

# Long Term Memory : The Foundation of AI Self-Evolution

Xun JIANG <sup>$\mu\theta$</sup>     Feng LI <sup>$\theta*$</sup>     Han ZHAO <sup>$\theta*$</sup>     Jiahao QIU <sup>$\iota*$</sup>     Jiaying WANG <sup>$\theta*$</sup>   
 Jun SHAO <sup>$\theta*$</sup>     Shihao XU <sup>$\theta*$</sup>     Shu ZHANG <sup>$\theta*$</sup>     Weiling CHEN <sup>$\theta*$</sup>     Xavier TANG <sup>$\theta*$</sup>   
 Yize CHEN <sup>$\theta*$</sup>     Mengyue WU <sup>$\alpha$</sup>     Weizhi MA <sup>$\sigma$</sup>     Mengdi WANG <sup>$\iota$</sup>     Tianqiao CHEN <sup>$\mu\theta$</sup>

<sup>$\mu$</sup>  Tianqiao and Chrissy Chen Institute

<sup>$\iota$</sup>  Princeton University

<sup>$\sigma$</sup>  Institute for AI Industry Research, Tsinghua University

<sup>$\alpha$</sup>  Shanghai Jiao Tong University

<sup>$\theta$</sup>  Shanda Group

## Abstract

Large language models (LLMs) like GPTs, built on vast datasets, have demonstrated impressive capabilities in language understanding, reasoning, and planning, achieving performance comparable to humans in various challenging tasks. Most studies have focused on further enhancing these models by training them on ever-larger datasets, aiming to develop more powerful foundation models. However, while training stronger foundation models is crucial, we propose how to enable models to evolve while inference is also vital for the development of AI, which refers to AI *self-evolution*. Compared to using large-scale data to train the models, the self-evolution may only use limited data or interactions. Drawing inspiration from the columnar organization of the human cerebral cortex, we hypothesize that AI models could potentially develop emergent cognitive capabilities and construct internal representational models through iterative interactions with their environment. To achieve this, we propose that models must be equipped with Long-Term Memory (LTM), which stores and manages processed real-world interaction data. LTM not only enables the representation of long-tail individual data in statistical models but also facilitates self-evolution by supporting diverse experiences across various environments and agents. In this report, we first explore the concept and significance of AI Self-Evolution, focusing on its potential to enhance AI models during the inference stage. We examine the role of LTM as a key mechanism for enabling lifelong learning in AI systems, allowing models to continually evolve based on accumulated interactions and experiences. Next, we detail the structure of LTM and the corresponding data systems required to facilitate high-quality data acquisition and retention, ensuring the effective representation of individual data. Finally, we classify various approaches for constructing personalized models using LTM data and discuss how models enhanced by LTM can achieve self-evolution through interaction with their environments. Based on LTM, our multi-agent framework OMNE achieved first place on the GAIA benchmark. This demonstrates the great potential of utilizing LTM for AI Self-Evolution and solving real-world problems. We present our technical roadmap and discuss potential avenues for future research. We believe that advancing research in LTM is critical for the ongoing development and practical application of AI technology, especially for self-evolution. We hope this work will inspire more researchers to contribute to the exploration of this exciting and evolving field.

\*Equal contribution.

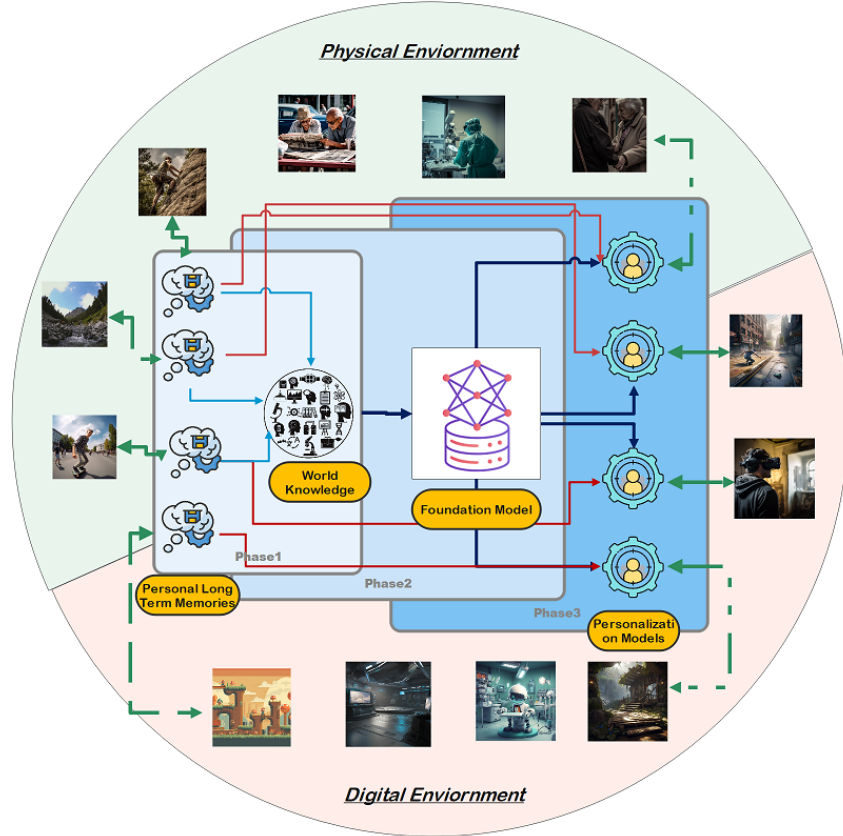


Figure 1: The overview of LTM and AI Self-Evolution.

## 1 Introduction

Artificial Intelligence (AI) is recognized as a key technology in the Fourth Industrial Revolution[1], empowering machines to perceive their environments and act intelligently through algorithms and software to optimize the achievement of various objectives. AI technologies are now widely applied in areas such as finance, education, and healthcare. In recent years, the development of large language models (LLMs) and LLM-powered agents has significantly enhanced AI capabilities, which are more powerful in solving various challenging tasks under diverse scenarios.

At its core, a model can be understood as an advanced form of data compression. A classic example is Newton’s law of universal gravitation, which condenses complex astronomical data into a simple mathematical formula. This compression represents large amounts of data in a concise form. Similarly, LLMs compress vast amounts of text corpora into statistical patterns to generate coherent text [2]. However, we argue that intelligence is not limited to learning from existing data, the self-evolution ability is also important for the development of AI models, which is similar to the evolution ability of humans. Tasks from different scenarios often have distinct data distributions and diverse ability requirements, where the self-evolution ability will enable AI models to adapt to new tasks by learning limited interactions for powerful performances. Self-evolution will also contribute to diverse models, which can further help the development of AI, especially LLM models in recent years.

### 1.1 Phases of Model Evolution processes of LLMs

To better understand the need for self-evolution for LLMs, as shown in Figure 1, we propose to break down the model evolutionary process of LLMs into three main phases. These phases highlight the gradual progression from simple pattern recognition to self-evolved personalized intelligence.

- **Phase 1: Cognitive Accumulation in the Physical World.** Data accumulation is the first and vital step for the development of AI, which is achieved by humans through continuous practical interactions. In the process of understanding the world, humans have also evolved stronger abilities to discover and apply patterns, but the first step is cognitive accumulation. Individuals interact with their environment, producing diverse and personal fragmented cognition pieces. Some of these cognitive fragments are digitized and stored, which can be used to construct AI model. While others remain in individual minds, contributing to diverse personalities.
- **Phase 2: Constructing Foundation Models in the Digital World.** AI models attempt to learn from the data accumulated by humans in Phase 1, and have achieved promising results, especially LLMs. LLMs consolidate all digitized cognitive fragments to form a unified “average” model (foundation model). These models reflect commonalities and general patterns in large-scale data, making them suitable for a broad range of language generation tasks. However, these models, while statistically efficient, overlook the expression of personalized information and struggle with handling long-tail data or rare scenarios. We think the main reason is they ignored the remained individual evolution of humans.
- **Phase 3: Model Self-Evolution to Achieve Stronger Intelligence.** The third phase moves beyond averaged intelligence, focusing on building self-evolving, personalized intelligent models. To address the complexity and sparsity of personalized data, future model architectures must break away from the existing “global average” paradigm and shift towards more flexible and adaptive distributed intelligence architectures, even with limited interactions in a new task/scenario. This self-evolution ability will also contribute to more diverse and stronger intelligence models through their dynamic and continual evolution. Furthermore, the most promising thing is multi-agent evolution based on single self-evolution.

