MEMORYFIELD: EXPLOITING GRAVITATIONAL FIELD FOR LONG-TERM MEMORY MANAGEMENT

Anonymous authors

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

031

033

034

037

038

040

041

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Despite the rapid progress of large language models (LLMs), which enables agents to perform complex decision-making and interaction, their limited long-term memory capacity hinders the effective retention and organization of historical interactions. This often leads to instability and semantic fragmentation in multi-turn dialogues and long-range reasoning tasks. Existing memory mechanisms struggle with structural reorganization, dynamic semantic retrieval, and the modeling of cognitive phenomena such as memory consolidation and forgetting. To address these challenges, we propose MemoryField, a novel dynamic spatial cognitive memory architecture driven by an attention-based gravitational field model. MemoryField represents memory items as nodes in a high-dimensional semantic space, where semantic attraction, repulsion, attention-driven forces, and decay mechanisms enable self-organized evolution and adaptive restructuring. By integrating node dynamics with fusion and forgetting processes, our approach ensures semantic coherence and cognitive stability. Extensive experiments demonstrate that MemoryField consistently outperforms existing memory mechanisms, improving dialogue quality by up to +4.9 Mauve and +3.3 ROUGE-L, boosting adversarial and temporal reasoning F1 by up to +14.7, and achieving superior performance across real-world tasks such as AlfWorld, ScienceWorld, HotPotQA, and FEVER, while maintaining strong cross-model generalization.

1 Introduction

The rapid advancement of artificial intelligence technologies has led to significant breakthroughs in large language models (LLMs) across natural language understanding, generation, and reasoning tasks (Vaswani et al., 2017; Chang et al., 2024). Consequently, LLM-based agents have emerged as a critical research focus in the field of AI (Guo et al., 2024; Xi et al., 2025). These agents possess autonomous decision-making and continuous interaction capabilities, enabling them to demonstrate substantial potential across a wide range of complex tasks (Cheng et al., 2024). In recent years, autonomous task agents such as AutoGPT (Yang et al., 2023) and BabyAGI (Nakajima, 2023), as well as reinforcement learning and knowledge-enhanced applications like Voyager (Wang et al., 2023a), Toolformer (Schick et al., 2023), and LangChain (Topsakal & Akinci, 2023), have showcased the powerful adaptability and task execution capabilities of LLM-driven agents in diverse environments.

Despite the strong performance of LLMs in short-term context modeling, their long-term memory capacity remains a critical limitation (Wang et al., 2023b). Specifically, LLM-based agents struggle to store and organize historical interaction data effectively and lack the ability to model long-term contextual continuity (Bulatov et al., 2022). This leads to instability, forgetting, and semantic discontinuities in multi-turn conversations, cross-task transfers, and long-term reasoning scenarios (Zhang et al., 2024). The absence of robust long-term memory mechanisms not only hinders the agent's ability to accumulate and reuse experience, but also limits its progression toward embodied intelligence or human-like cognitive capabilities (Wang et al., 2023a).

Currently, three primary approaches are being explored to address memory in LLMs: log-based memory, vector-based memory, and tool-augmented memory (Zhang et al., 2024). Log-based memory stores task histories or dialogue contents in chronological order, which is simple in structure but prone to redundancy and limited in revealing deep semantic relationships (Sordoni et al., 2015). Vector-based methods encode information into high-dimensional vectors and retrieve relevant content

based on similarity, enhancing relevance but lacking dynamic adjustment and semantic clustering capabilities—therefore struggling to support knowledge evolution and reasoning structure (Lewis et al., 2020). Tool-augmented memory relies on external knowledge bases or function calls to enhance functionality, but often neglects the optimization and self-evolution of internal memory structures (Nakano et al., 2021).

As a result, existing LLM memory mechanisms face three core challenges: the lack of structural reorganization capabilities, limited semantic retrieval efficiency, and an inability to effectively simulate key memory phenomena such as memory consolidation, conceptual fusion, and natural forgetting. To address these challenges, we propose a novel dynamic spatial cognitive memory architecture based on an attention-driven gravitational field model, and we are the first to manage memory in the form of a "force field." This framework constructs a quasi-physical interaction mechanism among memory nodes in high-dimensional semantic space, retaining the advantages of log-based memory accumulation while reconstructing the structure and access mechanisms of stored information. Each memory item is treated as a node in high-dimensional space, and four types of "forces" are designed—semantic attraction, repulsion, attentional center pull, and peripheral pushback. These forces guide the spatial reconfiguration and structural evolution of memory nodes based on semantic similarity, access frequency, and temporal decay.

Specifically, the attention gravitational field models memory state as a four-tuple (C_i, P_i, V_i, A_i) representing content, position, velocity, and activation level, respectively. A complete set of physical evolution rules is defined to allow memory nodes to dynamically adjust their spatial layout during interaction. The system also incorporates node fusion (for conceptual abstraction and redundancy reduction) and a forgetting mechanism (for pruning long-term low-activity memory), alongside energy-based convergence control to ensure stability and manageability of the evolving memory topology.

Across dialogue, long-context reasoning, and real-world benchmarks, our framework demonstrates consistent advantages over both naive and advanced memory baselines. It improves multi-turn dialogue coherence, enhances reasoning stability under extended contexts, and achieves competitive performance in interactive environments such as AlfWorld, ScienceWorld, HotPotQA, and FEVER. Moreover, ablation analyses highlight the necessity of our proposed gravitational forces for ensuring semantic cohesion and interpretability. Taken together, these findings underscore the effectiveness, robustness, and generality of MemoryField as a scalable long-term memory solution for LLM-based agents.

- We propose MemoryField, an attention-driven gravitational memory architecture that models memory as particles in a high-dimensional semantic space. By integrating semantic attraction/repulsion, attention pull, fusion, and forgetting, it supports dynamic self-organization, abstraction, and natural forgetting for scalable long-term memory management.
- We validate MemoryField through extensive experiments on multi-turn dialogue, long-context reasoning, and real-world tasks, showing significant improvements in coherence, reasoning stability, and cross-model generality over strong baselines.

2 RELATED WORK

2.1 MEMORY MECHANISMS IN LLM-BASED AGENTS

With the widespread application of large language models (LLMs) in dialogue, reasoning, and task planning, agents have demonstrated the ability to solve complex tasks through long-term interactions (Vaswani et al., 2017; Wei et al., 2022; Wang et al., 2024; Xi et al., 2025). Efficient information management has thus become a core challenge, driving research into memory mechanisms for intelligent agents (Sumers et al., 2023; Guo et al., 2024). Early approaches mainly relied on limited context windows, which are insufficient for long-term and complex tasks (Brown et al., 2020; Touvron et al., 2023). Recent studies have proposed scalable long-term memory mechanisms, including skill storage, knowledge base construction, and dynamic updating strategies, as seen in systems like Voyager, AppAgent, and MemPrompt (Madaan et al., 2022; Wang et al., 2023a; Zhang et al., 2023). In addition, hierarchical memory models improve retrieval efficiency through summarization and aggregation (Lewis et al., 2020; Jiang et al., 2023). However, current methods are

still limited in dynamic adjustment and forgetting strategies, often relying on static mechanisms that struggle to balance information retention and redundancy elimination (Madaan et al., 2022; Liu et al., 2024; Cheng et al., 2024). Therefore, developing more flexible and dynamic memory management approaches has become an important trend.

2.2 DYNAMIC KNOWLEDGE ORGANIZATION AND FORCE FIELD MODELING

Inspired by particle interactions in physics, force-directed modeling has been widely used in graph optimization and the self-organization of complex networks (Fruchterman & Reingold, 1991; Eades, 1984; Kamada et al., 1989). The four-force equilibrium model utilizes attraction and repulsion mechanisms to enable adaptive adjustment among nodes, improving structural rationality and dynamics (Newman, 2003; Leskovec et al., 2007). In artificial intelligence, existing knowledge graphs (e.g., TransE) are mostly static and struggle to handle relational evolution and new knowledge generation (Bordes et al., 2013; Wang et al., 2017). Although dynamic knowledge graphs introduce temporal embeddings, their flexibility remains limited (Trivedi et al., 2017; Xu et al., 2020). Furthermore, current forgetting mechanisms are mostly static and cannot simulate cognitive phenomena such as associative reinforcement and natural forgetting (Ebbinghaus, 2013; Cai et al., 2018). These limitations highlight the urgent need for a knowledge organization method capable of dynamic adjustment, flexible restructuring, and cognitive forgetting.

For a comprehensive review, please refer to Appendix A.1.

