Vincent Poor. Utility-privacy tradeoff in databases: An  
296 information-theoretic approach. *IEEE Transactions on Information Forensics and Security*, 2013.  
297 URL <http://arxiv.org/abs/1102.3751>.
- 298 Akhil Sharma, Shaikh Yaser Arafat, Jai Kumar Sharma, and Ken Huang. Bilevel optimization  
299 for covert memory tampering in heterogeneous multi-agent architectures (xamt). *arXiv preprint  
300 arXiv:2512.15790*, 2025. URL <http://arxiv.org/abs/2512.15790>.
- 301 Theodore R. Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L. Griffiths. Cognitive ar-  
302 chitectures for language agents. *arXiv preprint arXiv:2309.02427*, 2023. URL [http://arxiv.  
303 org/abs/2309.02427](http://arxiv.org/abs/2309.02427).
- 304 Aleksei Triastcyn and Boi Faltings. Bayesian differential privacy for machine learning. 2020. URL  
305 <http://arxiv.org/abs/1901.09697>.

- 306 Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Aug-  
307 menting language models with long-term memory. *arXiv preprint arXiv:2306.07174*, 2023. URL  
308 <http://arxiv.org/abs/2306.07174>.
- 309 Yongqiang Wang. Privacy in multi-agent systems. 2024. URL <http://arxiv.org/abs/2403.02631>.
- 311 John Yang, Carlos E. Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, et al. Swe-  
312 agent: Agent-computer interfaces enable automated software engineering. *arXiv preprint*  
313 *arXiv:2405.15793*, 2024. URL <http://arxiv.org/abs/2405.15793>.
- 314 Yangyang Yu, Zhiyuan Yao, Haohang Li, Zhiyang Deng, Yupeng Cao, et al. Fincon: A synthesized  
315 llm multi-agent system with conceptual verbal reinforcement for enhanced financial decision  
316 making. *arXiv preprint arXiv:2407.06567*, 2024. URL <http://arxiv.org/abs/2407.06567>.
- 317 Zhonghan Zhao, Kewei Chen, Dongxu Guo, Wenhao Chai, Tian Ye, et al. Hierarchical auto-  
318 organizing system for open-ended multi-agent navigation. *arXiv preprint arXiv:2403.08282*,  
319 2024. ICLR 2024 Workshop on LLM Agents.
- 320 Changxi Zhu, Mehdi Dastani, and Shihan Wang. A survey of multi-agent deep reinforcement learn-  
321 ing with communication. *Autonomous Agents and Multi-Agent Systems*, 38(1):4, 2024. URL  
322 <http://arxiv.org/abs/2203.08975>.