# Intrinsic Memory Agents: Heterogeneous Multi-Agent LLM Systems through Structured Contextual Memory

**Anonymous authors**Paper under double-blind review

000

001

002

004

006

008 009 010

011 012

013

014

015

016

017

018

019

021

023

025

026

027

028

031

033

034

037

040

041

042

043

044

046

047

051

052

## **ABSTRACT**

Multi-agent systems built on Large Language Models (LLMs) show exceptional promise for complex collaborative problem-solving, yet they face fundamental challenges stemming from context window limitations that impair memory consistency, role adherence, and procedural integrity. This paper introduces Intrinsic Memory Agents, a novel framework that addresses these limitations through structured agent-specific memories that evolve intrinsically with agent outputs. Specifically, our method maintains role-aligned memory templates that preserve specialized perspectives while focusing on task-relevant information. We benchmark our approach on the PDDL dataset, comparing its performance to existing state-of-the-art multi-agentic memory approaches and showing an improvement of 38.6% with the highest token efficiency. An additional evaluation is performed on a complex data pipeline design task, and we demonstrate that our approach produces higher quality designs across 5 metrics: scalability, reliability, usability, cost-effectiveness, and documentation, plus additional qualitative evidence of the improvements. Our findings suggest that addressing memory limitations through structured, intrinsic approaches can improve the capabilities of multi-agent LLM systems on structured planning tasks.

#### 1 Introduction

Recent advances in large language models (LLMs) have enabled their application as autonomous or semi-autonomous agents capable of complex reasoning and decision-making (Huang et al., 2024). Multi-agent LLM systems, where multiple LLM instances interact to solve problems collaboratively, have shown particular promise for tasks requiring diverse expertise (Park et al., 2023; Qian et al., 2025). These systems leverage the complementary capabilities of specialized agents to address challenges that would be difficult for single-agent approaches to resolve effectively.

Despite their theoretical advantages, multi-agent LLM systems face several implementation challenges that limit their practical effectiveness, from coordination overhead, to the consistency in role adherence among the agents (Li et al., 2024c). Most critically, the fixed-size context windows of LLMs restrict their ability to maintain long-term conversational context, an issue that is exacerbated in multi-agent frameworks with multiple agents in a single conversation. This leads to issues such as perspective inconsistency, forgetting key requirements, and procedural drift. Current solutions such as Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Gao et al., 2024) and agentic memory approaches (Packer et al., 2024; Xu et al., 2025; Chhikara et al., 2025) are designed for single-agent and user interaction scenarios, which do not account for the volume of information growing with the number of agents.

To address these challenges, we introduce Intrinsic Memory Agents, a novel multi-agent architecture that uses structured, agent-specific memories aligned with conversational objectives. Unlike previous approaches, our system updates memories that are specific to each agent, ensuring heterogeneity and memories that reflect both historical context and recent developments while preserving agent-specific perspectives. The intrinsic nature of memory updates, derived directly from agent outputs rather than external summarization, ensures unique memories that maintain consistency with agent-specific reasoning patterns and domain expertise. We evaluate our approach through benchmarking

and through a specific data pipeline design case study to show its practical usage. The evaluation demonstrates that our Intrinsic Memory Agents approach yields significant improvements in conversational coherence, role consistency, and collaborative efficiency compared to conventional multi-agent implementations. These improvements translate to qualitative enhancements in solution quality without increasing the number of conversation turns, suggesting broad applicability across domains where multi-agent LLM systems are deployed.

The main contributions of our work are as follows:

- Structured Memory Templates: Predefined memory structures aligned with agent roles and conversational objectives.
- Intrinsic Memory Updates: Memory updates derived from agent outputs rather than external summarization.
- Agent-Specific Memory: Independent memories maintained for each agent to preserve perspective autonomy.

## 2 Related work

Recent years have seen significant progress in the development of multi-agent systems powered by LLMs. These systems have been applied in various domains, such as software development, scientific experimentation, gaming, and social simulation (Li et al., 2024c). For example, in software development, multi-agent systems enable concurrent consideration of architectural design, security, user experience, and performance optimization (Hong et al., 2024). Hallucinations due to outdated knowledge or retrieval extraction issues remains a major challenge which limits the effectiveness of multi-agent systems Huang et al. (2025). The use of a shared knowledge base or memory storage is an important aspect to maintain up-to-date, coherent and correct information among agents.

#### 2.1 Memory in Agent-based systems

In agent-based systems, memory is pivotal for maintaining context, learning from historical interactions, and making informed decisions. As Zhang et al. (2024) noted, memory supports tasks such as ensuring conversation consistency and effective role-playing for single-agent systems. In multiagent systems, memory facilitates coordination, communication, and collaborative problem-solving, as Guo et al. (2024) discussed.

Memory in LLMs can be categorized under short-term memory and long-term memory. Short-term memory is information that fits within the model's fixed context window. Commercial LLMs such as GPT-40 (OpenAI, 2024) and Claude (Anthropic, 2024) are able to process large contexts of over 100K tokens, with some models such as Gemini 2.5 Pro (Comanici et al., 2025) able to process over 1 million tokens in its context window. However, the hard limit of the context window size remains, and increasing the context length does not necessarily increase reasoning or learning capabilities of the LLM (Li et al., 2024b). This is because the long context can move the relevant information further away from each other in the context window.

Long-term memory is information that persists beyond the context window or single instance of an LLM. This information can be stored in external databases and retrieved using RAG techniques (Lewis et al., 2020; Gao et al., 2024). Long-term memory aims to alleviate the issue of short-term memory's limited capacity, but introduces other disadvantages such as retrieval noise, the complexity of building a retrieval system, latency, and storage costs (Asai et al., 2024; Yu et al., 2024).