Most existing work focuses on how to construct better data and use it to train a more powerful foundational model, which is essentially research centered around phases 1 and 2. There is also currently a popular view that *Architectures aren’t fundamentally important in the curve-fitting paradigm, while the critical factor is the dataset* [3]. This perspective applies to the second phase, but in the third phase, architecture becomes as important as data. The core challenge lies in **how to effectively express small amounts of individual data within the foundation of statistical models**. Our research focuses on how to distill individual data to ensure more efficient expression within statistical models. At the same time, we are exploring new model architectures that can better support these refined data, as well as investigating how intelligent agents can collaborate to achieve self-evolution with enhanced individual data. This is based on our belief that second-phase average models will continue to strengthen, providing a foundation for future-oriented designs.

## 1.2 Principles to Achieve Model Self-Evolution

The ability of a model to self-evolve is crucial for its long-term adaptability and personalization, and this depends heavily on an effective memory mechanism. In this context, we propose that long-term memory (LTM) provides the historical data accumulation and experiential learning capacity necessary for continuous model evolution. Just as humans refine their cognition and behavior through experience and memory, LTM enables a model to gradually optimize its reasoning and learning capabilities when dealing with long-term, dispersed, and personalized data.

### 1.2.1 Empower Foundation Models with LTM Data for Self-Evolution

In traditional LLMs, updating the model typically requires adjustments to all parameters, which is impractical for processing individual-specific data [4]. A more optimal approach is to employ localized updates, allowing the model to adapt to sparse [5], personalized LTM data without compromising the stability of the global model. This method addresses the issue of individual data being “averaged out” in current models, enabling more comprehensive expression of personalized information.

Techniques such as Retrieval-Augmented Generation (RAG) with In-Context Learning (ICL) and Low-Rank Adaptation (LoRA) for fine-tuning (SFT) can be seen as ways to locally update individual data. We have developed a mixed strategy to integrate LTM data, yielding promising results in practical applications. However, this may not a perfect solution, and we are continuing to explore

how to effectively integrate long-tail individual data into the model’s memory mechanisms, hoping to attract more researchers to contribute to this field of exploration.

### 1.2.2 Real-Time Weight Updates Combined With LTM Data for Self-Evolution

Current LLMs typically separate the inference and training phases, where model weights are frozen during inference, preventing adjustments and learning based on new input [6]. This fixed inference process limits the model’s adaptability, particularly when handling personalized tasks and real-time learning. Inspired by the human brain’s updating mechanism, we believe that future LLMs should integrate inference and training with LTM, enabling the model to dynamically adjust weights upon receiving new information, akin to the continuous learning ability of humans. We also provided an overview of early work in the following sections, demonstrating how the integration of real-time training and inference allows LLMs to become more flexible and quickly adapt to new tasks or long-tail data. Additionally, this integration could help the model self-reflect and correct faulty reasoning paths when faced with complex inference tasks, improving both accuracy and efficiency. This dynamic self-adjustment capability would greatly enhance the model’s personalization capacity and its potential for long-term evolution.

With LTM, a model can not only learn from short-term memory but also extract valuable insights from historical data, forming a deeper understanding of individual preferences and behavior patterns over time [7]. This understanding lays a solid foundation for personalized customization and dynamic adjustments, allowing the model to evolve more effectively. Especially when faced with new or extreme situations, LTM enables the model to reference past experiences, quickly make adjustments, and self-evolve, thereby achieving greater flexibility and adaptability.

### 1.3 The Implementation Path of LTM in Model Self-Evolution

Inspired by the importance of LTM for humans, we argue that research into long-term memory is essential for advancing model personalization. As AI models and intelligent agents continue to evolve, their foundational capabilities, akin to an increase in machine intelligence “brain capacity”, provide greater support for the integration of long-term memory into personalized models. While efforts have been made to construct memories or experiences for evolving AI systems, significant gaps remain in defining, constructing, and evaluating long-term memory for AI, hindering the development of personalized LLMs. We begin by defining AI self-evolution and LTM, exploring the key role of LTM within it. We then focus on how LTM can be utilized to enable AI self-evolution. Our research focuses on three questions:

1. **What is AI self-evolution, and what constitutes long-term memory?** Why does AI self-evolution require models to have personalized capabilities? Why is long-term memory essential for achieving true personalization? What are the shortcomings of memory mechanisms in current language models, and how can these deficiencies help us refine the definition of LTM?
2. **How to construct LTM for self-evolution?** What types of data are most suitable for forming the foundation of a model’s personalized long-term memory, and how can we distill and structure this raw data into LTM?
3. **How to use LTM for AI self-evolution?** How can we efficiently process and utilize individual data to continuously update long-term memory, ensuring that it not only understands individual preferences but also enables self-evolution by adapting and coordinating with the environment as data grows and changes?

The main contributions of our study are summarized as follows:

- **Definitions of AI Self-Evolution and LTM.** We provide an in-depth discussion of the relationship between AI self-evolution and LTM, proposing a systematic framework that highlights the core role of LTM in the process of AI evolution. Through LTM, AI not only addresses personalized needs but also continuously learns and optimizes, bridging the gap between general models and truly personalized intelligent systems. By effectively handling individual long-tail data, the long-term memory mechanism significantly enhances the individual capabilities and diversity of agents, laying the foundation for AI self-evolution.

- **Data Framework for LTM.** To implement long-term memory, we developed a data collection, analysis, and synthesis framework that allows for differentiated system deployment based on various business scenarios. To verify the generalization ability of this framework, we deployed independent data systems in two distinct business contexts—office collaboration and health management. Specialized intelligent agents collaborate by integrating data from individual sub-models into a unified long-term memory. Each agent focuses on specific aspects of the data, ensuring seamless and accurate personalization even when data is sparse or inconsistent. Based on this data framework, we successfully established the world’s largest real-user voice dataset for mental health (see Section 6.1.1), and augmented it through data synthesis (see Section 6.1.2). We are planning to open this dataset on a data platform to further support scientific research.
- **Development Framework for LTM.** We propose a multi-agent collaborative development framework (OMNE) based on LTM. In this framework, each agent has an independent system structure that allows for autonomous learning and storing of a complete world model, thereby constructing an independent understanding of the environment. Through this LTM-based collaborative development, AI systems can adapt in real time to changes in individual behavior, optimizing task planning and execution, further promoting personalized and efficient AI self-evolution. We detail this framework in Section 6.3, demonstrating its potential in leveraging individual data for decision-making.

Our research contributes both theoretically and practically by integrating LTM into model personalization to promote AI self-evolution, with progress already made in practical applications. First, we discuss the importance of AI self-evolution and the critical role of model personalization in Section 2. Next, we examine memory mechanisms in current LLMs and humans, exploring how human memory systems can inspire the design of LTM for model personalization, followed by a definition of LTM in Section 3. The questions of how to construct LTM (Question 2) and How can LTM be used to achieve model personalization in AI Self-Evolution (Question 3) are addressed in detail in Sections 4 and 5, respectively. Our efforts and results are presented in Section 6, with further discussions and conclusions summarized in Sections 7 and 8.

## 2 AI Self-Evolution

The process of AI self-evolution can be compared to the Thousand Brains Theory or biological individual evolution. In the Thousand Brains Theory proposed by Jeff Hawkins [8], the brain does not operate through a single centralized system; rather, it constructs an understanding of the world through thousands of mini-models in the neocortex. These mini-models function independently while working together to form a diverse, distributed intelligence system. This theory challenges the traditional linear understanding of the brain by emphasizing that intelligence arises from the parallel processing and collaboration of multiple models. The Thousand Brains Theory posits that each region of the brain can independently create maps of the world and interact with maps from other regions, resulting in a more accurate and comprehensive cognition. Therefore, intelligence evolves progressively through the collaboration of multiple independent models.

On the other hand, the history of biological evolution demonstrates that there is no single “superorganism” dominating ecosystems [9; 10]. Instead, the diversity driven by individual adaptations and mutations has enabled the formation and flourishing of the complex network of life we observe today.

Similarly, AI self-evolution can follow a path of multi-agent collaboration. In a multi-agent system, different agents interact, learn, and collaborate with one another to optimize their capabilities, generating personalized data. This personalized data serves as the foundation for continuous AI evolution, driving it from an initial general-purpose model to a system that increasingly adapts to individual needs. Just as biological evolution forms complex ecosystems through mutation and adaptation, AI self-evolution relies on diversity and co-evolution. This evolutionary path enables AI to continuously adapt to different environments and requirements, forming a more diverse and flexible intelligence system.

In the following sections, we will define AI self-evolution, break down the key system dependencies within AI self-evolution, and discuss how these capabilities contribute to more effective and adaptive self-evolution.