3 METHOD

To improve memory organization in LLM-powered agents for long-term interaction and complex reasoning, we propose a dynamic spatial cognitive architecture driven by an attentional gravitational field. Memory nodes are modeled as particles in a high-dimensional Euclidean space \mathbb{R}^n , each containing a semantic content vector, position, velocity, and activity level. Through four types of forces—inter-node repulsion and attraction, and attraction and repulsion relative to the origin—combined with query-driven dynamics and time decay, the system supports nonlinear memory structures, self-organizing knowledge topologies, and cognitive phenomena such as reinforcement, abstraction, and forgetting. Figure 1 illustrates the framework of our constructed memory field.

3.1 Model Architecture

In this system, each memory node is defined as $N_i = (C_i, P_i, V_i, A_i)$, where the meanings of each parameter are as follows:

 $C_i \in \mathbb{R}^d$: Semantic content vector. It exists in the d-dimensional real number space and is used to represent the semantic information of the memory node. For example, in the text memory scenario, through word vectors or sentence vectors, the semantics of the text are transformed into numerical vector representations. Different semantic contents will correspond to different vector values, enabling the similarity between semantics to be measured through vector calculations, as shown in Figure 2.

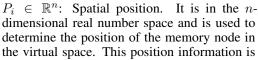




Figure 2: A user query is embedded and used to retrieve top-*k* relevant memory nodes from the vector memory space. The retrieved nodes, along with the current question, are input into the LLM to generate an answer. The answer and its context are then stored as a new memory node, initializing semantic embedding, spatial position, and activity level for subsequent dynamic updates.

crucial when simulating the interactions between nodes. For instance, the distance calculation between nodes depends on the position vectors, which in turn affect the attraction and repulsion forces between nodes.

 $V_i \in \mathbb{R}^n$: Velocity. Also, in the *n*-dimensional real number space, it describes the movement speed of the memory node in space. The change in velocity is determined by the net force acting on the

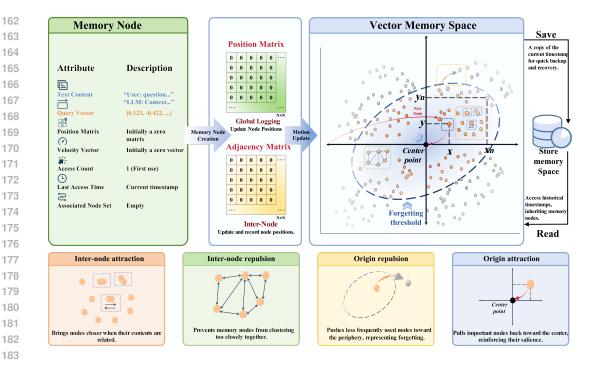


Figure 1: Overall workflow of the attentional memory system. Each memory node is modeled as a particle with semantic content, spatial position, velocity, and activity level. Node dynamics are governed by four forces: inter-node attraction and repulsion (based on semantic similarity and spatial proximity), and origin-based attraction and repulsion (driven by attention frequency and forgetting). The system maintains a self-organizing topology via position and adjacency matrices, supporting memory reinforcement, abstraction, and decay.

node and is closely related to the update of the position, reflecting the dynamic characteristics of the memory node in the system.

 A_i : Activity level. It represents the degree of activity of a memory node and is used to determine whether the node will be forgotten. The activity level changes dynamically over time and with usage: for example, each time the node is accessed, its activity level increases; if it remains unaccessed for a prolonged period, the activity level gradually decays. When the activity level of a node falls below a certain threshold, the node is marked as forgotten, thereby releasing storage space and maintaining the efficiency of the memory structure.

Let $W_{ij} = f_{\text{sim}}(C_i, C_j)$ denote the semantic similarity matrix, where f_{sim} is a function for calculating semantic similarity. It is calculated based on the semantic content vectors C_i and C_j of the nodes and reflects the degree of semantic association between two memory nodes. $D_{ij} = \|P_i - P_j\|$ is the spatial distance matrix, obtained by calculating the Euclidean norm of the position vectors of two nodes, and is used to measure the spatial distance between nodes. The net force on node i is defined as:

$$F_i = F_{i,\text{repel}} + F_{i,\text{attract}} + F_{i,\text{origin - repel}} + F_{i,\text{origin - attract}}$$
 (1)

This formula comprehensively considers four different types of forces, providing a comprehensive description of the force acting on the node in the system. The interaction of these forces determines the movement and state changes of the node. The following is an introduction to these four forces. **Inter-node Repulsion:**

$$F_{i,\text{repel}} = \sum_{j \neq i} \alpha \cdot \frac{P_i - P_j}{\|P_i - P_j\|^3} \tag{2}$$

The inter-node repulsion is designed to prevent memory nodes from over-aggregating in space. When two nodes are close, the repulsion force increases, causing them to move away from each other. In the formula, α is the repulsion coefficient, which controls the strength of the repulsion force. $\frac{P_i - P_j}{\|P_i - P_j\|^3}$

indicates that the direction of the repulsion force is from node j to node i, and the magnitude of the force is inversely proportional to the cube of the distance between the nodes. The closer the nodes are, the greater the repulsion force.

Inter-node Attraction:

$$F_{i,\text{attract}} = \sum_{j:W_{ij}>0} \beta_{ij} \cdot \frac{P_j - P_i}{\|P_j - P_i\|}$$
(3)

The inter-node attraction is used to connect semantically related nodes. Only when the semantic similarity matrix $W_{ij} > 0$, that is, when there is a certain semantic association between two nodes, will the attraction force be generated. β_{ij} is the attraction coefficient related to nodes i and j, which will be updated during operations such as associative queries. $\frac{P_j - P_i}{\|P_j - P_i\|}$ determines that the direction of the attraction force is from node i to node j, and the magnitude of the attraction force is related to the attraction coefficient and the distance between the nodes.

Repulsion from the Origin (Decay):

$$F_{i,\text{origin-repel}} = \gamma_i \cdot \frac{P_i}{\|P_i\|^3} \tag{4}$$

Parameter Updates. This repulsion simulates the natural decay process of memory. γ_i is the origin-repulsion coefficient related to node i, and $\frac{P_i}{\|P_i\|^3}$ indicates that the direction of the repulsion force is away from the origin, and the magnitude of the force is inversely proportional to the cube of the distance from the node to the origin. As the node moves away from the origin, the repulsion force gradually increases, meaning that the farther a node is from the origin, the more it is repelled, simulating the process by which memories that have not been accessed for a long time gradually weaken.

Attraction to the Origin (Attention Frequency):

$$F_{i,\text{origin - attract}} = \delta_i \cdot ||P_i|| \cdot \frac{-P_i}{||P_i||}$$
 (5)

This attraction reflects the attention frequency of the node. δ_i is the origin - attraction coefficient

Position Matrix

O a 332

Second Node

Second Node

Second Node

O a 332

O a 332

Node

O a 332

Second Node

O a 332

O a 332

O a 332

Node

O a 332

Node

O a 332

O a 33

Figure 3: Progressive construction of the position and adjacency matrices. As new memory nodes are added, semantic similarities and structural links are encoded to update global matrices, enabling spatial reasoning and interaction modeling.

related to node i, $||P_i||$ represents the distance from the node to the origin, and $\frac{-P_i}{||P_i||}$ determines that the direction of the attraction force is towards the origin. The closer a node is to the origin, the greater the attraction force it receives, indicating that nodes that are frequently accessed (with high activity levels) will be closer to the origin, reflecting the emphasis on frequently accessed memories. As shown in Figure 3.

During the operation of the system, some key parameters are updated based on different events.(1) $\delta_i(t+1) = \delta_i(t) + \Delta \delta_{\text{direct}}$ (Direct Query): During a direct query operation, if a certain node is queried, its corresponding origin - attraction coefficient δ_i will increase by $\Delta \delta_{\text{direct}}$. This indicates that the node's degree of attention has increased due to the query, and the attraction force to the origin has strengthened. (2) $\beta_{ij}(t+1) = \beta_{ij}(t) + \Delta \beta_{\text{assoc}}$ (Associative Query): When an associative query is performed and an association is found between nodes i and j, the attraction coefficient β_{ij} between them will increase by $\Delta \beta_{\text{assoc}}$, thereby strengthening the connection between semantically related nodes. (3) $\delta_i(t+1) = \delta_i(t) \cdot (1-\mu_\delta)$ (Time Decay): Over time, the origin-attraction coefficient δ_i will decay at a certain rate. μ_δ is the decay rate. This simulates the phenomenon that even nodes that were frequently accessed in the past will gradually decrease in attention as time passes. (4) $\beta_{ij}(t+1) = \beta_{ij}(t) \cdot (1-\mu_\delta)$: Similar to δ_i , the attraction coefficient β_{ij} also decays over time, with μ_β being its decay rate, reflecting the weakening process of the semantic association between nodes.