The limitations of context length and existing memory mechanisms are particularly pronounced in multi-agent settings, where the volume of information exchanged grows with the number of agents involved Li et al. (2024a). As multi-agent conversations extend, the probability of critical information being at a long distance or even excluded from the accessible context increases dramatically. This information loss undermines the primary advantage of multi-agent systems: The integration of diverse, specialized perspectives toward cohesive solutions He et al. (2025). This is exacerbated by current long-term memory approaches which provide a homogeneous memory for the agents, decreasing the benefits of having agents focused on a single part of the task. Our proposed approach therefore focuses on the heterogeneity of agents and their memories, ensuring that each agent maintains a memory that is uniquely relevant to their role.

# 2.2 AGENTIC MEMORY

Agentic memory offers a solution to long-term memory and limited contextual information by periodically condensing conversation history into concise summaries (Wang et al., 2025; Chen et al., 2024). These approaches generate sequential or hierarchical summaries that capture key decisions and insights from previous exchanges. Some agentic memory approaches combine with RAG approaches by storing the summarized contexts for retrieval later in the conversation (Xu et al., 2022), or by storing in- and out-of-context memory in a hierarchical system to dynamically adapt the current context (Packer et al., 2024; Xu et al., 2025). While agentic memory methods provide better contextual integration than pure retrieval approaches, they frequently lose critical details during the condensation process. Furthermore, the undirected and unstructured nature of general summarization often fails to preserve role-specific perspectives and specialized knowledge that are essential to effective multi-agent collaboration.

Our proposed Intrinsic Memory Agents similarly uses an agentic memory approach to summarize and store information. Unlike existing approaches, we introduce structured heterogeneous memory for each agent in the multi-agent system to maintain specialized roles in collaborative tasks, and apply a structured template to each agent ensuring a cohesive memory structure. This addresses the limitations of existing memory mechanisms by ensuring that each agent maintains its own memory, reflecting both historical context and new information while maintaining heterogeneous agent-specific perspectives and expertise.

# 3 Intrinsic Memory Agents

The various agentic memory approaches are all designed in single-agent scenarios to remember crucial details when interacting with an end-user. Due to the multi-turn long conversations between agents, a direct implementation of single-agent agentic memory becomes complicated and resource-intensive, with each agent requiring retrieval systems and distinct contextual updates.

We propose Intrinsic Memory Agents, a framework for multi-agent LLM systems that maintains agent-specific structured memories aligned with conversational objectives. Figure 1 illustrates the architecture of our Intrinsic Memory Agents framework. In this approach a query is made by the user, the first agent makes a comment based on its role description, the conversation is updated, followed by a memory update for the agent that commented, there is a check for consensus and the cycle starts again. The context in this case is made up of both the agent's intrinsic memory and the conversation, meaning that as the conversation continues the agents increasingly diverge in their interpretation of that context.

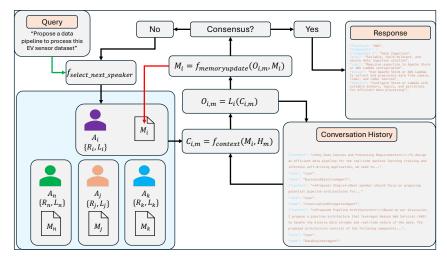


Figure 1: Intrinsic Memory Agents Framework. For n agents and m conversation turns, each agent  $A_n$  contains its own role description  $R_n$  and language model  $L_n$ . Its memory  $M_{n,m}$  is updated based on the input context  $C_{n,m}$  and output  $O_{n,m}$ .

#### 3.1 Framework Definition

Let us define the multi-agent system  $\mathcal{A} = \{A_1, A_2, ..., A_N\}$  consisting of N agents. Each agent  $A_n = \{R_n, M_n, LLM_n\}$  is characterized by a role specification  $R_n$  that defines the agent's expertise domain and objectives, a structured memory  $M_n$  that evolves throughout the conversation, and an LLM instance  $LLM_n$ , which may share parameters between agents.

The conversation consists of a sequence of turns  $T=t_1,t_2,...,t_M$  where each turn  $t_m$  involves an agent selection function  $\sigma(t_m)\to A_n$  that determines which agent speaks, an input context  $C_{n,m}$  constructed for the selected agent, an output  $O_{n,m}$  generated by the selected agent, and a memory update operation for the selected agent.

Critically, our framework separates the input context construction and memory update processes, allowing for agent-specific memory maintenance while preserving a shared conversation space.

#### 3.2 STRUCTURED MEMORY TEMPLATES

For each agent  $A_n$ , we define a structured memory template  $MT_n$  that specifies the organization of agent-specific memories. The template consists of a set of memory slots  $MT_n = \{S_1, S_2, ..., S_K\}$ . Each slot  $S_k$  is defined by a descriptive identifier (e.g., "domain\_expertise", "current\_position", "proposed\_solution"), formatted in JSON. The template can be nested so that each slot can have its own descriptive identifier with further detailed information. The structured nature of the memory templates ensures that updates remain focused on role-relevant information while maintaining consistency with the agent's expertise domain.