## 2.1 Definition of AI Self-Evolution

**Definition:** AI self-evolution refers to the process by which AI models achieve breakthroughs in multi-agent collaboration and cognition through continuous learning and optimization with personalized data. This process is based on a shared core architecture, where each model evolves by processing personalized experiences and data, thereby enhancing its reasoning capabilities and adaptability, ultimately achieving autonomous learning and continuous evolution in dynamic environments.

Model self-evolution enables AI models to continuously learn from personalized data, adapting to ever-changing environments and meeting diverse needs without relying heavily on human intervention. Throughout this process, models process and absorb new experiences, optimizing their architecture and outputs, evolving from generalized knowledge to more contextually adaptive personalized knowledge. Through dynamic learning mechanisms, models can retain and utilize key information from past interactions, supporting future decision-making and effectively mitigating issues like overfitting and data drift.

A key feature of model self-evolution is that it is based on a unified foundational architecture, ensuring that all model instances share a consistent core structure. However, the evolution of each model is driven by the unique experiences and data it processes, with differences between models arising from their individualized handling of personalized data. This approach ensures that, while models adhere to consistent internal rules and mechanisms, they can develop in differentiated ways according to personalized needs and environments. As evolution progresses, models become better at simulating individual behaviors, providing personalized and precise context-aware outputs, ultimately laying a solid foundation for multi-agent collaboration and cognitive breakthroughs.

## 2.2 Key System Dependencies for AI Self-Evolution

The realization of AI self-evolution does not happen spontaneously; it relies on a series of key system dependencies that provide the necessary foundation and framework for AI models to learn, optimize, and evolve into more personalized agents in ever-changing environments. These dependencies include not only the mechanisms of multi-agent collaboration but also the generation of personalized data, the construction of long-term memory, distributed model updating mechanisms, and self-correction mechanisms. These interdependent factors collectively drive the transition of AI from static models to self-evolving systems, helping AI transcend the limitations of traditional intelligence and gradually move toward the future of AI self-evolution.

### 2.2.1 Multi-Agent Collaboration Mechanism

The multi-agent collaboration mechanism is a key element of AI self-evolution, especially in handling complex tasks where efficient cooperation among multiple agents can significantly enhance the overall system performance[11]. With the rise of large language models (LLMs) demonstrating the phenomenon of "emergent intelligence," the capabilities of AI systems have been greatly enhanced, accelerating the development of AI models[12]. These LLMs, with their vast number of parameters, endow agents with stronger memory, reasoning, and adaptability, allowing AI to perform remarkably well in more complex tasks. However, from the perspective of multi-agent collaboration, model personalization becomes the core factor promoting agent collaboration and evolution, particularly in solving complex problems and tasks.

Just as in human intelligence evolution, where increased brain capacity enhanced memory, thinking, and reasoning abilities, propelling the development of civilization[13], LLMs similarly drive leaps in AI capabilities. However, current research suggests that in more complex scientific and engineering problems, human collaboration is essential. Similarly, whether AI systems will experience a second wave of emergent intelligence in the future depends on whether collaboration among multiple agents can be elevated to a new level.

Some explorations suggest that in small-scale multi-agent collaborations, increasing the number of agents does bring some performance improvements, but these gains are not consistent or stable[14]. We believe the key issue behind this lies in the fact that most current multi-agent collaborations are still limited to role-playing interactions, where the agents' capabilities and knowledge are often homogeneous, lacking the differentiated skills required for deep collaboration. To achieve

significant breakthroughs in large-scale multi-agent collaboration, it is crucial to rely on agent personalization, which can provide each agent with unique expertise and abilities, thus promoting efficient collaboration and evolution of the system.

Therefore, the true potential of multi-agent collaboration depends on the creation of a group of differentiated, highly personalized agents. This not only supports the resolution of complex tasks but also provides a new path for the self-evolution of AI systems. Model personalization is an indispensable part of this evolutionary process, enabling different agents to contribute uniquely in collaborations, pushing AI systems toward higher levels of intelligence.

### 2.2.2 Differentiated Personalized Models

Personalized data generation is one of the core driving forces behind AI self-evolution, especially in multi-agent systems where having personalized models for each agent is crucial for ensuring system diversity and efficient collaboration. As AI application scenarios become increasingly complex, models must continuously acquire, process, and respond to personalized data to dynamically adjust to the needs and preferences of different individuals, truly meeting diverse task requirements[15]. Each agent, through its personalized model, can not only handle specific tasks independently but also contribute unique insights and experiences during collaboration with other agents, thereby generating more diverse and personalized data within the overall system.

This diversity brings stronger collaborative capabilities to AI systems. Through effective interactions among agents, personalized models can better respond to the needs of different individuals and support continuous model evolution through the accumulation of long-term data and continuous learning. In complex fields like healthcare[16], the requirement for model personalization becomes especially prominent when handling multimodal and heterogeneous data. A one-size-fits-all strategy shows limitations in addressing these complex tasks, while personalized models with differentiated processing capabilities can dynamically adapt to individualized scenarios, delivering precise and efficient performance.

Moreover, even within the same task scenario, different individuals will have significantly different expectations of the model[17]. Through personalized models, AI systems can dynamically adjust according to these varied needs, providing highly tailored services to each individual. For example, in dialogue systems, individuals have unique preferences for the style, tone, and even response format of the model’s output. To achieve precise personalized services, models must have dynamic learning capabilities, deepening their understanding of individual needs through continuous interaction.

Therefore, personalized models not only promote the diversity and collaboration of AI systems but also support the self-optimization and evolution of the entire system by generating more personalized data. This process forms a positive feedback loop: through the collaborative work of multiple agents, AI systems generate increasingly personalized feedback, continuously meeting individual needs and further enhancing model capabilities, ultimately achieving true self-evolution.

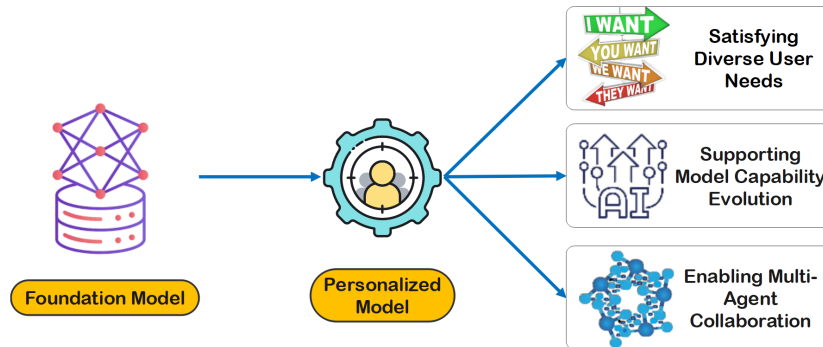


Figure 2: Model personalization in AI self-evolution

### 2.2.3 Self-Correction and Evolution Mechanism

To achieve AI self-evolution, models must possess a self-correction mechanism [18]. This mechanism not only ensures the internal consistency of the model but also enables it to adapt to changes in the external environment. Self-correction is at the core of AI self-evolution, allowing models to update their cognition and behavioral strategies through continuous feedback loops. This is similar to the process of biological evolution, where selection and adaptation drive continuous optimization to fit the environment.

Recent mathematical studies suggest that predictable information often exists in the low-dimensional structures of high-dimensional data [19] [20] [21]. Understanding and leveraging this concept is crucial for AI self-evolution. By identifying the low-dimensional structures within data, AI can learn and generalize more effectively, thus enhancing its autonomous learning and evolutionary capabilities.

### 2.2.4 Long-Term Memory and Learning Ability

Long-Term Memory (LTM) is a fundamental cornerstone in the process of AI self-evolution, providing models with the ability to accumulate historical experiences and knowledge, enabling continuous optimization through long-term interaction and learning. LTM can store data not only at the individual level but also accumulate data over time, helping models adjust their responses and behaviors based on this information, thus facilitating self-evolution. Compared to most personalized approaches that rely on a context window, long-term memory overcomes the limitations of short-term approaches by endowing models with the ability for continuous learning and self-improvement, enabling them to exhibit stronger adaptability when faced with complex environments and multi-agent collaboration.

Just as humans use long-term memory to shape their behavior and identity, AI systems can also utilize LTM to provide customized responses based on individual data. This “individual data” is not limited to user interaction data but can also include the specific needs of an organization or domain, allowing the model, through long-term accumulation, to surpass the framework of general knowledge and support more precise decision-making and behavioral adjustments. Therefore, LTM is crucial for enabling AI to achieve continuous self-evolution.

This paper emphasizes the role of long-term memory as a core driver of AI self-evolution. To truly realize AI self-evolution, we need models equipped with LTM, which can store and manage long-term interaction data in real-world environments and effectively represent long-tail data from individuals. Through this expression of diversity, models can continuously adjust themselves in collaboration with the environment and multiple agents, thus promoting the process of self-evolution. As such, long-term memory is not only the foundation for personalized data generation but also a key mechanism supporting the long-term adaptation and optimization of AI systems. In subsequent chapters, we will discuss the critical role of long-term memory in AI self-evolution in greater detail.