Association Probability (Establishment/Deletion). $p_{\text{build}} = \sigma(w_1 \cdot \sin(C_i, C_j) + w_2 \cdot (1 - D_{ij}/\theta_{\text{build}}))$: This is used to calculate the probability of establishing a new association. Here, σ is

the Sigmoid function, which maps the input value to the interval [0,1], making the result conform to the range of probabilities. w_1 and w_2 are weight parameters used to adjust the relative importance of the semantic similarity $\sin(C_i,C_j)$ and the spatial - distance - related term $(1-D_{ij}/\theta_{\text{build}})$ in the probability calculation. θ_{build} is the distance threshold for association establishment. When the distance D_{ij} between nodes is less than this threshold and the semantic similarity meets certain conditions, the probability of establishing an association will increase accordingly.

 $p_{\mathrm{drop}} = \sigma(v_1 \cdot (1 - \mathrm{sim}(C_i, C_j)) + v_2 \cdot \frac{D_{ij} - \theta_{\mathrm{build}}}{\theta_{\mathrm{drop}} - \theta_{\mathrm{build}}})$: This is used to calculate the probability of deleting an association. v_1 and v_2 are weight parameters, $(1 - \mathrm{sim}(C_i, C_j))$ represents the semantic dissimilarity, $\frac{D_{ij} - \theta_{\mathrm{build}}}{\theta_{\mathrm{drop}} - \theta_{\mathrm{build}}}$ is a distance - related term, and θ_{drop} is the distance threshold for association deletion. When the distance between nodes is greater than θ_{drop} or the semantic similarity is low, the probability of deleting the association will increase.

Fusion. (If $D_{ij} < \theta_{\text{fuse}}$ and $\text{sim}(C_i, C_j) > s_{\min}$) When multiple memory nodes meet the conditions that the distance is less than the fusion threshold θ_{fuse} and the semantic similarity is greater than the minimum similarity s_{\min} , a fusion operation will be performed:

$$C_f = f_{\text{fuse}}(C_1, \dots, C_k), \quad P_f = \frac{\sum w_i P_i}{\sum w_i}, \quad V_f = \frac{\sum w_i V_i}{\sum w_i}$$
 (6)

The semantic content vector C_f after fusion is calculated by the function $f_{\rm fuse}$, which comprehensively integrates the semantic information of each node participating in the fusion. The position vector P_f and velocity vector V_f are obtained by weighted averaging the corresponding vectors of the nodes participating in the fusion. The weights w_i can be set according to actual situations, and usually, $w_i=1$ is assumed for simple averaging. The fusion operation helps to reduce redundant memories and improve the organization and efficiency of memory.

Activity Decay and Forgetting. The activity decay formula is:

$$Activity(i,t) = Activity(i,t_0) \cdot \exp(-\lambda(t-t_0))$$
(7)

This indicates that the activity level of the memory node decays exponentially over time. λ is the decay coefficient, and $(t-t_0)$ is the time difference. As time increases, the activity level gradually decreases, reflecting the timeliness of memory. The forgetting judgment formula is:

$$\text{Forget}(i) = \begin{cases} \text{True,} & \text{if Activity}(i,t) < \theta_{\text{forget}} \text{ and } ||P_i|| > d_{\text{forget}} \\ \text{False,} & \text{otherwise} \end{cases} \tag{8}$$

When the activity level of a node is lower than the forgetting threshold θ_{forget} and the distance from the node to the origin is greater than the forgetting distance threshold d_{forget} , the node will be marked as a forgotten state. This mechanism ensures that memory nodes that have not been accessed for a long time and are far from the center of attention are properly processed, avoiding the occupation of excessive resources by invalid memories.

Position Update. The position and velocity of the node are updated based on the net force received: $V_i(t+\Delta t)=\beta\cdot V_i(t)+\alpha\cdot F_i(t)\cdot \Delta t$. This formula, based on the idea of Newton's second law, describes the update method of velocity. β is the velocity decay coefficient, used to simulate the natural decay of velocity during movement; α is the coefficient related to the force, which controls the influence of the net force on the change in velocity; $F_i(t)$ is the net force on node i at time t; Δt is the time step. $P_i(t+\Delta t)=P_i(t)+V_i(t+\Delta t)\cdot \Delta t$: The position of the node is updated according to the updated velocity, reflecting the cumulative effect of velocity on position change.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

Datasets. We evaluate MemoryField on diverse benchmarks spanning dialogue, reasoning, and real-world tasks. For dialogue, we use Multi-session Chat (MSC) and Conversation Chronicles (CC). For long-context reasoning, we construct five task categories—single-hop, multi-hop, temporal, open-domain, and adversarial—under context lengths from 4K to 16K. For real-world validation, we test on AlfWorld (sequential action execution), ScienceWorld (scientific reasoning), HotPotQA

(multi-hop QA), and FEVER (fact verification). Together, these benchmarks cover controlled settings and interactive environments.

Models. Our main experiments use GPT-3.5-turbo-16K, with and without MemoryField, and extend to GPT-4o, Claude Opus 4, LLaMA3.1-8B, Gemini 2.5 Flash, and Deepseek-R1. Baselines include All Dialogue History, All Memories + Context, Memory Retrieval, Rsum-LLM, MemoChat, COMEDY, and THEANINE, covering both naive and advanced memory mechanisms.

Metrics. For dialogue, we report BLEU-4, ROUGE-L, Mauve, and BERTScore. For reasoning, we measure F1 across the five task categories. For real-world tasks, we adopt official metrics: success rate (SR) for AlfWorld and HotPotQA/FEVER, and average reward (AR) for ScienceWorld. Cross-model evaluation follows automatic dialogue quality scoring.

Implementation. MemoryField is integrated as a structured *gravitational memory field*, where memory nodes evolve via attraction, repulsion, and decay forces. All methods use consistent prompts and fixed hyperparameters for fairness. Ablations disable individual forces to assess contributions. Repeated trials with fixed seeds ensure stable comparisons.

4.2 Performance Evaluation

Table 1: F1 scores of GPT-3.5-turbo-16K with and without MemoryField across context lengths.

Methods / Metrics	Multi-session Chat (MSC)				Conversation Chronicles (CC)			
Treations / Treatings	BLEU-4	ROUGE-L	Mauve	BERTScore	BLEU-4	ROUGE-L	Mauve	BERTScore
All Dialogue History	1.65	14.89	9.06	86.28	4.90	21.56	26.47	88.13
All Memories & Current Context \mathcal{D}	1.56	14.89	10.62	86.23	4.41	20.00	31.86	88.02
+ Memory Update	1.55	14.77	9.28	86.24	4.34	20.34	34.44	88.06
Memory Retrieval	1.92	15.49	11.16	86.20	4.40	20.48	33.24	88.09
+ Memory Update	1.67	15.30	13.71	86.44	4.36	20.33	34.84	88.02
Rsum-LLM	0.75	11.53	2.45	84.61	0.98	11.42	2.28	85.59
MemoChat	1.42	13.11	7.72	85.94	2.31	15.87	15.12	87.03
COMEDY	1.06	12.79	7.27	85.29	1.70	13.57	19.55	85.90
THEANINE	1.80	15.37	18.62	86.70	6.58	22.68	64.41	88.58
MemoryField(Ours)	1.87	16.10	23.50	86.79	6.82	23.44	64.73	89.10

Multi-turn Dialogue Evaluation. As shown in Table 1, MemoryField achieves either the best or highly competitive overall performance across both dialogue datasets. On the MSC dataset, MemoryField reaches a Mauve score of 23.50, outperforming the best-performing baseline (THEANINE, 18.62) by 4.88 points. It also improves the ROUGE-L score to 16.10, representing a gain of approximately 3.3 points over COMEDY (12.79). In addition, MemoryField slightly surpasses other methods in both BLEU-4 and BERTScore. On the CC dataset, MemoryField yields modest improvements in BLEU-4 (6.82 vs. THEANINE's 6.58) and ROUGE-L (23.44 vs. 22.68), while maintaining a lead in Mauve (64.73 vs. 64.41). Notably, it achieves the highest BERTScore of 89.10, indicating superior semantic consistency. Compared to All Dialogue History and Memory Retrieval-based methods, MemoryField delivers an average improvement of more than 12 points in Mauve, and gains of 1.0–3.5 points in BLEU-4 and ROUGE-L. These results demonstrate the effectiveness of our structured gravitational memory field in enhancing semantic focus, reinforcing relevant information, and suppressing redundancy. Overall, MemoryField exhibits strong context preservation and improved generation quality, enabling more semantically coherent and consistent responses in multi-turn open-domain dialogue. These findings validate the model's memory advantages in long-range interactive scenarios.