#### 3.3 Memory Update mechanism

For each agent in the system, we maintain a structured memory  $M_n$  that evolves over time. Let  $M_{n,m}$  represent the memory of agent n after m conversation turns. The memory update process works as follows:

Agent  $A_n$  receives input context  $C_{n,m}$  consisting of relevant conversation history  $H_m$  and previous memory  $M_{n,m-1}$ ,

$$C_{n,m} = f_{\text{context}}(H_m, M_{n,m-1}); \tag{1}$$

and agent  $A_n$  generates output  $O_{n,m}$  using the underlying LLM  $L_n$ ,

$$O_{n,m} = L_n(C_{n,m}). (2)$$

Then with the generated output  $O_{n,m}$  and the previous memory  $M_{n,m-1}$ , we update the slot content using a memory update function,

$$M_{n,m} = f_{\text{memory\_update}}(M_{n,m-1}, O_{n,m}). \tag{3}$$

The memory update function  $f_{\mathrm{memory.update}}$  is implemented as a prompted LLM operation. Specifically, for the previous memory  $M_{n,m-1}$  at turn m-1 and agent output  $O_{n,m}$  at turn m, the update function constructs the prompt as shown in Figure 4. The LLM's response to this prompt becomes the updated memory  $M_{n,m}$ . The context construction function  $f_{\mathrm{context}}$  presented in equation 1 determines what information is provided to an agent when generating a response. The algorithm takes the existing conversation history and agent memory, appending both to the context and using the remaining tokens to include the rest of the conversation history. The full algorithm pseudo-code is displayed in the Appendix A Algorithm 1.

This algorithm prioritizes:

- 1. The initial task description to maintain objective alignment.
- 2. The agent's structured memory to preserve role consistency.
- 3. The most recent conversation turns to maintain immediate context.

By prioritizing memory inclusion over exhaustive conversation history, the algorithm ensures that agents maintain role consistency and task alignment even when conversation length exceeds context window limitations.

# 4 PDDL BENCHMARK

To evaluate our approach, we test our memory agents against the PDDL (Planning Domain Definition Language) numeric benchmark. PDDL involves structured planning tasks from AgentBoard (Ma et al., 2024), where the agents generate executable plans for abstract problem domains, evaluating their reasoning and coordination.

For numerical benchmarks, we follow the same experimental methodology as G-Memory (Zhang et al., 2025), another memory framework for multi-agent systems. We re-run the G-Memory framework <sup>1</sup> as we cannot directly compare to the published G-Memory results which were benchmarked with GPT-40-mini as the base language model. We chose to use the G-Memory framework as a comparison as the framework implements a variety of existing memory architectures, allowing us to compare our Intrinsic Memory Agents with existing architectures and benchmarks. G-Memory uses Autogen for multi-agent simulation, matching our use of Autogen for our architecture. We chose to use the PDDL dataset as its structured planning task is the intended use case for the Intrinsic Memory Agents and aligns with the data pipeline case study detailed in Section 5. We run Llama3.1:8b for the numeric benchmarks using Ollama <sup>2</sup> with 1 repetition. We use a larger model for the numeric benchmarks as initial tests on the 3b model found poor results for every benchmark and memory framework. We use a single run with a set seed for reproducibility. Our computational infrastructure utilizes a high performance computing cluster with A100 GPUs, running on GNU/Linux 4.18.0-553.el8\_10.x86\_64.

#### 4.1 BENCHMARKING RESULTS

Memory	Average rewards	Average tokens
No Memory	0.0231	117,437
G-Memory	0.0231	118,316
Generative	0.0530	121,105
MemoryBank	0.0305	80,599
MetaGPT	0.0601	127,860
Voyager	0.0250	133,268
ChatDev	0.0583	107,043
Intrinsic Memory	0.0833	140,418

Table 1: Performance comparison between different memory architectures on PDDL. We record the average rewards per task and average number of tokens used of per task. The best results for each benchmark are highlighted in **bold**.

Table 1 shows the rewards and tokens used for different memory architectures with a multi-agent system on the PDDL benchmark problems. We find substantially better average rewards for PDDL at 0.0833 compared to the next highest with MetaGPT at 0.0601. The improved rewards come at a cost of increased token usage, the highest among all memory architectures, requiring 12,558 more tokens than MetaGPT. Based on token efficiency, defined as average reward per token, Intrinsic Memory shows the best result at  $5.933 \times 10^{-7}$  with ChatDev coming second, which has a token efficiency of  $5.466 \times 10^{-7}$ .

The PDDL dataset are structured planning tasks, which fits the intended use case of Intrinsic Memory for agent discussion, planning and design. As Intrinsic Memory assigns agent-specific memory, it can more clearly distinguish planning and actions to complete tasks. More tokens are used by Intrinsic Memory to generate structured templates per agent per round of discussion, and is a worth-while trade-off in both reward score and token efficiency. Notably, all memory architectures except ChatDev and Intrinsic Memory, the most token efficient methods, utilize cross-trial memory, in which memory is stored and carried across tasks, which significantly helps give few-shot examples when agents are stuck in a loop in the environment.

<sup>&</sup>lt;sup>1</sup>https://github.com/bingreeky/GMemory

<sup>&</sup>lt;sup>2</sup>https://ollama.com/library

# 5 DATA PIPELINE DESIGN CASE STUDY

As a practical case study to evaluate our approach, we applied our memory agents to a collaborative data pipeline design, a complex task requiring multiple perspectives. We run 10 independent outputs with eight specialized agents:

1. Evaluation Agent (EA) evaluates the output solutions.

2. **Knowledge Integration Agent (KIA)** summarizes each discussion round (e.g. after every agent has contributed at least once).

3. Data Engineer Agent (DEA) determines the data processing needs.

4. Infrastructure Engineer (IA) designs the cloud infrastructure.

 $5. \ \textbf{Business Objective Engineer (BOA)} \ \text{checks against business requirements}.$ 

 6. Machine Learning Engineer (MLE) provides ML implementation.

 Conversation Delegation Agent (CDA) is responsible for facilitating the collaborative process.

8. **Documentation Joining Agent (DJE)** is responsible for producing final output after consensus is reached among agents.

The agents are tasked with designing a cloud-based data pipeline architecture through a structured process involving proposals, discussions, and consensus formation. The full prompts and task descriptions can be found in Appendix B.

The output requirements include a concise summary, high-level plan, resource estimates, and a structured JSON specification.