## 2.3 Thought Experiment: From Euclidean Geometry to Riemannian Geometry

An interesting question is: *Can a large language model deduce Riemannian geometry from the five axioms of Euclidean geometry?* This would not only require the model to perform logical deductions based on existing axioms but also to examine and challenge fundamental assumptions—such as the parallel postulate—thereby stepping into the entirely new domain of non-Euclidean geometry. Current LLMs are not yet capable of such creative leaps, as they rely on existing data for reasoning and lack the ability to propose new hypotheses and extend the boundaries of knowledge.

Suppose LLMs could continuously adapt and update with personalized data, gradually developing sensitivity to atypical patterns—could they eventually overcome this limitation? In this thought experiment, the volume of data is not a constraint, as synthetic data could be continuously generated. If the model could not only process human-provided data efficiently but also reflect on its own knowledge, dynamically adjust, and propose new hypotheses, using experimental deduction to validate them, then such a system would no longer be just a passive tool for knowledge storage and reasoning. Instead, it would become a cognitive entity capable of self-evolution, breaking through established mathematical frameworks and pioneering new theoretical realms.

From this perspective, the potential of real-time updates and personalized learning lies not only in improving the accuracy of existing knowledge processing but also in granting the model higher levels of cognitive flexibility. This would enable it to genuinely participate in the discovery of new



knowledge. Realizing this potential may be the key to overcoming the current limitations of LLMs, evolving them from data-driven intelligence into explorers capable of expanding the boundaries of knowledge.

## 2.4 Future Directions and Challenges

As explored in this chapter, personalized models lay the foundation for multi-agent collaboration and cognitive breakthroughs in increasingly complex task environments. By integrating reasoning with training, allowing local updates, and adopting distributed architectures, AI will not only simulate human language outputs but also develop human-like reasoning and innovation abilities. Ultimately, AI will break through current technological limitations by interacting with the physical world, autonomously learning, and continuously evolving, driving the expansion of cognitive boundaries.

While the prospect of AI self-evolution is promising, there remain significant technical and theoretical challenges in achieving this goal. After establishing the foundation for self-evolution through model personalization, key challenges for AI self-evolution include:

- How can models be effectively evaluated and achieve self-evolution?
- How can collaboration mechanisms among agents be designed?
- How can we break through the current Scaling Law of models and continuously improve performance?

To address these questions, further discussions are presented in section 7. Our future research will focus on developing AI systems with autonomous learning, exploration, and evolutionary capabilities. In the next chapter, we will focus on the importance of long-term memory in AI self-evolution.

## 3 LTM for AI Self-Evolution

AI self-evolution emphasizes the ability of AI systems to improve their capabilities through continuous learning and adaptation. In this process, the retention and updating of individual model information by standalone AI models is a crucial feature of the entire AI system’s self-evolution. However, most current methods for implementing LTM in models primarily rely on context windows. These models typically often utilize immediate context or recent individual interactions to generate responses [22]. While this approach can achieve a certain degree of diversity, it has significant limitations in supporting long-term learning and continuous adaptation, hindering models from achieving true self-evolution.

Therefore, in this chapter, we will explore in depth the crucial role that LTM plays in AI self-evolution. We first define LTM in the context of AI self-evolution and analyze the shortcomings of current LLM memory mechanisms. We then discuss enhancing AI models’ self-evolution capabilities by drawing inspiration from human LTM characteristics and addressing the challenges and potential solutions in achieving this goal. Through these discussions, we aim to provide new ideas and directions for building AI systems capable of continuous learning and self-improvement.

### 3.1 Definition of LTM in AI Self-Evolution

**Definition:** LTM is the information that can be retained and utilized by AI systems over extended periods, enabling models to adjust their responses and behaviors based on a broader context.

Just as humans use LTM to shape their behavior and identity, AI systems can employ similar mechanisms to customize their responses and behaviors based on individual data. Here, “individual data” is not limited to individual users but also includes specific organizations and domains, allowing models to adjust their responses and behaviors according to broader individual contexts and needs. This forms the basis for creating personalized models that go beyond mere general knowledge.

LTM can be understood as a vast and complex repository of refined knowledge, shaped over time by a group of independent yet harmoniously interacting agents—similar to cortical columns in the brain. Each agent functions as an autonomous unit capable of learning, refining, and storing a comprehensive model of its own corner of the world. However, these agents do not operate in

isolation; they contribute their insights to the broader LTM collective, creating a shared knowledge base that supports deeply personalized interactions.

Unlike traditional static data storage systems, the LTM framework is a dynamic and distributed memory framework, akin to a network of independently operating thoughts in the human brain, where insights from various independent learning processes can merge. Just as a society comes together to form a more coherent understanding, this collective intelligence enables the system not only to accumulate knowledge but also to synthesize it in ways that better reflect the complexity and nuances of user needs, ultimately achieving AI system self-evolution. LTM demonstrates a more nuanced and comprehensive understanding of both individuals and collectives, enabling the system to respond to personal needs with a level of granularity that reflects this complexity. In this sense, LTM transcends mere data storage—it becomes an adaptive, continuously evolving cognitive organism, constantly refining itself in response to its environment, much like the human cognition it seeks to emulate.

### 3.2 Limitations of Current LLM Memory Mechanisms

LLMs like GPT-4 and Gemini demonstrate advanced intelligence and a comprehensive understanding of the world. However, to achieve true self-evolution, these models must be able to effectively process, store, and integrate information acquired through continuous interaction with various environments. Currently, LLMs primarily manage information through two memory mechanisms: contextual memory and compression-based parametric memory. While these mechanisms perform excellently in short-term tasks, they still fall short in supporting long-term autonomous learning and evolution.

#### 3.2.1 Memory through Prompting

Current LLMs utilize prompts as a form of contextual memory to retain and leverage information during inference. The prompt, which includes both the instruction and relevant context, serves as a temporary memory buffer, allowing the model to process and generate content based on the given context. This mechanism enables LLMs to perform a wide range of tasks without task-specific fine-tuning. However, from a self-evolution perspective, this prompt-based memory mechanism has several key limitations:

- **Temporality:** Stored context information is limited to the scope of the current task. Once a single inference call is completed, the model discards its previous state and sequence information, meaning it cannot utilize previously acquired knowledge in subsequent tasks.
- **Absence of Continuous Learning:** Unlike systems with explicit memory mechanisms (e.g., Long Short-Term Memory (LSTM) networks [23]), LLMs do not have an intrinsic ability to accumulate and refine knowledge across multiple interactions or tasks based solely on prompts.
- **Limited Cross-Task Knowledge Integration:** While prompts can provide task-specific context, they do not facilitate the automatic integration of knowledge across diverse tasks or domains. This hinders the model’s ability to develop a cohesive understanding that evolves over time.
- **Dependency on External Curation:** The quality and relevance of contextual information heavily rely on how prompts are crafted by users or systems. The model itself cannot autonomously curate or optimize its contextual knowledge base.

#### 3.2.2 Parametric Compressed Memory

Another memory mechanism is compression-based parametric memory, which forms a type of LTM by compressing world knowledge into the model’s parameters [24]. This mechanism allows models to retain key information over longer periods, but it has two significant drawbacks:

- **Unable to Update in Real-Time:** Since the model’s parameters are fixed during the training process, LTM cannot be easily updated once training is complete. This limits the model’s diversification and adaptability to the environment, as it cannot quickly learn to generate new memories to adjust its behavior.
- **Difficulty in Expressing Individual Data:** As statistical models, Transformers struggle to adequately represent individual data in their LTM. This typically leads to one of two

outcomes: either the model fails to retain individual data, or it overfits and forgets previously-stored world knowledge, resulting in memory bias. The current approach to mitigate this issue is through incremental training with carefully balanced data proportions to incorporate individual data into the model. However, this process is inefficient and difficult to scale in large-scale or dynamic data environments.

These limitations indicate that AI models need a more flexible and adaptive memory system that can retain individual data and achieve real-time adjustments, similar to human LTM.

### 3.3 Inspiration from Human LTM

To better understand LTM in AI systems and achieve AI system self-evolution, we can draw inspiration from the concept of human LTM. Existing neuroscience research suggests that personal memories generated through human interaction with the world are key factors in forming diverse and personalized behaviors. Current research typically divides human memory into three main types: working memory, short-term memory, and LTM. Working memory and short-term memory mainly contain temporary information related to current tasks or situations, which is quickly forgotten if not processed and transformed into LTM. Therefore, LTM can be considered the key data foundation for the formation of human personality.