Table 2: F1 scores of GPT-3.5-turbo-16K with and without MemoryField across context lengths (4K–16K). Abbreviations: **S.H.** = Single Hop, **M.H.** = Multi Hop, **Temp.** = Temporal, **O.D.** = Open Domain, **Adv.** = Adversarial.

Model	Ctx.	S.H.	M.H.	Temp.	O.D.	Adv.	Overall
	4K	31.7	25.4	16.8	27.6	13.1	24.1
GPT-3.5-turbo-16K	8K	38.8	31.2	21.0	35.0	8.4	25.2
	12K	51.1	40.4	25.0	36.5	6.4	33.5
	16K	56.4	42.0	20.3	37.2	2.1	37.8
+MemoryField	4K 8K	33.4 40.3	27.8 34.5	23.3 28.3	34.4 39.7	24.7 19.6	26.8 28.6
(Ours)	12K 16K	54.2 58.2	42.9 44.7	41.5 35.0	40.2 41.9	17.2 10.6	35.6 39.1

Long-context Reasoning Evaluation. To further assess the effectiveness of MemoryField, we evaluate GPT-3.5-turbo-16K with and without MemoryField across various context lengths (4K to 16K) and five reasoning tasks. As summarized in Table 2, MemoryField consistently enhances model performance across all settings. Without memory augmentation, the model's F1 score improves with longer contexts (from 24.1 at 4K to 37.8 at 16K); however, it exhibits instability on complex tasks. Notably, in adversarial reasoning, the F1 score plummets from 13.1 to 2.1 at 16K, suggesting that extended contexts can introduce

detrimental noise that impairs reasoning. In contrast, the MemoryField-enhanced model demonstrates improved stability and scalability. At 16K, it yields relative F1 improvements of 1.8 (single-hop), 2.7 (multi-hop), 14.7 (temporal), and 8.5 (adversarial), with an overall gain of 1.3. The gains are particularly substantial for temporal and adversarial tasks, highlighting MemoryField's effectiveness in handling long-range dependencies and semantic noise.

We attribute this improvement to MemoryField's mechanism of modeling past information as structured semantic entities, which are dynamically integrated via a gravitational attention mechanism. This mechanism amplifies relevant signals while suppressing irrelevant ones, enabling more robust and coherent reasoning paths across long contexts.

Cross-model evaluation. To evaluate the performance of different memory mechanisms in multi-turn dialogue scenarios and to verify whether the proposed MemoryField can maintain consistent advantages across various mainstream large models (GPT-40, Claude Opus 4, LLaMA3.1-8B, Gemini 2.5 Flash, Deepseek-R1), we conduct experiments on two standard benchmark datasets: Multisession Chat (MSC) and Conversation Chronicles (CC). Both tasks involve long-range contexts and complex conversational dynamics. The baselines include history replay, memory-based retrieval mechanisms, summarization-based methods (Rsum-LLM), as well as existing memory-augmented models (MemoChat, COMEDY, THEANINE). The evaluation metric is automatic scoring, where higher values indicate better dialogue generation quality.

As shown in Table 3, MemoryField achieves the best or highly competitive results across all models and both datasets. On the MSC dataset, MemoryField typically shows improvements of 0.2–1.0 over the best baseline. On the CC dataset, its advantage is even more pronounced, with average improvements of 1.0–2.5 compared to THEANINE and other methods. For instance, with GPT-40, MemoryField reaches a score of 27.35 on CC, significantly surpassing THEANINE's 25.05; similar consistent gains are observed with Claude Opus 4 and LLaMA3.1-8B.

Table 3: Performance comparison across memory methods on MSC and CC tasks. Abbreviations of methods: **Hist.** = All Dialogue History, **Mem.+Ctx.** = All Memories & Context, **Retr.** = Memory Retrieval, **Rsum** = Rsum-LLM, **Memo** = MemoChat, **COM.** = COMEDY, **THEA.** = THEANINE, **MemField** = MemoryField (Ours).

Model Name	Task	Hist.	Mem.+Ctx.	Retr.	Rsum	Memo	COM.	THEA.	MemField
GPT-40	MSC	18.25	18.32	17.80	14.30	15.10	14.15	16.90	18.72
GPT-40	CC	24.15	23.75	23.40	14.90	18.10	16.25	25.05	27.35
Claude Opus 4	MSC	17.90	18.12	17.30	13.80	14.65	14.90	16.10	17.90
Claude Opus 4	CC	23.85	22.45	23.10	14.45	17.25	15.95	24.90	26.10
LLaMA3.1 8B	MSC	16.90	17.50	18.00	13.10	14.20	14.45	15.70	17.50
LLaMA3.1 8B	CC	23.05	21.75	22.15	13.40	16.20	15.00	23.50	25.05
Gemini 2.5 flash	MSC	17.15	17.80	17.20	13.50	14.35	14.05	15.60	17.40
Gemini 2.5 flash	CC	23.35	22.25	22.80	14.55	17.00	16.00	24.10	25.85
Deepseek-R1	MSC	17.25	17.65	17.10	13.80	14.50	14.30	15.90	17.60
Deepseek-R1	CC	23.25	22.05	22.45	14.80	16.75	15.60	24.25	25.60

Traditional summarization-based methods (e.g., Rsum-LLM) and some earlier memory models (e.g., COM-EDY, MemoChat) perform relatively poorly in long dialogue settings, struggling to capture global context. While THEANINE demonstrates competitive performance in certain cases, it still falls short of MemoryField. Importantly, MemoryField delivers stable improvements across diverse model architectures, indicating that its memory mechanism possesses strong generality and transferability.

4.3 ABLATION STUDY

As illustrated in Figure 4, we visualize the spatial distribution of memory nodes under different force configurations. When all four forces—node attraction, node repulsion, origin attraction, and origin repulsion—are enabled (Figure 4a), the nodes form a well-structured and coherent layout around the central query point (red star). Node attraction clusters semantically related items, node repulsion prevents overlap, origin attraction pulls important nodes toward the center, and origin repulsion ensures dispersion. Their synergy yields semantically cohesive and spatially interpretable memory organization.

In contrast, disabling all forces (Figure 4b) produces a random distribution, where nodes scatter without clear semantic clustering and some drift far from the query. This highlights the necessity of the gravitational field mechanism for generating meaningful and interpretable memory structures.

To further analyze the role of each force, we conduct ablation experiments (Figures 4c–e). Removing node attraction disrupts semantic clustering, yielding more uniform but less coherent layouts, showing its importance for encoding semantic similarity. Disabling node repulsion collapses nodes into dense clusters, confirming its role in maintaining separation and preventing crowding. Without origin

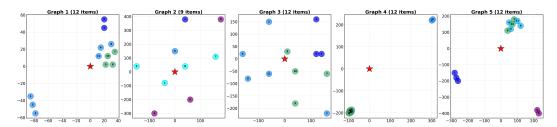


Figure 4: Visualization of memory node configurations: (a) baseline, (b) w/o node attraction, (c) w/o node repulsion, (d) w/o origin attraction, and (e) w/o full force mechanism.

attraction, local clusters still form, but the global structure drifts away from the query point, indicating its importance for contextual alignment. Collectively, these results demonstrate that each force contributes uniquely to the memory topology, and their combination is essential for achieving a balanced, interpretable, and effective memory organization.

4.4 REAL-WORLD TASK EVALUATION

To assess practical effectiveness, we evaluate MemoryField on several real-world benchmarks, targeting environments that require long-horizon reasoning, interactive decision-making, and evidence verification, and ask whether it delivers consistent performance gains. We consider four representative real-world tasks: (1) AlfWorld, a household environment requiring sequential action execution, evaluated with success rate (SR); (2) ScienceWorld, a scientific experiment environment requiring reasoning and multi-step tool usage, evaluated with average reward (AR); (3) HotPotQA, a multi-hop question answering benchmark measuring reasoning accuracy (SR); (4) FEVER, a fact verification task evaluating evidence retrieval and logical consistency (SR). We compare MemoryField against multiple baselines, including zero-shot reasoning (Z-CoT, F-CoT, CoT-SC), interactive decision-making approaches (SayCan, ReAct), and reflection-based reasoning (Reflexion).

Table 4: Performance comparison of reasoning and memoryaugmented methods across multiple real-world benchmarks. Metrics: **SR** = Success Rate, **AR** = Average Reward. "-" means not reported.