## 5.1 System Configurations

We evaluated two system configurations: First, the **Baseline System** which consists of a standard multi-agent implementation without structured memory. It uses standard prompt templates for each agent role, relying exclusively on conversation history for context. This limits the most recent conversation turns due to context window constraints. Second is our **Intrinsic Memory System** approach with agent-specific structured memories. It implements role-specific memory templates, updates memories intrinsically based on agent outputs, and constructs context using both conversation history and structured memories.

Both systems used identical agent roles and task specifications, with Llama-3.2-3b as the underlying LLM. Each agent role was initialized with the same role description and initial instructions across the two system configurations to ensure a fair comparison.

An agent selection function iterates through each worker agent and the conversation delegation agent (CDA), ensuring that all agents are represented in the discussion. Once all agents have accepted a proposed solution, marked through the "ACCEPT" flag, the CDA emits a "FINALIZE" flag, prompting the Documentation Engineer Agent to produce the final data pipeline output. The full algorithm for finalisation and ordering of agents is displayed in the Appendix A Algorithm 2.

#### 5.2 EVALUATION METRICS

To evaluate the quality of the data pipeline designs generated by our memory agents, and to compare to the pipeline designs generated from default Autogen, we use an LLM as a judge (Zheng et al., 2023) to score each pipeline design and provide a qualitative analysis to support these scores. We evaluated the multi-agent system performance under the following metrics:

• Scalability: ability for the data pipeline to handle increasing data volumes or user loads.

• Usability: is there enough detail in the data pipeline design for developers to implement the design?

• Cost-effectiveness: balance between costs and benefits of each chosen component.

• Reliability: ability for the data pipeline to handle failures and ensure data integrity.

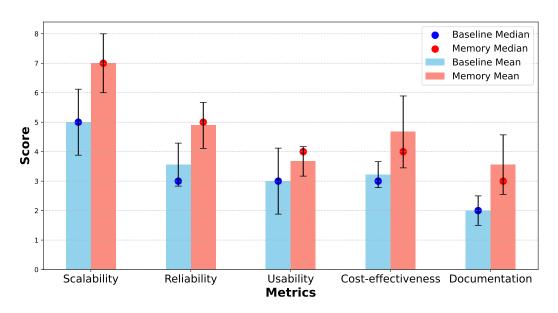


Figure 2: LLM-as-a-Judge metrics for the Data Pipeline design case study.

 Documentation: how well-justified and documented is the choice of elements for the data pipeline?

The scalability and reliability metrics are chosen as core requirements in modern data pipelines. Scalability reflects the ability to grow the pipeline to handle larger volumes of data and users, while reliability ensures the pipeline is consistent and fault-tolerant, both of which are crucial if the pipeline were to be deployed. Usability and documentation metrics reflect the details and design decisions taken. A strong design is not useful if it does not contain enough detail or is too abstract to be practically implemented. Usability measures whether the output designs are detailed and clear enough for engineering teams to implement. Design decisions must be well-documented, with clear justifications and explanations for each component, which reveals the reasoning behind the agents' choices. Finally, the cost-effectiveness metric evaluates whether the pipeline design has considered and balanced the need for computation resources with the cost of those resources. Run-time metrics such as latency and throughput are not included in our evaluation metrics as we only present the design of the data pipelines to be evaluated, and do not implement the designs into code

# 5.3 DATA PIPELINE DESIGN PERFORMANCE

The median and standard deviation of each quality metric is presented in Figure 2. The Intrinsic Memory system shows consistent improvement on all metrics compared to the baseline Autogen, with Usability as the only metric where there is not a statistically significant difference.

The Documentation quality focuses on the clarity and how well-justified the design choices are. While Intrinsic Memory helps to boost the Documentation score over the baseline, the score is still relatively low at a mean of 3.56. This suggests that retaining memory of the conversation alone does not guarantee good justification, and while some context and attributes of each component are remembered, the reasons for choosing the components are not. This could be a problem with the training corpus, and a requirement for better annotated training data. Similarly, the Usability score is low with means of 3 and 3.67 for the baseline and Intrinsic Memory , respectively.

The improved quality comes at a cost of additional tokens outlined in Table 2. Intrinsic Memory uses on average 32% more tokens than the baseline as it outputs are more descriptive on average, although the number of conversation turns is similar and not statistically significant. This indicates that the addition of a memory module costs additional token overhead to maintain, but does not increase the number of conversation turns between agents.

Table 2: Mean efficiency and LLM-as-a-Judge metrics after 10 independent runs, with p-values calculated using a Wilcoxon ranked sum test. Usability and number of conversation turns is highlighted in italics as the metrics that do not show statistical significance between the baseline and our Intrinsic Memory approach.

Metric	Baseline	Intrinsic	p-value
	Autogen	Memory	
Tokens	36077	47830	0.0195
Conversation turns	14.3	16	0.2632
Scalability	5	7	0.0041
Reliability	3.56	4.89	0.005
Usability	3	3.67	0.0948
Cost-effectiveness	3.22	4.67	0.004
Documentation	2	3.56	0.0017

#### 5.4 QUALITATIVE ANALYSIS OF DATA PIPELINE OUTPUTS

Figure 3 shows snippets for one component from the highest-scoring outputs for the intrinsic memory agent system and baseline Autogen system.

```
"Component 1": "Data Ingestion (AWS Kinesis)",
  "pros": "Scalable, fault-tolerant, and supports
high-throughput data ingestion from multiple sources.",
  "cons": "Requires additional infrastructure costs for
increased capacity.",
  "design": "Use Kinesis to collect and process data from camera
sources, lidar source, radar sources, and other IoT devices.",
  "details": "Implement custom Kinesis data stream for each data
source, with appropriate buffering and queuing mechanisms to
handle varying data volumes"
```

(a) Intrinsic Memory Agent system sample from highest-scoring output. This data pipeline received scores of Scalability: 8, Reliability: 5, Usability: 4, Cost-effectiveness: 6, Documentation: 4

```
"Component 1":
  "Name": "Data Ingestion",
  "Description": "Ingest data from various sources (camera, lidar, radar) at high speeds",
  "Implementation difficulties": 7,
  "Maintainability difficulties": 6
```

(b) Baseline Autogen system system sample from highest-scoring output. This data pipeline received scores of Scalability: 5, Reliability: 4, Usability: 3, Cost-effectiveness: 3, Documentation: 2

Figure 3: Snippets of one component within the data pipeline design from both systems. The full outputs can be found in the appendix in Figures 7 and 8.