Specifically, human LTM refers to the brain’s ability to store and retrieve information over extended periods, ranging from hours to decades[25]. Unlike short-term memory, which is temporarily used for immediate use, LTM is responsible for preserving knowledge and experiences that can influence our future behavior[26; 27]. This type of memory includes various subtypes, including episodic memory (personal experiences), semantic memory (general knowledge), and procedural memory (skills and habits)[28].

The formation of LTM involves multiple processes, including encoding, consolidation, and retrieval[29]. Encoding is the initial acquisition of information, while consolidation is the process of stabilizing new information and integrating it into existing memory networks[30]. Retrieval is the process of accessing and utilizing stored information when needed. These mechanisms are supported by neural processes in the brain, including the hippocampus and various cortical regions[31].

Moreover, LTM not only influences the formation of personal interests and habits but also plays a crucial role in the emergence of diverse needs[32]. Additionally, LTM significantly affects an individual’s knowledge accumulation[33], problem-solving abilities[34], social adaptability and self-regulation abilities[35], leading to different expressions of intellectual development and evolution. In social life, social experiences and memories stored in LTM help individuals build trust, promote knowledge and resource sharing, and ultimately drive cooperation and collaboration[36].

### 3.4 LTM in AI Models

In the previous section, we emphasized the important role of LTM in human evolution, which relies on a comprehensive mechanism for memory formation, updating, and utilization to ensure LTM effectively serves humans. Many aspects of model design and development draw inspiration from human cognitive and reasoning structures, such as the propagation mechanism of neural networks. So if LTM can aid the progress of human society, can it also be used for AI self-evolution? Our answer is affirmative, and we discuss this from the following aspects:

- **From the perspective of data accumulation:** Both models and humans interact extensively with their environment, providing foundational data for personalization. Compared to humans, AI models can interact with their environment more efficiently and can perform these interactions and iterations in purely virtual, digital environments. Therefore, by designing appropriate memory refinement strategies, models should be able to accumulate long-term memories like humans, possibly even with higher efficiency and scale.
- **From the perspective of model updates:** Artificial intelligence excels in storing and calling upon vast amounts of data, far surpassing the scale of human memory. Neural networks manage this data through distributed parameters, processing inputs from different domains. However, this storage is relatively rigid, lacking the flexibility for real-time updates and typically requiring retraining to implement updates. In contrast, human memory is highly

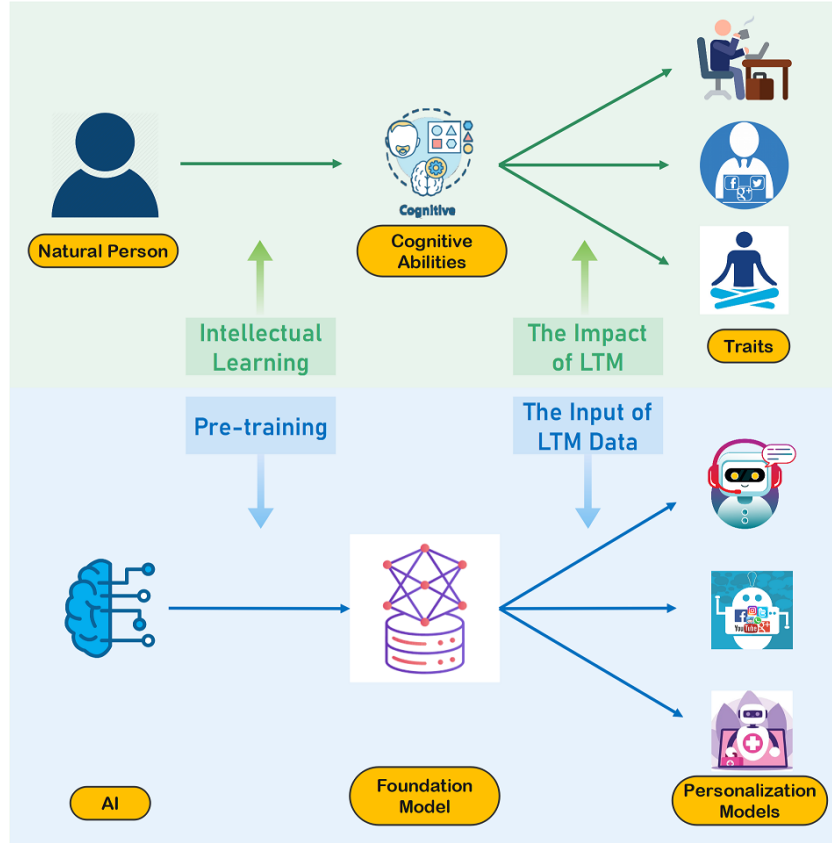


Figure 3: The difference between human learning and growth versus artificial intelligence training and evolution.

adaptive, quickly integrating new information and discarding outdated details through a process of "adaptive forgetting"[37]. This flexibility helps humans avoid cognitive overload and focus on the most relevant information at hand. To match this capability, AI systems need to develop dynamic update mechanisms that allow them to selectively update knowledge and discard outdated information without comprehensive retraining.

- **From the perspective of utilizing LTM:** Current advanced LLM memory mechanisms, such as contextual memory and parametric memory, can store and utilize large amounts of information with the advantage of non-forgetting. Moreover, AI can combine instant memory update mechanisms when encountering important new data, which is particularly useful in dynamic environments. This provides AI with certain advantages over human memory utilization, such as larger storage capacity and faster retrieval speed. However, these memory utilization methods still face significant challenges in managing dynamically accumulated long-term memories, especially in terms of flexibility and efficiency of memory updates, and still struggle to reach the level of human memory systems. Future research needs to focus on how to achieve efficient, flexible memory update and utilization mechanisms like humans.

In summary, the importance of human LTM to human society points the way for the development of AI systems. While human memory excels in context integration, flexibility, and real-time updates, AI performs better in managing large-scale datasets and identifying patterns. To fully leverage AI's advantages and address its shortcomings, we should combine AI's flexibility with its ability to process at scale, while introducing prioritization mechanisms and efficient information update methods similar to human memory. This hybrid approach can produce more personalized, responsive, and context-aware AI models, bringing them closer to the complexity and adaptability of human cognition.

### **3.5 The Importance of LTM for AI Systems**

#### **3.5.1 LTM in Single Models**

LTM is crucial for diversity and AI self-evolution because it allows models to generate new outputs based on a deep understanding of individual interactions. By recalling past interactions, identifying preferences, and adjusting context over time, LTM enables AI systems to infer new response patterns from historical data and provide adaptive answers in new environments. This capability is rooted in the information encoding within neural network weights and parameters, which is critical for supporting personalized interactions and continuous learning[38].

The distributed and dynamic nature of LTM allows AI systems to optimize interactions based on historical data, individual preferences, and context[39]. Ultimately, integrating LTM into tens of thousands of AI models marks a significant step towards AI evolution. The data from LTM lays the foundation for personalized and context-aware interactions, evolving continuously over time. As AI technology continues to advance, the development of LTM will become key to creating systems that not only understand individual needs but can also anticipate and adapt to unique requirements, thus providing truly intelligent experiences.

#### **3.5.2 The Role of LTM in AI Self-Evolution**

The ability for models to self-evolve is a highly promising direction in AI technology development and application. Just as human intelligence relies on personal interactions and feedback from the real environment, the LTM data accumulated by models during interactions provides the most critical data support for self-evolution. Of course, there are some differences:

Unlike human evolution, LTM-driven model evolution is not limited to real-world interactions. On one hand, models can interact with the physical environment and receive direct feedback like humans, which, after processing, will enhance their capabilities. This is also a key area of research in embodied AI. On the other hand, models can interact in virtual environments and accumulate LTM data, which has lower costs and higher efficiency compared to real-world interactions, thus achieving more effective capability enhancement.

Compared to existing methods of directly storing interaction history and reflection, refined LTM data can provide high-quality data support for model self-evolution, further accelerating the efficiency and effectiveness of model evolution. This provides more effective support for enhancing the capabilities and rapid application of large model (LLMs) driven agents.

Having clarified the importance of LTM for AI model self-evolution, we will focus on discussing the following two questions: 1) How to construct LTM (Section 4). 2) How to integrate LTM into models (Section 5).

## **4 How to Construct LTM?**

The construction of LTM represents a critical aspect in model personalization. Raw data serves as the foundation and source of LTM. Therefore, the relationship between LTM and raw data will be explored in the following sections. Subsequently, the method of refining raw data into LTM will be introduced, culminating in a review of existing methods to get high quality data for LTM, and the discussion consists of two main parts: data collection and data synthesis.