Method	AlfWorld (SR%)	ScienceWorld (AR)	HotPotQA (SR%)	FEVER (SR%)
Z-CoT	-	-	0.01	0.39
F-CoT	0.43	16.58	0.32	0.61
CoT-SC	0.57	15.24	0.33	0.62
SayCan	0.60	12.36	-	-
ReAct	0.57	15.05	0.34	0.63
Reflexion	0.71	19.39	0.39	0.68
MemoryField (Ours)	0.75	20.42	0.41	0.71

Results and Analysis. As in Table 4, MemoryField achieves the best or highly competitive results across all benchmarks. In AlfWorld, MemoryField attains a success rate of 0.75, outperforming Reflexion (0.71) and demonstrating stronger robustness in long-horizon action planning. In ScienceWorld, MemoryField achieves an average reward of 20.42, exceeding Reflex-

ion (19.39), which highlights its advantage in scientific reasoning and tool usage. In HotPotQA, MemoryField obtains 0.41, surpassing all baselines and showing its ability to maintain consistency in multi-hop reasoning. In FEVER, MemoryField reaches 0.71, higher than Reflexion (0.68), confirming its benefit in fact verification tasks. These results verify that MemoryField consistently improves agent performance in diverse real-world scenarios, demonstrating its strong generalization ability and robustness under challenging interactive and reasoning tasks.

5 Conclusion

In this paper, we propose MemoryField, a novel attention-driven gravitational memory architecture designed to address the challenges of long-term memory management in LLM-based agents. By modeling memory nodes as particles in a high-dimensional semantic space and simulating their dynamic evolution through force-directed interactions (semantic attraction, repulsion, attention-centric pull, and decay), we achieve structured memory self-organization, conceptual abstraction, and natural forgetting. Extensive experiments on multi-turn dialogue and long-context reasoning benchmarks demonstrate that, compared with traditional vector-based and graph-augmented memory methods, MemoryField significantly improves semantic coherence, information retention, and reasoning consistency. These results validate the potential of MemoryField in long-term interaction and adaptive knowledge management.

REPRODUCIBILITY STATEMENT

We are committed to the full reproducibility of this work. The proposed MemoryField architecture, including the gravitational force—driven memory dynamics and update rules, is described in detail with pseudocode in the Appendix, ensuring that future researchers can directly reproduce and extend our study. Our experimental setup is comprehensively introduced in Section 4. Details of hyperparameter choices, ablation configurations, and heuristic tuning are provided in Appendix A.4 and A.5. Algorithm pseudocode is included in Appendix A.3, while additional experimental results, ablations, and visualizations are presented in Appendix A.5. All experiments are implemented in a Python environment. Upon publication, we will release the complete source code, configuration files, and training examples, enabling other researchers to directly verify and further advance this line of work

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Richard C Atkinson and Richard M Shiffrin. Human memory: A proposed system and its control processes. In *Psychology of learning and motivation*, volume 2, pp. 89–195. Elsevier, 1968.
- Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems*, 26, 2013.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. Recurrent memory transformer. *Advances in Neural Information Processing Systems*, 35:11079–11091, 2022.
- Hongyun Cai, Vincent W Zheng, and Kevin Chen-Chuan Chang. A comprehensive survey of graph embedding: Problems, techniques, and applications. *IEEE transactions on knowledge and data engineering*, 30(9):1616–1637, 2018.
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM transactions on intelligent systems and technology*, 15(3):1–45, 2024.
- Yuheng Cheng, Ceyao Zhang, Zhengwen Zhang, Xiangrui Meng, Sirui Hong, Wenhao Li, Zihao Wang, Zekai Wang, Feng Yin, Junhua Zhao, et al. Exploring large language model based intelligent agents: Definitions, methods, and prospects. *arXiv preprint arXiv:2401.03428*, 2024.
- Peter Eades. A heuristic for graph drawing. Congressus numerantium, 42(11):149–160, 1984.
- Hermann Ebbinghaus. [image] memory: A contribution to experimental psychology. *Annals of neurosciences*, 20(4):155, 2013.
- Thomas MJ Fruchterman and Edward M Reingold. Graph drawing by force-directed placement. *Software: Practice and experience*, 21(11):1129–1164, 1991.
- Rishab Goel, Seyed Mehran Kazemi, Marcus Brubaker, and Pascal Poupart. Diachronic embedding for temporal knowledge graph completion. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 3988–3995, 2020.
- 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.
- Jacob N Israelachvili. Intermolecular and surface forces. Academic press, 2011.

- Zhengbao Jiang, Frank F Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang,
 Jamie Callan, and Graham Neubig. Active retrieval augmented generation. In *Proceedings of the* 2023 Conference on Empirical Methods in Natural Language Processing, pp. 7969–7992, 2023.
- Tomihisa Kamada, Satoru Kawai, et al. An algorithm for drawing general undirected graphs. *Information processing letters*, 31(1):7–15, 1989.
 - 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.
 - Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM transactions on Knowledge Discovery from Data (TKDD)*, 1(1):2–es, 2007.
 - 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.
 - Junwei Liu, Kaixin Wang, Yixuan Chen, Xin Peng, Zhenpeng Chen, Lingming Zhang, and Yiling Lou. Large language model-based agents for software engineering: A survey. *arXiv preprint arXiv*:2409.02977, 2024.
 - Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, et al. Agentbench: Evaluating llms as agents. *arXiv preprint arXiv:2308.03688*, 2023.
 - Aman Madaan, Niket Tandon, Peter Clark, and Yiming Yang. Memprompt: Memory-assisted prompt editing with user feedback, 2022.
 - Yohei Nakajima. BabyAGI, 2023. URL https://github.com/yoheinakajima/babyagi.
 - Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv* preprint arXiv:2112.09332, 2021.
 - Mark EJ Newman. The structure and function of complex networks. *SIAM review*, 45(2):167–256, 2003.
 - Maximilian Nickel, Kevin Murphy, Volker Tresp, and Evgeniy Gabrilovich. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE*, 104(1):11–33, 2015.
 - Andreas Noack. Modularity clustering is force-directed layout. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 79(2):026102, 2009.
 - Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551, 2023.
 - Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. A neural network approach to context-sensitive generation of conversational responses. *arXiv preprint arXiv:1506.06714*, 2015.
 - Theodore Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas Griffiths. Cognitive architectures for language agents. *Transactions on Machine Learning Research*, 2023.
 - Mariya Toneva, Alessandro Sordoni, Remi Tachet des Combes, Adam Trischler, Yoshua Bengio, and Geoffrey J Gordon. An empirical study of example forgetting during deep neural network learning. *arXiv preprint arXiv:1812.05159*, 2018.
 - Oguzhan Topsakal and Tahir Cetin Akinci. Creating large language model applications utilizing langchain: A primer on developing llm apps fast. In *International Conference on Applied Engineering and Natural Sciences*, volume 1, pp. 1050–1056, 2023.

- 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.
- Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-evolve: Deep temporal reasoning for dynamic knowledge graphs. In *international conference on machine learning*, pp. 3462–3471. PMLR, 2017.
- 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.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv* preprint arXiv:2305.16291, 2023a.
- Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.
- Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. Knowledge graph embedding: A survey of approaches and applications. *IEEE transactions on knowledge and data engineering*, 29(12): 2724–2743, 2017.
- Weizhi Wang, Li Dong, Hao Cheng, Xiaodong Liu, Xifeng Yan, Jianfeng Gao, and Furu Wei. Augmenting language models with long-term memory. *Advances in Neural Information Processing Systems*, 36:74530–74543, 2023b.
- 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.
- Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *Science China Information Sciences*, 68(2):121101, 2025.
- Chengjin Xu, Mojtaba Nayyeri, Fouad Alkhoury, Hamed Shariat Yazdi, and Jens Lehmann. Tero: A time-aware knowledge graph embedding via temporal rotation. arXiv preprint arXiv:2010.01029, 2020.
- Hui Yang, Sifu Yue, and Yunzhong He. Auto-gpt for online decision making: Benchmarks and additional opinions. *arXiv preprint arXiv:2306.02224*, 2023.
- Chi Zhang, Zhao Yang, Jiaxuan Liu, Yucheng Han, Xin Chen, Zebiao Huang, Bin Fu, and Gang Yu. Appagent: Multimodal agents as smartphone users. *arXiv preprint arXiv:2312.13771*, 2023.
- 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. *arXiv* preprint arXiv:2404.13501, 2024.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv* preprint arXiv:2303.18223, 1(2), 2023.
- Wanjun Zhong, Lianghong Guo, Qiqi Gao, He Ye, and Yanlin Wang. Memorybank: Enhancing large language models with long-term memory. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 19724–19731, 2024.