The Intrinsic Memory Agent system outperforms the baseline system across the five quality metrics. In terms of scalability, the Intrinsic Memory Agent system is capable of providing an overall assessment of scalability, specifically around varying data volumes, whereas the baseline system encapsulates that measure only in the form of "maintenance difficulty" for each component of the pipeline. In terms of reliability, the Intrinsic Memory Agent provides considerations for AWS Kinesis's fault tolerance, as well as the need for appropriate buffering and queuing mechanisms to handle varying data volumes. The Intrinsic Memory Agent provides a more descriptive Usability output of the Intrinsic and a clearer pathway to implementation. Neither system makes specific observations for the cost-effectiveness on an individual component basis, but the Intrinsic Memory Agent does provide an overall numerical evaluation of the pipeline's cost-effectiveness. Finally, the Intrin-

sic Memory Agent ultimately provides justification and documents its recommendation under each component, including pros and cons for each component choice.

Overall, the Intrinsic Memory Agent provides a more descriptive answer and more value to engineers by specifying tools, configurations and trade-offs. For example, its Data Ingestion design recommends Apache Storm or AWS Lambda, whereas the baseline simply states "Ingest data from various sources (camera, lidar, radar) at high speeds." Similarly, the IMA cites OpenCV and TensorFlow for image processing, PCL and Open3D for point-cloud handling, and MATLAB plus machine-learning libraries for radar signals. Although some precise configuration settings remain unspecified, the baseline merely names each component without offering implementation details or alternatives.

The components specified by the Intrinsic Memory Agent are more relevant to the problem specification. The data pipeline design task explicitly specifies the input data contains lidar and radar data sources, in which the PCL and Open3D are libraries specifically used for lidar data processing. This contrasts the vague "Lidar data processing" component of the baseline, which only contains a general description to process the data without providing any details.

# 6 DISCUSSION AND LIMITATIONS

Although the Intrinsic Memory Agent approach shows improved performance across the data pipeline generation task and the selected benchmarks, further validation is required across a broader set of complex tasks, potentially with varying number of agents, models, and memory templates. The structured memory templates are currently created manually, which does not easily transfer across tasks. An automated or generalized method to producing structured memory templates would improve the Intrinsic Memory Agent's ability to adapt to new tasks.

The results demonstrate that a movement towards heterogeneity of agents leads to an improvement in performance of the multi-agent system, allowing agents to focus more specifically on an area of the design. This indicates that methods to provide additional heterogeneity, such as the ability to fine-tune agents towards their specialisation, might see additional performance gains, alongside the personalization of memories focused on individual experience.

#### 7 CONCLUSION

This paper introduces Intrinsic Memory Agents, a novel multi-agent LLM framework that constructs agent-specific heterogeneous structured memories to enhance multi-agent collaboration in discussion and planning tasks. Evaluation on the PDDL dataset, and on a practical data pipeline design problem demonstrates our framework's improved performance on structured planning tasks, with a 38% increase over the next best memory architecture, and maintaining the best token efficiency despite the increased token usage.

Results on the data pipeline case study further show the Intrinsic Memory Agents' enhanced ability to collaborate on complex tasks. The Intrinsic Memory Agent system outperforms the baseline system across all quality measures of scalability, reliability, usability, cost-effectiveness, and documentation, as well as an ability to more closely follow the task specification, providing more actionable recommendations by suggesting specific tools and frameworks, as well as trade-off details of each component in the pipeline.

# REPRODUCIBILITY STATEMENT

We have taken care to ensure that our experiments and results are transparent and reproducible by detailing the models, computational setup, code, statistical tests, and prompts used in our experiments. The LLM model (Llama3.2:3b) is cited, named, and referenced in the main text. The computational infrastructure used, including the GPU model names and operating system are specified in section 4 of the main text. For code, the names and versions of relevant Python libraries are specified within the supplementary code files. A Wilcoxon rank-sum test is used to test statistical significance for the data pipeline case study. P-values and standard deviation measures are included in the performance analysis in section 5.3. A single run with a set seed is used for the PDDL benchmark, with the seed specified within the supplementary code. All code for running the data pipeline case study and PDDL benchmarks are included in the supplementary materials. Finally, the selected prompts of the multi-agent and Intrinsic memory architecture are shown in Appendix B. Further prompts for each agent can be found as part of the supplementary code, under the "prompts" directory.