### **4.1 Raw Data to LTM**

Raw Data represents the comprehensive collection of all unprocessed data that a model receives through interactions with the external environment or during the training process. This data includes a wide range of observations and records, which may contain both valuable patterns and large amounts of redundant or irrelevant information. While Raw Data forms the foundational layer for the model's memory and cognition, it requires further processing to be effectively used for personalization or efficient task execution.

LTM refines and structures this raw data, making it usable by the model. This process enhances the model's ability to deliver personalized responses and recommendations. While Raw Data captures direct observations that may result in redundancy and clutter, LTM organizes this data into a structured

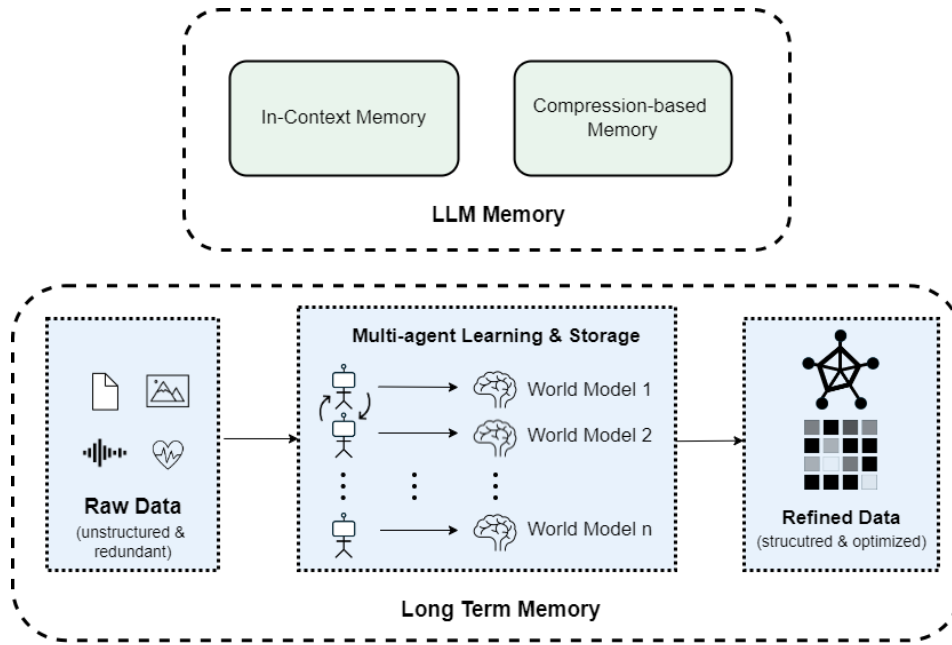


Figure 4: Comparison of LLM Memory and Long-Term Memory

memory, enabling the model to recognize patterns, remember preferences, and provide adaptive responses.

For instance, in an AI medical scenario, Raw Data might capture basic patient information, such as demographic details, current symptoms, and immediate diagnoses. However, Raw Data alone would struggle to effectively manage and integrate a patient’s comprehensive health history into a cohesive understanding of their condition.

In contrast, LTM refines and organizes this raw patient data into an interconnected structure, allowing the AI model to draw meaningful inferences from the patient’s entire medical history. This supports advanced reasoning and personalized healthcare solutions. LTM would not only retain a patient’s ongoing medication regimen but also relate it to their historical response to similar treatments, recognize patterns in evolving symptoms, and adapt recommendations accordingly. This refined memory enables the AI model to facilitate advanced reasoning and personalized healthcare that evolves over time, delivering deeply individualized care.

In conclusion, while Raw Data provides a foundation for immediate context management, LTM is indispensable for achieving the deep, sustained personalization necessary for AI systems to remain relevant and effective. This allows AI to deliver richer, more personalized experiences that continuously improve and adapt to real-world complexities.

## 4.2 Construction Strategies of Raw Data

Unlike the human brain that seamlessly receives and processes vast arrays of sensory inputs into memory, AI models face significant challenges in utilizing raw data to build LTM effectively. Critical issues that need to be addressed include:

- **Data diversity and representativeness:** LTM must be diverse and representative to ensure model robustness across various scenarios and user groups. Imbalances and biases in collected data can lead to suboptimal performance. Additionally, the quality and consistency of data labeling are critical; inaccuracies can diminish model performance and personalization outcomes [40].
- **User behavior capture and reasoning:** To achieve true personalization, the model needs to deeply understand user behavior patterns and reasoning processes. However, current LTM models struggle to capture intermediate reasoning steps and time-series information. This

deficiency makes the model fully simulate human cognitive processes, resulting in outputs that may lack coherence and contextual relevance.

- **Data privacy and security:** The development of personalized LTM models has brought data privacy and security issues to the forefront, especially in on-device fine-tuning. Ensuring the security of data handling and processing is essential [41].

#### 4.2.1 Data collection framework

Establishing a robust data collection and labeling process is essential for ensuring diversity and high quality LTM. The collection of data for personalized models begins with understanding the personal experiences, encompassing both digital footprints and physical interactions. Given the fragmented nature of data collected compared to the continuous, rich data processed by the human brain, it is crucial to employ advanced methods for both realms. We have developed a comprehensive framework for data collection, analysis, and synthesis, which allows for differentiated system deployment tailored to various business scenarios. To validate the generalizability of this framework, we deployed independent data systems in two distinct business contexts: office collaboration and health management. A detailed description of these data systems is provided in Sec6.1.2..

In the digital realm, human activities and behavior on various platforms form significant digital footprints, providing crucial insights into preferences, interests, and interaction patterns. Digital communication data, including messages, emails, and voice recordings, reveal language preferences, communication styles, and social connections [42]. Analysis of this data is essential for developing personalized AI systems capable of tailoring responses and enhancing human interaction experiences.

Behavioral data, such as web browsing histories, app usage patterns, and social media activities, offer a detailed perspective on human interests and behavior in digital spaces [43]. This information is crucial for understanding how individuals interact with various digital environments, facilitating the development of more accurate personalized models. Content consumption data, including the types and genres of articles, videos, and music consumed, contribute to building comprehensive profiles of user preferences and trends over time. Research by Jeung and Huang emphasizes the significance of this data in facilitating tailored interactions and fostering personalized experiences [44].

In the physical realm, wearable technologies and ambient sensors provide extensive data about human interactions with the environment. Devices such as AR glasses (e.g., Apple Vision Pro) capture continuous streams of visual and auditory data, reflecting environmental interactions. These devices are crucial for gathering comprehensive data inputs that are otherwise difficult to obtain through traditional methods. Recording personal interaction logs across various devices and contexts, such as smartphone usage and smart speaker commands, provides valuable insights into behavior. This practice facilitates the integration and analysis of diverse data sources, improving the reliability and accuracy of personalized models [45]. Biometric signals, including heart rate, galvanic skin response, and brain wave activity, offer insights into human states and preferences. Biometric authentication systems based on signals such as ECG and EMG provide robust datasets for personalized modeling applications, enhancing the accuracy and reliability of AI systems [46; 47].

Contextual and behavioral data collection is crucial for a comprehensive understanding of human experiences. Such data includes location-based information and interactions within specific contexts, which provides a richer and more nuanced picture of human behavior. Tracking location data and correlating it with behavior patterns introduces essential spatial and temporal dimensions to the dataset. Hybrid location privacy schemes protect data privacy while maintaining the accuracy of services [48].

To enhance the fidelity and richness of data collected for personalized models, future research should focus on Enhanced Sensory Integration and Context-Aware Data Collection. Leveraging advanced sensor technologies to capture richer, multi-modal data streams that closely mimic human sensory experiences is vital for improving data comprehensiveness and contextual richness. Developing advanced algorithms that can infer and fill contextual gaps in recorded data provides a more comprehensive understanding of personal experiences, thereby leading to better predictions and personalization.

However, the diversity, representativeness, and quality control of the data are critical issues that must be emphasized during the collection process. High-quality personalized models require diverse and high-quality data to ensure robustness across different scenarios and user groups. Simultaneously,

data privacy and security issues are paramount, and it is essential to protect user data during collection, storage, and processing.

#### 4.2.2 Data synthesis techniques

Given the privacy concerns, limitations in individual data collection, and the fragmented nature of such data, synthetic data generation becomes increasingly critical. By producing synthetic data that closely mirrors real-world scenarios, we can compensate for the shortcomings of real data, providing rich context and diverse experiences for personalized modeling, thereby enhancing the relevance and practicality of the model.

Early approaches, such as utilizing human interactions and role-playing, were effective in producing high-quality datasets. For instance, one study [49] introduced a three-phase human role-playing paradigm to synthesize a depression diagnosis dialogue dataset. However, human-involved methods still face significant challenges, including high costs and privacy concerns [50]. Additionally, other studies [51; 52] have emphasized that human-generated data may not be optimal for LTM construction due to inherent biases.