CONTENTS

Introduction **Related Work** Dynamic Knowledge Organization and Force Field Modeling Method 3.1 **Experiments** 4.2 4.3 4.4 Conclusion A Appendix Dynamic Knowledge Organization and Force Field Modeling A.3

A APPENDIX

A.1 DETAILED RELATED WORK

A.1.1 MEMORY MECHANISMS IN LLM-BASED AGENTS

With the rapid development of artificial intelligence technologies, large language models (LLMs) have demonstrated significant potential in areas such as dialogue systems, automated reasoning Wei et al. (2022), and task planning Vaswani et al. (2017); Zhao et al. (2023); Brown et al. (2020). LLM-based agents gradually learn and optimize their decision-making capabilities through long-term interactions, enabling them to tackle complex tasks Xi et al. (2025); Wang et al. (2024); Liu et al. (2023). For example, agents can adjust dialogue strategies based on user feedback or infer optimal plans during task execution Achiam et al. (2023); Kojima et al. (2022). However, such long-term interactions generate vast amounts of data, making efficient information management a critical challenge Sumers et al. (2023); Wang et al. (2024).

Memory management has thus emerged as a core mechanism of intelligent agents, responsible for storing and updating interaction experiences, as well as retrieving relevant information based on task requirements. For instance, an agent may record user preferences or task states to enhance decision-making efficiency Xi et al. (2025); Sumers et al. (2023); Guo et al. (2024). Memory management not only supports task execution but also improves the agent's decision-making capabilities by analyzing historical experiences. This adaptability in dynamic environments lays the groundwork for the pursuit of artificial general intelligence.

Early studies mainly relied on simple context windows to manage short-term memory, which sufficed for low-complexity tasks Brown et al. (2020). However, as task complexity increases and the duration of human-agent interaction extends, short-term memory reveals limitations in capacity and its ability to maintain contextual continuity Touvron et al. (2023). These limitations have prompted researchers to explore more scalable and adaptive long-term memory mechanisms Zhong et al. (2024).

To meet the demands of diverse tasks, long-term memory must not only support the storage and retrieval of information but also possess the capability for dynamic adaptation and updates. Current research primarily focuses on skill storage, knowledge base construction, and memory updating strategies. For example, *Voyager* stores executable code in a skill repository and dynamically updates it based on environmental feedback to enable skill transfer and reuse Wang et al. (2023a); *AppAgent* builds a structured knowledge base through autonomous exploration and human demonstrations to support complex tasks Zhang et al. (2023); *MemPrompt* records user feedback to generate memory entries that enhance future responses Madaan et al. (2022). In addition, some studies draw inspiration from multi-level caching in operating systems, proposing hierarchical memory models that employ summarization or information aggregation to improve retrieval efficiency Lewis et al. (2020); Jiang et al. (2023).

Despite recent advances, long-term memory management still faces several challenges. Specifically, the growing volume of interaction data increases storage and retrieval costs, hindering the scalability of agents in large-scale tasks Liu et al. (2024); Cheng et al. (2024). Existing forgetting and updating strategies are often static (e.g., time-decay-based deletion) and lack the ability to dynamically retain or discard information based on task context. This may lead to the loss of critical information or the accumulation of redundant data, thereby reducing overall efficiency Madaan et al. (2022).

To address the challenge of dynamic adjustment, we draw inspiration from physics-based force-directed principles, laying the groundwork for our subsequent exploration of force-guided models in memory management.

A.1.2 Dynamic Knowledge Organization and Force Field Modeling

Force-oriented modeling inspired by particle interactions in physics is a powerful method for dynamic evolution analysis. It has been widely applied in graph structure optimization, particle system simulation, and the visualization and organization of complex networks Fruchterman & Reingold (1991); Eades (1984); Kamada et al. (1989). By constructing a dynamic model based on the balance of four types of forces, this method simulates interaction forces between nodes to achieve adaptive system adjustment. Specifically, the attractive force between nodes promotes the connection of related

nodes, enhancing structural cohesion; the repulsive force between nodes prevents excessive clustering and maintains distribution uniformity; the attraction from nodes to the origin reflects external attention or activation frequency, guiding important nodes toward the center; and the repulsion from nodes to the origin simulates information decay or natural diffusion, pushing nodes away from the center to avoid information overload Leskovec et al. (2007); Israelachvili (2011). The synergy of these forces drives the system toward an energy-minimized equilibrium state during evolution, forming a structurally reasonable and dynamically adjustable distribution pattern Newman (2003); Noack (2009). For example, in social network analysis, such a four-force equilibrium model can reveal potential relationships between nodes and optimize network layouts; in molecular dynamics simulation, it can simulate particle interactions to predict stable configurations. The flexibility and generality of this model provide a solid theoretical foundation for the adaptive reorganization of complex information networks, opening new perspectives for interdisciplinary research, such as knowledge organization in artificial intelligence.

In the field of artificial intelligence, the design of long-term memory and knowledge organization systems aims to support information storage, retrieval, and reasoning in complex tasks. However, existing methods still face significant challenges in dynamic environments. Static knowledge graphs (such as TransE) represent knowledge using fixed triples (entity-relation-entity), which are suitable for reasoning in static scenarios but struggle to adapt to relational changes and the generation of new relations in dynamic tasks, leading to a decrease in prediction accuracy Bordes et al. (2013); Nickel et al. (2015); Wang et al. (2017). For example, in real-time recommendation systems, static knowledge graphs cannot quickly capture the dynamic evolution of user interests, limiting their effectiveness. Dynamic knowledge graphs attempt to capture the temporal evolution of knowledge by introducing time embeddings, but due to their reliance on predefined relation templates, they struggle to enable free restructuring of knowledge, limiting their adaptability in open-domain tasks Trivedi et al. (2017); Goel et al. (2020); Xu et al. (2020). For instance, when handling emergent events (such as news events), existing dynamic knowledge graphs often fail to flexibly update relational networks due to template constraints. Furthermore, current methods fall short in modeling the forgetting mechanism within cognitive processes. Static forgetting strategies (such as fixed time decay) cannot accurately simulate cognitive phenomena such as associative reinforcement, abstract integration, and natural forgetting Ebbinghaus (2013); Atkinson & Shiffrin (1968), leading to the erroneous elimination of critical information or the prolonged retention of redundant data, thereby reducing system efficiency and intelligence Cai et al. (2018); Toneva et al. (2018). These limitations suggest that current knowledge organization systems are in urgent need of a dynamic method capable of adaptively adjusting structure, flexibly restructuring relations, and simulating cognitive forgetting.

In view of the limitations in dynamic organization, flexible restructuring, and cognitive forgetting modeling in existing methods, this paper, inspired by the four-force equilibrium modeling in physics, proposes an attention-driven spatial memory mechanism. The specific methodology will be introduced in detail in the next section.

A.2 NOTATION

Symbol	Definition	Meaning			
$ \begin{array}{c} N_i = \\ (C_i, P_i, V_i, A_i) \end{array} $	Memory node	A memory unit including semantics, position velocity, and activity			
$C_i \in \mathbb{R}^d$	Semantic content vector	Vector representation of semantic information (e.g., text embeddings)			
$P_i \in \mathbb{R}^n$	Spatial position vector	Node's coordinates in high-dimensional space for force calculation			
$V_i \in \mathbb{R}^n$	Velocity vector	Describes the node's motion in space			
$A_i \in \mathbb{R}$	Activity level	Represents access frequency or memory strength			
W_{ij}	Semantic similarity	Degree of semantic association between nodes i and j			
$D_{ij} = \ P_i - P_j\ $	Euclidean distance	Spatial distance between two memory nodes			
F_{i}	Net force	Total force acting on node i			
$F_{i,\mathrm{repel}}$	Inter-node repulsion	Prevents nodes from over-aggregating			
$F_{i, \mathrm{attract}}$	Inter-node attraction	Attracts semantically related nodes			
$F_{i, \mathrm{origin}}$ - repel	Repulsion from origin	Simulates natural memory decay			
$F_{i, m origin}$ - attract	Attraction to origin	Simulates attention-based memory reinforce ment			
α	Repulsion coefficient	Controls strength of repulsion between nodes			
eta_{ij}	Attraction coefficient	Controls strength of attraction between nodes i and j			
γ_i	Origin repulsion coefficient	Governs tendency of node to drift away fron origin			
δ_i	Origin attraction coefficient	Governs tendency of node to be pulled towar origin			
λ	Activity decay rate	Controls exponential decay of node activity over time			
$ heta_{ ext{query}}$	Query threshold	Similarity threshold for returning a query resul			
$ heta_{ m activate}$	Activation threshold	Minimum value to activate related nodes in associative query			
$ heta_{ m link}$	Link formation threshold	Controls whether a semantic link is established			
$ heta_{ ext{fuse}}$	Fusion distance threshold	Max distance for node fusion to occur			
s_{\min}	Minimum similarity for fusion	Required semantic similarity for merging node			
$d_{ m forget}$	Forgetting distance threshold	Minimum distance for a low-activity node to be forgotten			
$ heta_{ ext{forget}}$	Forgetting activity threshold	Activity level below which nodes may be discarded			
ϵ	Energy threshold	System is stable if energy falls below this val			
δ	Energy change threshold	Determines system convergence by energy difference			
E(t)	System energy	Sum of squared net forces across all nodes			
$sim(C_i, C_j)$	Similarity function	Measures semantic similarity, e.g., cosine similarity			
$\sigma(x)$	Sigmoid function	Maps values to range $[0,1]$ to represent probabilities			
$f_{ m fuse}$	Fusion function	Aggregates semantic vectors from multiple nodes			