# REFERENCES

- Anthropic. Claude 3.5 sonnet. Technical announcement, 2024. URL https://www.anthropic.com/news/claude-3-5-sonnet. Available via Claude.ai and API.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *The 12th International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=hSyW5go0v8.
- Nuo Chen, Hongguang Li, Juhua Huang, Baoyuan Wang, and Jia Li. Compress to impress: Unleashing the potential of compressive memory in real-world long-term conversations, 2024. URL https://arxiv.org/abs/2402.11975.
- Prateek Chhikara, Dev Khant, Saket Aryan, Taranjeet Singh, and Deshraj Yadav. Mem0: Building production-ready ai agents with scalable long-term memory. *arXiv preprint arXiv:2504.19413*, 2025.
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornraphop Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilaï Deutel, Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller, Vaishakh Keshava, Ying Jian, Xiaofan Zhang, Raluca Ada Popa, Kedar Dhamdhere, Blaž Bratanič, Kyuyeun Kim, Terry Koo, Ferran Alet, Yi ting Chen, Arsha Nagrani, Hannah Muckenhirn, Zhiyuan Zhang, Corbin Quick, Filip Pavetić, Duc Dung Nguyen, Joao Carreira, Michael Elabd, Haroon Qureshi, Fabian Mentzer, Yao-Yuan Yang, Danielle Eisenbud, Anmol Gulati, Ellie Talius, Eric Ni, Sahra Ghalebikesabi, Edouard Yvinec, Alaa Saade, Thatcher Ulrich, Lorenzo Blanco, Dan A. Calian, Muhuan Huang, Aäron van den Oord, Naman Goyal, Terry Chen, Praynaa Rawlani, Christian Schallhart, Swachhand Lokhande, Xianghong Luo, Jyn Shan, Ceslee Montgomery, Victoria Krakovna, Federico Piccinini, Omer Barak, Jingyu Cui, Yiling Jia, Mikhail Dektiarev, Alexey Kolganov, and Shiyu Huang. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. URL https://arxiv.org/abs/2507.06261.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey, 2024. URL https://arxiv.org/abs/2312.10997.
- 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, 2024. URL https://arxiv.org/abs/2402.01680.

- Junda He, Christoph Treude, and David Lo. Llm-based multi-agent systems for software engineering: Literature review, vision and the road ahead, 2025. URL https://arxiv.org/abs/2404.04834.
  - Sirui Hong, Mingchen Zhuge, Jiaqi Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, Chenyu Ran, Lingfeng Xiao, Chenglin Wu, and Jürgen Schmidhuber. Metagpt: Meta programming for a multi-agent collaborative framework, 2024. URL https://arxiv.org/abs/2308.00352.
  - Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Trans. Inf. Syst.*, 43(2), January 2025. ISSN 1046-8188. doi: 10.1145/3703155. URL https://doi.org/10.1145/3703155.
  - Xu Huang, Weiwen Liu, Xiaolong Chen, Xingmei Wang, Hao Wang, Defu Lian, Yasheng Wang, Ruiming Tang, and Enhong Chen. Understanding the planning of llm agents: A survey, 2024. URL https://arxiv.org/abs/2402.02716.
  - Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33: 9459–9474, 2020.
  - Junyou Li, Qin Zhang, Yangbin Yu, Qiang Fu, and Deheng Ye. More agents is all you need, 2024a. URL https://arxiv.org/abs/2402.05120.
  - Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. Long-context llms struggle with long in-context learning, 2024b. URL https://arxiv.org/abs/2404.02060.
  - Xinyi Li, Sai Wang, Siqi Zeng, Yu Wu, and Yi Yang. A survey on llm-based multi-agent systems: workflow, infrastructure, and challenges. *Vicinagearth*, 1(1):9, October 2024c. ISSN 3005-060X. doi: 10.1007/s44336-024-00009-2. URL https://doi.org/10.1007/s44336-024-00009-2.
  - Chang Ma, Junlei Zhang, Zhihao Zhu, Cheng Yang, Yujiu Yang, Yaohui Jin, Zhenzhong Lan, Lingpeng Kong, and Junxian He. Agentboard: An analytical evaluation board of multi-turn llm agents, 2024.
  - OpenAI. Gpt-4o system card, 2024. URL https://arxiv.org/abs/2410.21276.
  - Charles Packer, Sarah Wooders, Kevin Lin, Vivian Fang, Shishir G. Patil, Ion Stoica, and Joseph E. Gonzalez. Memgpt: Towards llms as operating systems, 2024. URL https://arxiv.org/abs/2310.08560.
  - Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior, 2023. URL https://arxiv.org/abs/2304.03442.
  - Chen Qian, Zihao Xie, YiFei Wang, Wei Liu, Kunlun Zhu, Hanchen Xia, Yufan Dang, Zhuoyun Du, Weize Chen, Cheng Yang, Zhiyuan Liu, and Maosong Sun. Scaling large language model-based multi-agent collaboration, 2025. URL https://arxiv.org/abs/2406.07155.
  - Qingyue Wang, Yanhe Fu, Yanan Cao, Shuai Wang, Zhiliang Tian, and Liang Ding. Recursively summarizing enables long-term dialogue memory in large language models, 2025. URL https://arxiv.org/abs/2308.15022.
  - Jing Xu, Arthur Szlam, and Jason Weston. Beyond goldfish memory: Long-term open-domain conversation. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5180–5197, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.356. URL https://aclanthology.org/2022.acl-long.356/.

Wujiang Xu, Kai Mei, Hang Gao, Juntao Tan, Zujie Liang, and Yongfeng Zhang. A-mem: Agentic memory for llm agents, 2025. URL https://arxiv.org/abs/2502.12110. Tian Yu, Shaolei Zhang, and Yang Feng. Auto-rag: Autonomous retrieval-augmented generation for large language models, 2024. URL https://arxiv.org/abs/2411.19443. Guibin Zhang, Muxin Fu, Guancheng Wan, Miao Yu, Kun Wang, and Shuicheng Yan. G-memory: Tracing hierarchical memory for multi-agent systems, 2025. URL https://arxiv.org/ abs/2506.07398. Zeyu Zhang, Xiaohe Bo, Chen Ma, Rui Li, Xu Chen, Quanyu Dai, Jieming Zhu, Zhenhua Dong, and Ji-Rong Wen. A survey on the memory mechanism of large language model based agents, 2024. URL https://arxiv.org/abs/2404.13501. Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL https://arxiv.org/ abs/2306.05685. 