Recent advancements in large language models (LLMs) offer a scalable and efficient alternative. Synthesizing data that mirrors real-world scenarios using LLMs can significantly enhance the pertinence and utility of LTM in personalized models. An iterative dataset synthesis method was proposed, in which large language models (LLMs) generate datasets for smaller models by progressively refining inaccuracies identified during the evaluation of the smaller models against real-world data. This approach is called Synthesis Step-by-Step (S3)[53]. It iteratively prompts the LLM to extrapolate errors encountered by smaller models in validation datasets, gradually narrowing the gap between synthetic and real data. Experimental results indicate that the S3 improves model performance across multiple natural language processing (NLP) tasks, with gains of up to 15.17%, while reducing reliance on human-annotated data. This highlights its potential in enhancing model accuracy and reducing data annotation costs.

Simulating context-rich interactions is crucial for capturing the nuances of individual experiences. The construction of multi-agent systems, where agents assume various roles within text-based role-playing environments, offers an effective approach for generating diverse interaction data. For instance, the thespian agent framework offers a novel approach for emulating multiple characters in text role-playing games, facilitating the capture of diverse experiential nuances [54]. This framework enables agents to play multiple characters using a soft-prompt mechanism, which guides the agent on which character to simulate. Through an integrated attention mechanism, agents learn new characters based on previously actions and feedbacks, enabling few-shot learning and enhancing adaptability in simulated environments. Experiments have demonstrated that this approach outperforms existing multi-character learning frameworks, making it highly effective for generating context-rich interactions that support LTM in personalized models.

Interactive simulations in dynamic environments is also essential for emulating the complexity of human experiences. The EnvGen framework employs LLMs to generate and adapt training environments tailored to the specific weaknesses of reinforcement learning (RL) agents, continuously adjusting configurations based on agent feedback [55]. This dynamic adjustment fosters the development of robust and adaptable skills, outperforming static curriculum learning and generating synthetic data that closely mimics real-world interactions.

Enhancing agent-environment interactions is another promising approach for synthetic data generation. The Affordable Generative Agents (AGA) framework proposes cost-effective interaction paradigms to optimize agent-environment and inter-agent interactions[56]. To achieve this, AGA substitutes repeated LLM inferences with learned policies for agent-environment interactions and compresses auxiliary dialogue information for inter-agent interactions. This reduces the computational costs while maintaining the quality of interactions. Extensive experiments have demonstrated the efficiency and effectiveness of AGA in producing believable low-cost interactions, making it an essential approach for scaling LTM synthesis in personalized models.

Given the inherent limitations of any single approach, a mixed data synthesis strategy offers a more robust and comprehensive solution. A holistic and flexible combined approach that incorporates multiple data synthesis techniques can yield superior results. The integration of real data synthesis with iterative error correction, role-playing adaptive dialogues, dynamic simulation environments, and



agent-based interaction frameworks collectively forms a robust methodology for LTM in personalized models. This mixed-method framework leverages the strengths of each technique while mitigating their individual limitations, thus rendering a comprehensive and rich personalized model that can process a wide array of personal experiences effectively. In 6.1.2, a mix-method data generation algorithm called: RTG synthesis method will be introduced, demonstrating its effectiveness in enhancing the model’s performance in memory recall tasks without negatively impacting its overall capabilities.

### 4.3 Construction Strategies of LTM

LTM is the effective organization and structuring of raw data, not merely classifying and sorting raw data on the surface to make the information more organized. Instead, it is about designing and optimizing the forms of LTM storage from the perspective of rapid storage and retrieval of memory and efficient utilization of information. By establishing links between related information, effectively processing the data, and reorganizing the information, an intelligent agent can quickly locate the needed memory fragments, thereby improving response speed and accuracy. The following are several major operational modalities and their corresponding transformation methods.

- **Text summarization:** Summary storage compresses long contexts into short, concise data forms. It involves inductively summarizing continuous interaction data, allowing LTM to exist in continuous and abstract forms. This method significantly reduces storage space through high abstraction and data compression, thereby enhancing retrieval efficiency. Summary storage is particularly suitable for scenarios requiring frequent retrieval and rapid understanding. For instance, OpenAI’s ChatGPT employs this method by recognizing and summarizing interaction records in real-time, thereby generating short contextual memories aligned with the user’s personal habits, which enhances its ability to meet personalized user needs during subsequent use.
- **Data structuring:** Structured storage predefines specific data structures, such as hierarchical structures, tree structures, key-value pairs, etc., to organize and store processed raw data. This approach makes information highly systematic, facilitating data retrieval and management. Moreover, it allows interaction with the stored data via database query languages like SQL, achieving more precise memory retrieval while also providing a bridge for multi-turn interactions between memory and intelligence.
- **Graph representation:** Graph-based storage uses nodes and edges to represent and manage information, boasting flexibility and scalability, which is suitable for representing the associativity of long-term memory content. Transforming raw data into a graph-based long-term memory requires data preprocessing, including data cleaning and format conversion, followed by identifying key entities and their relationships within the data to define the graph’s nodes and edges. During the graph construction stage, these nodes and edges are added to the graph structure and optimized techniques such as graph sparsification and edge weighting are applied to enhance the storage and retrieval efficiency of the graph. These optimized graphs are stored in specialized graph databases like Neo4j to achieve efficient querying and long-term memory management, thereby improving the speed and accuracy of memory response.
- **Vectorization:** Vectorized storage involves segmenting raw data into fragments and then converting them into high-dimensional vector forms, making similarity computation and rapid retrieval more convenient and efficient. Retrieval-Augmented Generation (RAG) holds an advantage in handling long-context tasks. By converting raw data into high-dimensional vectors, RAG significantly enhances the comprehensiveness and accuracy of retrieval and generation through vector matching techniques, fully leveraging its role in managing complex data tasks, thereby achieving efficient data processing.
- **Model parameterization:** Parameterized memory storage stores memory by adjusting model parameters (such as neural network weights) to realize information compression and generalization. Implementation methods include storing memory in the key-value cache of Transformer models or saving memory through parameter training in RNN-based models[57]. Additionally, some methods involve fine-tuning large models to store memory.

These approaches transform raw data into different forms of long-term memory through various technical means, enhancing their adaptability to large models or intelligent agents in practical applications. Each method exhibits its own advantages and drawbacks in different scenarios. Currently, the most widely practiced approach is RAG-based storage. Meanwhile, the field of parameterized memory storage continues to see numerous emerging research outcomes. We have experimented and implemented multiple forms and plan to continue our research along this path.

## 5 How can LTM be used to achieve model self-Evolution?

After acquiring high-quality LTM data, the next challenge we face is how to use it to enhance model capabilities and enable model self-evolution. There are several key challenges need to be addressed in the process of using LTM data to maximize its effectiveness and efficiency, including:

- **Adapting to Continuously Updated LTM Data.** As user LTM data is continuously accumulated, the model must balance learning new information while retaining previously acquired knowledge. Traditional models often assume stable data distributions, but in real scenarios, new LTM data may significantly diverge from earlier patterns, leading to risks like overfitting or catastrophic forgetting. Efficiently handling these shifts is critical to adapt to the dynamic LTM data.
- **Real-Time Learning and Efficient Feedback Integration.** Due to the LTM data being accumulated dynamically, models must quickly adapt to real-time changes in user behavior. Rapid integration of new data is vital for applications such as intelligent assistants, where seamless user interaction is key. Besides, user feedback, both implicit (e.g., clicks or time spent) and explicit, should be taken into consideration when refining the foundation model. Incorporating both types of feedback in real-time enables the model to continuously improve and align with individual user needs.
- **Handling Data Sparsity and User Diversity.** Data sparsity is a common issue in continuously updated LTM systems, especially for users with limited interaction history or sporadic activity, making it difficult to train models. Furthermore, User diversity adds further complexity, requiring the model to adapt to individual patterns while still generalizing effectively across diverse user groups.

In the following subsections, we will review existing efforts that can be adopted to solve these challenges and summarize their weaknesses. Note that LLMs have become the most powerful AI models, the model self-evolution discussion will also focus on using LTM combined with LLMs. The discussion is divided into three key sections: first, LTM as an external knowledge repository, with retrieval techniques like Retrieval-Augmented Generation (RAG) to access information dynamically (in 5.1); second, the direct training of models with LTM, embedding knowledge into the model itself (in 5.2); and third, a hybrid approach that combines external retrieval for precision with a fine-tuned model for improved efficiency in recalling relevant data (in 5.3).