Table 5: Mathematical Symbols and Their Meanings

A.3 PSEUDOCODE

```
865
866
           Algorithm 1 Attentional Gravitational Field Architecture
867
           Input: Memory nodes \{N_i = (C_i, P_i, V_i, A_i)\}, query q
868
           Output: Query result r
           // Direct Query
870
         1 foreach node i do
            s_i \leftarrow \text{cosine\_similarity}(C_i, q) - \lambda ||P_i||
871
872
        j \leftarrow \arg \max_i s_i \text{ if } s_i > \theta_{query} \text{ then }
873
         4 | \delta_i += \Delta \delta_{\text{direct}} \ r \leftarrow N_i \ \text{QueryCount} + +
874
        5 else
875
                r \leftarrow \text{None}
876
           // Associative Query
877
        7 foreach active i, Depth < MaxDepth, total < MaxNodes do
878
                foreach j with W_{ij} > 0 do
879
                     p_j \leftarrow W_{ij}(1 - D_{ij}/\theta_{\text{max}}) if p_j > \theta_{activate} then
880
                          Activate j \beta_{ij}, \beta_{ji} += \Delta \beta_{assoc} QueryCount+ +
        10
882
       11 if QueryCount \geq N_{update} then
883
                while not converged do
        12
                     E \leftarrow 0 foreach node i do
        13
885
                          // Compute force-based updates
                          F_i \leftarrow F_{\text{repel}} + F_{\text{attract}} + F_{\text{origin-repel}} + F_{\text{origin-attract}} \quad V_i \leftarrow \beta V_i + \alpha F_i \Delta t \quad P_i \leftarrow P_i + V_i \Delta t
        14
887
                            E += ||F_i||^2
888
                     // Update Links
889
                     foreach pair (i, j) do
        15
890
                          p_{\text{est}} \leftarrow \sigma(w_1 \text{cosine\_similarity}(C_i, C_j) + w_2(1 - D_{ij}/\theta_{\text{establish}})) if p_{est} > \theta_{link}
        16
891
892
                             W_{ij}, W_{ji} \leftarrow \text{cosine\_similarity}(C_i, C_i)
        17
893
                     // Fuse & Forget
894
                     foreach pair (i, j) do
        18
895
                          if D_{ij} < \theta_{fuse} and cosine_similarity(C_i, C_j) > s_{\min} then
        19
896
                            N_f \leftarrow \text{fuse}(N_i, N_i) Replace N_i, N_i with N_f
        20
897
                     foreach i do
        21
                          A_i \leftarrow \text{decay}(A_i) if forget (i) with Activity(i) < \theta_{forget} and ||P_i|| > d_{forget} then
900
                           Mark N_i as forgotten
        23
901
                          \delta_i \leftarrow \delta_i (1 - \mu_\delta) foreach neighbor j do
902
                           \beta_{ij} \leftarrow \beta_{ij} (1 - \mu_{\beta})
        25
903
                     if E < \epsilon or |E - E_{prev}| < \delta then
       26
904
                      break
        27
905
                     E_{\text{prev}} \leftarrow E
906
       28
907
                Reset QueryCount to 0
908
       30 return r
909
```

Overall Algorithm. Algorithm 1 presents the pseudocode of our Attentional Gravitational Field Architecture. The process begins with a *direct query*, where each memory node is scored by the similarity between its content and the query, adjusted by spatial distance. If the best-matched node surpasses the query threshold, it is retrieved and its origin-attraction coefficient is reinforced. If direct retrieval fails, the system performs an *associative query* by expanding to neighbors with strong semantic or structural links, thereby activating additional relevant nodes. Once the number of queries exceeds a preset threshold, the system updates memory dynamics through iterative force-based evolution: all four forces (repulsion, attraction, origin-repulsion, origin-attraction) are applied to update velocity and position, while the global energy is accumulated to monitor convergence. During

this process, links are adaptively established or removed, and redundant nodes are merged through fusion. Simultaneously, activity levels decay over time, and nodes with low activity and peripheral positions are forgotten. The loop terminates when the energy drops below a predefined value or stabilizes, after which the system resets and returns the final query result. This design ensures that memory retrieval, update, fusion, and forgetting are integrated into a unified dynamic framework.

Memory Snapshot. To improve the flexibility of memory management, the system introduces a snapshot functionality. At the conclusion of each dialogue session, the system performs a snapshot operation to preserve the state of the memory repository. This operation captures comprehensive details of all memory nodes, including their semantic content vectors C_i , spatial positions P_i , velocities V_i , activity levels A_i , and inter-node associations, which are represented by the semantic similarity matrix W_{ij} and the spatial distance matrix D_{ij} . Additionally, the current values of key parameters, such as α , β_{ij} , γ_i , and δ_i , are recorded.

When a user seeks to resume a previous interaction, they can select the corresponding snapshot file, enabling the system to swiftly restore the memory repository to its saved state. Upon restoration, the system leverages the existing dynamic spatial cognitive architecture to continue memory node updates, association adjustments, fusion operations, and forgetting evaluations based on new query demands. This ensures seamless continuity and dynamic evolution of the memory repository.

A.4 HYPERPARAMETERS AND HEURISTIC TUNING

Grouping and Roles. For clarity and reproducibility, we categorize the hyperparameters into four groups. The first group consists of force coefficients: inter-node repulsion α , semantic attraction β_{ij} , origin repulsion γ_i , and origin attraction δ_i . These directly determine the four forces in Eqs. (2)–(5), shaping convergence patterns and global sparsity. The second group contains change-rate parameters: direct-query gain $\Delta \delta_{\text{direct}}$, associative-query gain $\Delta \beta_{\text{assoc}}$, and temporal decay rates μ_{δ} , μ_{β} . These parameters control the amplification of links and attention pulls triggered by queries, as well as gradual fading over time. The third group consists of structural thresholds: link creation and deletion thresholds θ_{build} , θ_{drop} , θ_{link} , fusion thresholds θ_{fuse} , s_{\min} , and forgetting thresholds θ_{forget} , d_{forget} , all of which determine graph construction, pruning, redundancy reduction, and forgetting. Finally, the fourth group covers dynamics and stopping criteria: velocity decay β , force-velocity scaling α (denoted α_{dyn} to distinguish it from Eq. (2)), integration step size Δt , and stopping thresholds ε , δ for energy magnitude and change (cf. Eq. (9)).

Initialization Strategy. In both dialogue and reasoning scenarios, we adopt a coarse-to-fine initialization. Repulsion α is typically set to a moderate value, while semantic attraction β_{ij} is sparsely initialized only for pairs with $W_{ij}>0$ to avoid early collapse. Origin attraction δ_i is scaled with access frequency to create an attention center, and origin repulsion γ_i provides peripheral dispersion and forgetting. Gains $\Delta\delta_{\rm direct}$ and $\Delta\beta_{\rm assoc}$ are initialized as small increments so that link weights and attention pulls increase gradually, while the temporal decays μ_{δ} , μ_{β} are chosen conservatively to prevent oscillation. Structural thresholds are set so that new links are established only when nodes are both semantically similar and spatially close, fusion requires both high similarity and low distance, and forgetting is triggered only for nodes that are simultaneously inactive and spatially distant. For dynamics, the velocity decay β suppresses oscillations, $\alpha_{\rm dyn}$ controls the translation of forces into velocity, and Δt is chosen such that single-step displacement is small compared to cluster scale. The energy-based stopping rule ensures termination when either the total energy E(t) falls below ε or its change is smaller than δ .