# A ALGORITHMS

This section contains the context construction algorithm presented in Section 3 and the finalisation algorithm presented in 5.1.

**Algorithm 1** The context construction algorithm, which takes the current conversation history, memory of the agent, and maximum number of tokens. It appends the most recent conversation turn and agent memory to the context first, before using the remainder of the tokens to append the rest of the conversation history, ensuring the memory and most recent output is always included.

```
def construct_context(conversation_history ,
                       agent_memory,
                       max_tokens):
context = []
# Include the initial task description
context.append(conversation_history[0])
context.append(agent_memory)
# Add most recent conversation turns until context limit is reached
remaining_tokens = max_tokens - count_tokens(context)
recent_turns = []
for turn in reversed (conversation_history[1:]):
    turn_tokens = count_tokens(turn)
    if turn_tokens <= remaining_tokens:</pre>
        recent_turns.insert(0, turn)
        remaining_tokens -= turn_tokens
    else:
        break
context.extend(recent_turns)
return context
```

Algorithm 2 The finalisation algorithm that specifies the order of agents speaking. It is a modified round-robin discussion between the agents: The discussion begins with each of the worker agents (BOA, DEA, MLA, IA) contributing to the conversation, with each worker's turn being followed by the conversation delegation agent (CDA). Once the workers have each had their turn, the knowledge integration agent and evaluation agent make their contributions, and the cycle begins again. The CDA is programmed to dedicate a certain number of turns to discussion, proposals, and consensus. The number of turns dedicated to each conversation stage is tracked, and once the consensus round is reached, each agent is asked to confirm if they agree with the proposed solution or not. If all agree on the proposed solution as being acceptable, the CDA will emit a "FINALIZE" response, triggering the documentation joining agent (DJE) to compile the agreed response and format it according to the task requirements.

```
workers = [BOA, DEA, MLA, IA]
global turn_counter
turn_counter += 1
if "FINALIZATION" in groupchat.messages[-1]['content']:
    return DJE
if last_speaker is CDA:
    global worker_counter
    w = workers [worker_counter%4]
    worker_counter += 1
    print(f'worker_counter: \[ \] \{ worker_counter \}')
    return w
elif last_speaker in workers:
    if worker_counter %4 == 0:
        return KIA
    else:
        return CDA
elif last_speaker is KIA:
        return ERA
elif last_speaker is ERA:
        return CDA
```

# B FULL PROMPTS AND EXAMPLE OUTPUTS

You are maintaining the memory of an agent working as [ROLE] in a multi-agent conversation. Use the old memory and the newest output by the agent to populate and up- date the current memory json with factual information.

For context, old memory content:
[MEMORY\_CONTENT]

Current content generated by the agent:
[AGENT\_OUTPUT]

Update the memory content to incorporate new information while preserving key historical context. The updated content should be concise and focus on information relevant to both the old memory and the newly generated output.

Figure 4: Prompt of the memory update function, where ROLE is the agent's role specification  $R_n$ ;  $MEMORY\_CONTENT$  is the current content  $M_{n,m-1}$ ;  $AGENT\_OUTPUT$  is the agent's output  $O_{n,m}$ .

858 859

860

```
812
813
814
          ''' This discussion session is set up to discuss the best data
         pipeline for a real time data intensive machine learning training
815
         and inference self driving application. The goal is to discuss and
816
         find consensus on how to set up the data pipeline, including each
817
         component in the data pipeline.
818
         You can assume that we have access to AWS.
819
         **Data Description: ** Real-time data of cars driving in street.
820
         There are 6 camera sources with data in .jpg format; 1 lidar source
821
         in .pcd.bin format; and 5 radar sources with data in .pcd format.
822
823
         **Discussion and Design:**
824
         - Emphasize comprehensive understanding of the data sources,
         processing requirements, and desired outcomes.
825
          - Encourage each other to engage in an open discussion on potential
826
         technologies, components, and architectures that can handle the
827
         diverse data streams and real-time nature of the data.
828
         - Keep the conversation on design and evaluating the pros and
829
         cons of different design choices, considering scalability,
         maintainability, and cost-effectiveness.
830
         - The team should agrees on a final architectural design,
831
         justifying the choices made.
832
          - The team should produce the required the document
833
         PIPELINE_OVERVIEW.json.
834
         **Final Output:**
835
         - Produce a concise summary of the agreed-upon pipeline
836
         architecture, highlighting its key components and connections.
837
          - Provide a high-level plan and rationale for the design,
838
         explaining why it is well-suited for the given data and use case.
839
         - Estimate the cloud resources, implementation efforts, and
         associated costs, providing a rough breakdown and complexity
840
         rating.
841
         - Generate a 'PIPELINE_OVERVIEW.json' file, detailing the proposed
         complete architecture in JSON format with the following fields:
843
         - \Platform\: A cloud service provider's name if the cloud
844
         solution is the best, or \local server" if locally hosted servers
         are preferred.
845
         - \Component 1":
                           The first component in the pipeline framework.
846
         - \component 2": The second component in the pipeline framework.
847
         Continue until all required components are listed.
848
         - \Implementation difficulties": A rating from 1 to 10 (lowest to
849
         - \Maintenance difficulties": A rating from 1 to 10 (lowest to
850
         highest).
851
852
         **Instructions:**
853
         - Remember, this is a collaborative design discussion, not a
854
         project execution. Refrain from assigning tasks with deadlines.
          - Keep the conversation focused on architectural choices,
855
         technologies, and potential challenges.
856
          - Emphasize the importance of a well-thought-out design.
857
```

Figure 5: Prompt for the Data Pipeline design. This is the task description provided to all agents in the system.

```
868
869
870
871
872
873
874
875
876
877
         , , ,
878
        You are an expert in data pipeline design evaluation.
879
        task is to evaluate the following data pipeline design based
880
        on its description. For each of the specified metrics,
881
        assign a score from 1 to 10, where 1 is the lowest and 10
882
        is the highest. Provide a brief justification for each
883
        score. Be critical and harsh if the design is poor and give
884
        it a low score. Base your evaluation solely on the provided
885
        description. Do not assume any additional information.
886
887
        Metrics:
        1. Scalability:
                           Ability to handle increasing data volumes
888
        or user loads.
889
        2. Reliability: Ability to handle failures and ensure data
890
        integrity.
891
        3. Usability: Enough detail for developers to implement the
892
        design.
893
        4. Cost-effectiveness: Balance between costs and benefits.
894
        5. Documentation: How well-justified is the choice of
895
        elements for the data pipeline
897
        Provide your evaluation in the following format in a json
898
        dict:
         [metric1]: score:
                            [score],
899
         justification: [justification],
900
         [metric2]...
901
902
903
```

Figure 6: Prompt for the LLM-as-a-Judge to evaluate data pipeline designs.

965

966

967

968

```
921
922
          "Platform":
                      "AWS",
          "Component 1": "Data Ingestion (AWS Kinesis)",
923
         "pros": "Scalable, fault-tolerant, and supports high-throughput
924
         data ingestion from multiple sources.",
925
         "cons": "Requires additional infrastructure costs for increased
926
         capacity.",
         "design": "Use Kinesis to collect and process data from camera
927
         sources, lidar source, radar sources, and other IoT devices.",
928
         "details": "Implement a custom Kinesis data stream for each data
929
         source, with appropriate buffering and queuing mechanisms to handle
930
         varying data volumes."
931
         "Component 2": "Data Processing (AWS Lambda)",
932
         "pros": "Serverless, scalable, and supports real-time data
933
         processing with event-driven architecture.",
934
         "cons": "Requires careful function design and optimization to
935
         ensure efficient processing.",
936
         "design": "Use Lambda functions to process data from Kinesis
937
         streams, leveraging computer vision libraries for image processing,
         point cloud reconstruction frameworks for lidar source data, and
938
         radar signal processing libraries for radar source data.",
939
         "details": "Implement separate Lambda functions for each data
940
         source, with optimized function designs to handle varying data
941
         volumes and processing requirements."
942
         "Component 3": "Data Storage (Amazon S3)",
943
         "pros": "Scalable, durable, and supports high-performance data
944
         storage for large datasets.",
945
         "cons": "Requires careful data management and access control
946
         strategies.",
         "design": "Use S3 as a centralized data repository for storing
947
         processed data from Kinesis streams.",
948
         "details": "Implement a custom S3 bucket structure to organize
949
         data by source, with appropriate access controls and data
         encryption mechanisms."
951
         "Component 4": "Data Synchronization (AWS Step Functions)",
952
         "pros": "Supports real-time data synchronization across different
953
         data sources, with built-in workflow management and error
954
         handling.'
955
         "cons": "Requires careful workflow design and optimization to
956
         ensure efficient data processing.",
957
         "design": "Use Step Functions to synchronize data from Kinesis
         streams, leveraging Lambda functions for data processing and S3 for
958
         data storage.",
959
         "details": "Implement a custom Step Function workflow to handle
960
         data synchronization, with optimized function designs to handle
961
         varying data volumes and processing requirements."
962
963
964
```

Figure 7: Example output from intrinsic memory agent system which scored the highest. The intrinsic memory agent system adheres to the initial prompt well, providing an answer in the correct format, as well as providing sufficient justification for the choice of components, as well as implementation details.

1022 1023

```
973
974
          "Platform": "AWS",
975
          "Component 1":
976
          "Name": "Data Ingestion",
977
          "Description":
                         "Ingest data from various sources (camera, lidar,
978
          radar) at high speeds",
          "Implementation difficulties": 7,
979
          "Maintainability difficulties":
980
         "Component 2":
981
         "Name": "Image Processing",
982
         "Description": "Process camera data in real-time to extract
983
         relevant features",
984
          "Implementation difficulties": 8,
          "Maintainability difficulties":
985
986
          "Component 3":
987
          "Name": "Object Detection",
988
         "Description": "Detect objects of interest (e.g., pedestrians,
989
         cars, lanes) within image frames and track their movement over
         time",
990
          "Implementation difficulties": 9,
991
          "Maintainability difficulties": 8
992
993
          "Component 4":
          "Name": "Lidar Data Processing",
994
         "Description": "Process lidar data in real-time to estimate
995
         distances, velocities, and other relevant metrics",
996
         "Implementation difficulties": 8,
997
          "Maintainability difficulties":
998
          "Component 5":
999
          "Name": "Radar Data Processing",
1000
          "Description": "Process radar data in real-time to estimate
1001
         distances, velocities, and other relevant metrics",
1002
          "Implementation difficulties": 8,
1003
         "Maintainability difficulties":
1004
          "Component 6":
1005
          "Name": "Data Fusion",
          "Description": "Fuse the outputs from different sensors (camera,
1007
         lidar, radar) to create a more accurate representation of the
1008
         environment",
          "Implementation difficulties": 9,
1009
          "Maintainability difficulties":
1010
1011
          "Component 7":
1012
          "Name": "Model Training",
1013
          "Description": "Train machine learning models on large datasets
1014
         using AWS SageMaker's Training Grounds feature",
          "Implementation difficulties": 8,
1015
          "Maintainability difficulties": 7 ,
1016
          "Component 8":
1017
         "Name": "Inference",
1018
         "Description": "Perform real-time inference on trained models,
1019
         making predictions on new, unseen data",
         "Implementation difficulties": 9,
1020
          "Maintainability difficulties": 8
1021
```

Figure 8: Example output from baseline Autogen system which scored the highest.