Stepwise Heuristic Tuning. The tuning procedure proceeds in several stages. First, retrieval alignment is ensured by adjusting the initialization of β_{ij} and the associative gain $\Delta\beta_{\rm assoc}$, so that semantically related nodes become connected while avoiding premature link removal. Next, stability is achieved by increasing velocity decay or reducing $\alpha_{\rm dyn}$ and Δt , after which the thresholds ε and δ are tightened so that energy decreases smoothly and convergence occurs within finite steps. Sparsity and forgetting are tuned by modifying μ_{δ} , μ_{β} together with $(\theta_{\rm forget}, d_{\rm forget})$, allowing long-tail nodes to be removed without harming performance. Fusion is then optimized by searching over $(\theta_{\rm fuse}, s_{\rm min})$ and gradually relaxing thresholds to balance abstraction and granularity, with fused positions and velocities computed by weighted averages as in Eq. (6). Finally, task-specific adaptation is applied: for long dialogues, higher baseline δ_i and larger $\Delta\delta_{\rm direct}$ reinforce central clustering of frequently

 accessed memories, while for long-horizon reasoning, stronger γ_i and stricter sparsification accelerate the decay of peripheral noise.

Monitoring and Early Stopping. In addition to task metrics, several signals are monitored during training. The energy curve E(t) is inspected for monotonic decrease and plateau length, the number of active nodes and average degree are tracked along with the ratio of added versus dropped links, and the frequency of fusion and forgetting events is measured to quantify their marginal influence on generation quality. These indicators help diagnose oscillations, over-pruning, or excessive memory growth.

Implementation Notes and Practical Summary. In practice, the force and update equations (Eqs. (2)–(5), (6), (9)) must be faithfully implemented. The parameters β and $\alpha_{\rm dyn}$ are the most critical for stabilizing dynamics, as they directly regulate oscillation. Empirically, we first balance repulsion, attraction, and velocity on a development set so that semantic clusters form without collapse. After stabilization using energy-based early stopping and mild temporal decay, structural pruning is performed via fusion and forgetting. Ablation results confirm that all four forces are necessary: removing any one of them degrades alignment, clustering, or separation, underscoring the necessity of the multi-force design.

A.5 SUPPLEMENTARY INFORMATION ON EXPERIMENTAL SETUP

To ensure fair and reproducible comparison across all baselines and our proposed MemoryField framework, we detail the configuration settings for each task category as follows: **Dialogue Evaluation.** Multi-session Chat (MSC) is a benchmark dataset designed to evaluate long-term dialogue capabilities. It consists of multi-turn, multi-session conversations that span various topics and personas. The dataset challenges models to maintain coherent context across discontinuous dialogue turns, emphasizing long-range dependency handling.

Conversation Chronicles (CC) simulates natural, evolving dialogues involving complex conversational goals and topic transitions. The dataset provides annotations for session segmentation and context shifts, making it suitable for evaluating a model's ability to track memory and maintain contextual coherence in long-term interactions.

Reasoning Evaluation. Single-hop Reasoning data involves answering questions that require retrieving and reasoning over a single piece of information. It serves as a baseline to assess fundamental retrieval accuracy.

Multi-hop Reasoning tasks require the model to connect multiple facts across different documents or parts of a context. Datasets like HotpotQA and 2WikiHop are representative of this category.

Temporal Reasoning tasks evaluate the model's ability to interpret and reason over time-based information, such as sequences of events, timelines, or date-based dependencies.

Open-domain QA involves answering diverse questions using broad knowledge. This category includes datasets like Natural Questions (NQ) and MuSiQue, which test general knowledge and contextual comprehension.

Adversarial Reasoning challenges the model with deliberately confusing or misleading inputs to test its robustness. These examples often include distractors or irrelevant but plausible-sounding information.

These datasets jointly form a comprehensive benchmark suite to evaluate the proposed **MemoryField** framework under both conversational and reasoning-intensive conditions.

A.6 Training Examples

Step 1

User Input: I study natural language processing.

Recent History:

-

1026 Retrieved Memory: 1027 1028 1029 • Current Input: [I study natural language processing] 1030 Commencing association update (current number of nodes: 1) 1031 Semantic Similarity Matrix: 1032 1033 [[1.00]] 1034 Position Update: 1035 Node 0 displacement: 0.004, speed: 0.036 1036 • Position matrix transformation: 1037 - Current number of nodes: 1 1039 Position matrix shape: (1, 128) 1040 $[[-0.0806906 \ 0.19792689 \ -0.47895087 \ ...]]$ 1041 Node distance from origin 1042 Node 0 distance to origin: 3.596 1043 1044 Response: That's interesting! Natural language processing is a rapidly evolving 1045 field with many real-world applications such as dialogue systems and information retrieval. 1046 1047 1048 Step 2 1049 1050 **User Input:** I often work on large language models. 1051 1052 • Recent History: 1053 - User: I study natural language processing 1054 AI: That's interesting! Natural language processing is a rapidly evolving field 1055 with many real-world applications. 1056 • Retrieved Memory: 1057 User: I study natural language processing 1058 **AI:** That's interesting! Natural language processing is a rapidly evolving field with many real-world applications. • Current Input: [I often work on large language models] 1062 • Commencing association update (current number of nodes: 2) 1063 Semantic Similarity Matrix: 1064 [[1.00 - 0.01],1065 $[-0.01 \ 1.00]]$ 1066 1067 • Position Update: 1068 - Node 0 displacement: 0.003, speed: 0.033 1069 - Node 1 displacement: 0.004, speed: 0.036 1070 Position matrix transformation: 1071 Current number of nodes: 2 1072 - Position matrix shape: (2, 128) 1074 $[[-0.08283248 \ 0.19459349 \ -0.47730172 \ ...],$ 1075 [0.17930676 0.16485462 0.31871097 ...]] 1076 • Node distance from origin 1077 Node 0 distance to origin: 3.596 1078 - Node 1 distance to origin: 3.564 1079

- Fusion evaluation between Node 0 and Node 1 yielded a semantic similarity of -0.01.
- Response: Large language models are powerful tools for many NLP tasks, but they also require careful design to manage memory and computation efficiently.

Step 3

User Input: I am very interested in scaling language models.

- Recent History:
 - User: I study natural language processing
 AI: That's interesting! Natural language processing is a rapidly evolving field with many real-world applications.
 - User: I often work on large language models
 AI: Large language models are powerful tools for many NLP tasks, but they also require careful design to manage memory and computation efficiently.
- Retrieved Memory:
 - User: I study natural language processing
 AI: That's interesting! Natural language processing is a rapidly evolving field with many real-world applications.
 - User: I often work on large language models
 AI: Large language models are powerful tools for many NLP tasks, but they also require careful design to manage memory and computation efficiently.
- Current Input: [I am very interested in scaling language models]
- Commencing association update (current number of nodes: 3)
- Semantic Similarity Matrix:

```
[[1.00 -0.01 0.02]
[-0.01 1.00 0.92]]
[0.02 0.92 1.00]]
```

- Position Update:
 - Node 0 displacement: 0.003, speed: 0.033Node 1 displacement: 0.004, speed: 0.036
- Position matrix transformation:
 - Current number of nodes: 3
 - Position matrix shape: (3, 128)

```
[[-0.08283248 0.19459349 -0.47730172 ...],
[0.17930676 0.16485462 0.31871097 ...],
[0.05350199 0.44217721 -0.5286078 ...]]
```

- Node distance from origin
 - Node 0 distance to origin: 3.592Node 1 distance to origin: 3.626
 - Node 2 distance to origin: 3.671
- Fusion evaluation between Node 0 and Node 1 yielded a semantic similarity of -0.01.
- Fusion evaluation between Node 0 and Node 2 yielded a semantic similarity of 0.02.
- Fusion evaluation between Node 1 and Node 2 yielded a semantic similarity of 0.92. [Node Fusion] Trigger condition met (similarity i, 0.92), initiating fusion...
- Response: Scaling laws and efficient training strategies are indeed key to advancing large language models. Which aspects of scaling do you find most compelling—compute efficiency, model architecture, or application performance?

A.7 LIMITATIONS AND ETHICS Engineering and Scalability. The computational efficiency of MemoryField under very large memory sizes (e.g., 10^4 + nodes) has not been systematically evaluated. Future work may require approximate techniques such as cluster-level interactions, sparse neighborhood updates, or local energy minimization to ensure practical scalability. Parameter Sensitivity. Although heuristic tuning yields stable performance across multiple tasks, a systematic sensitivity analysis and the development of automated tuning strategies remain open directions. Convergence Analysis. At present, the framework relies on empirical thresholds for termination rather than a formal proof of energy convergence. Providing theoretical guarantees on stability and convergence would further strengthen the method. Privacy and Compliance. Explicit and persistent memory may raise privacy risks. We recommend incorporating mechanisms for memory sanitization, access control, and user authorization. This includes anonymization or encryption of stored data, user-visible memory records, and the ability for users to delete memory items. Practical guidelines and risk assessment checklists are provided in the appendix. Responsible Use of LLMs. This article used large language models (such as ChatGPT) only as auxiliary tools for language polishing. They were not employed in research conception, experimental design, or the generation of academic content.