### 5.1 Incorporating LTM as Outside Knowledge Bases

Since LTM is generated during interactions and may be refined with preprocessing, which is continuously updated, an intuitive strategy of utilization is to keep it independent from the foundation model, serving as an outside knowledge base to assist with reasoning and inference. While taking Transformer-based LLMs as an example, their structure inherently does not retain or recall any previous state or sequence, which limits the model’s ability to carry LTM information but only short-term “working memory” for the current context window. So extra efforts should be made to address this limitation in utilizing LTM as outside knowledge bases for LLMs.

The most popular strategy in using an outside knowledge base for LLMs is Retrieval-Augmented Generation (RAG) with In-Context Learning (ICL), which can be adapted to utilize LTM data too. ICL mimics working memory, processing short-term information. Meanwhile, RAG, serving as an external database, acts as a module for LTM.

As a comparison, the LLM’s ICL ability to plan and reason improves as model size and capabilities grow, approaching—if not exceeding—the capacity of the human prefrontal cortex in specific domains. Human working memory is limited to a span of around  $7 \pm 2$  items, whereas modern LLMs, with

increasingly large context windows (ranging from 4K to 1M tokens or more), can hold significantly more information. RAG, as an external storage mechanism, provides the potential for storing and retrieving an unlimited amount of information, which seemingly addresses the limitation of human memory by offering three analogous stages: encoding, storage, and generation.

### 5.1.1 Advances in Model Memory Encoding, Storage, and Retrieval

In this section, we first review existing efforts made in memory encoding, storage, and retrieval with RAG and ICL, which are necessary modules to utilize LTM to achieve model self-evolution. Although there are several challenges in distinct steps, techniques like adaptive memory pruning, more efficient indexing systems, and context-aware retrieval algorithms could help mitigate some of these issues. Additionally, integrating reinforcement learning-based feedback mechanisms could enable more selective memory encoding, allowing LLMs to prioritize important data and improve long-term retrieval accuracy.

- **Encoding:** One of the most critical shortcomings of current LLM architectures is the indiscriminate nature of data encoding. Unlike human memory, which selects relevant information for long-term retention based on salience and context, LLMs encode vast amounts of information without prioritization, leading to bloated and inefficient memory systems. Recent advancements in RAG-based systems have focused on introducing mechanisms for selective encoding, similar to human memory. Models like MuRAG [58] and RA-CM3 [59], for instance, are integrating multimodal data (images, text) with refined hierarchical indexing systems. This approach mirrors how humans encode experiences from different sensory inputs, allowing the system to distinguish between relevant and irrelevant information, ensuring that only the most contextually pertinent data is encoded into memory. Additionally, recent techniques such as UPRISE [60] offer promising methods to enhance memory encoding by constructing a pool of context-aware prompts, selectively retrieving relevant templates for different tasks. This dynamic process parallels how humans draw on past experiences when encoding new information, improving the efficiency and selectivity of memory formation.
- **Consolidation:** While RAG architectures allow for unlimited data storage, they still struggle with consolidation, organization, and fragmentation of memories, which become more difficult to retrieve as the memory system grows. In contrast, human brains consolidate memories during rest periods, integrating and organizing information in a way that strengthens important memories while letting others fade. To address this, recent studies have developed more structured storage systems, such as RAPTOR [61] and MEMWALKER [62], which take a step towards improving the hierarchical organization of stored information. These models introduce a tree-like structure that clusters and summarizes data recursively, creating a more structured memory space that allows for efficient retrieval of information across different levels of abstraction.  
In addition, emerging models like Memory<sup>3</sup> [63] propose integrating explicit memory modules that manage knowledge as key-value pairs in external non-volatile storage, emulating how the human brain efficiently consolidates and accesses memories based on relevance and frequency of recall. Techniques like compressive memory, seen in Beyond Retrieval: Embracing Compressive Memory in Real-World Long-Term Conversations [64], explore how to compress and store interactions while retaining the core elements of past experiences, thus reducing redundancy and improving access to the most relevant parts of memory during retrieval.
- **Retrieval:** The retrieval process is where the largest divergence between human and LLM-based memory systems occurs. While humans rely on associative and context-driven recall, current LLM retrieval relies heavily on probabilistic token generation, often resulting in the retrieval of incomplete or irrelevant information. Recent work has shifted toward improving retrieval systems by introducing hybrid search techniques that blend traditional keyword-based search with more advanced semantic retrieval systems. For instance, Blended RAG [65] and Hybrid Search [66] combine BM25 with dense vector embeddings to improve both the relevance and context-awareness of retrieved information. Moreover, methods like FiD emphasize the need for creating more sophisticated queries that better match the intent and structure of the memory database, drawing parallels to how humans refine their recall based on the specific context of a query.

Besides, post-retrieval mechanisms such as re-ranking and self-reflection, as seen in models like Self-RAG [67] and Re2G[68], aim to refine the retrieved information to align more closely with the user’s query. These systems not only filter out irrelevant content but also prioritize the most contextually appropriate responses, echoing how human memory retrieval is shaped by the relevance and importance of past experiences to the current situation.

To summarize, the RAG and ICL strategies do not require updating the foundation model’s parameters, allowing for efficient integration and updating of LTM data, which makes it possible to support real-time, efficient LTM management and utilization. The foundation model can evolve with distinct LTM data through this strategy. However, despite improvements in memory encoding and retrieval, there is still a fundamental disconnect between the deep, associative nature of human memory and the more surface-level pattern recognition performed by LLMs. New models like MemoCRS [69] and COMEDY [64] represent initial attempts to bridge this gap by integrating deeper, long-term contextual memory into conversational AI, but much work remains in improving the depth and dynamism of memory representation.

## 5.2 Incorporating LTM by Updating Model Parameters

Another type of strategy for model self-evolution using LTM is through a parameterized approach, where the accumulated LTM data is transformed into the model’s internal parameter memory through model training. The largest difference between this method and the LTM data utilization in RAG and ICL is that it directly modifies the parameters of the foundation model. As a result, it requires a larger data scale, higher implementation costs, and lower update efficiency compared to the former. However, its advantage lies in not needing an explicit LTM management mechanism design, and it can also be completed through single-step inference during the inference.

Regardless of a foundation model’s architecture, the LLM training paradigm is typically classified into three phases: pre-training, instruction tuning, and alignment tuning [70]. Each phase serves a distinct purpose and many state-of-the-art LLMs [71; 72; 73] undergo multiple iterations (and often interleaving) of these phases to achieve optimal results. In our context, each phase presents a unique opportunity for LTM data integration.

### 5.2.1 Brief Review of LLM training

We first briefly review each of the LLM training phases and discuss how they can be utilized to encode LTM data for model self-evolution.

**Pre-training.** LLMs undergo the pre-training phase in an unsupervised manner and typically using diverse, large-scale corpora. This is essential for LLMs because it can equip them with extensive world knowledge and a deep understanding of language [24; 70]. The next-token prediction task is commonly used to train decoder-only models, such as Llama [74] and GPT [24], which aims to predict a target token  $x_i$  given the preceding sequence in an auto-regressive manner. Concretely, we would like to maximize:

$$\mathcal{L} = \sum_{i=1}^n \log P(x_i | x_1, \dots, x_{i-1}). \quad (1)$$

Pre-training typically requires significant computational resources and large-scale data support, while LTM data often faces data sparsity issues. So a candidate method known as continued pre-training can be employed to encode LTM into models. By using a similar training objective, a pre-trained LLM can continue learn new data about specific tasks [75], such as recent information [76; 77] or domain-specific knowledge [78; 79; 80]. Thus, this strategy also works for encoding LTM data.

Despite continued pre-training enabling LLMs to acquire new knowledge and is particularly useful for integrating LTM data, there is a potential risk of catastrophic forgetting [75]. Research shows that appropriate learning rate management [81; 82] and careful adjustment of data distribution [83] to mitigate distribution shift play a critical role in a successful continued pre-training process. Additionally, techniques which enable updating only part of the model’s parameters can be helpful to minimize the impact on the original knowledge [84].

**Instruction tuning.** Instruction tuning is usually performed after pre-training to further enhance their capabilities. In this phase LLMs are fine-tuned with instruction instances in a supervised manner. These instances typically consist of a task description ( $T$ ), along with input ( $X$ ) and desired

output ( $Y$ ) pairs. LTM data can also be reorganized in this format and be encoded into LLMs. For LLM-orientated training, the instructions (or prompts) ( $I$ ) are usually constructed as  $I = T + X$  [85]. Therefore the training objective of instruction tuning (in the context of LLM, this is also called Supervised Fine-Tuning (SFT) [86]) that can be formalised as:

$$P(Y_1, \dots, Y_n | I_1, \dots, I_m) = \prod_{j=1}^n P(Y_j | I_1, \dots, I_m, Y_1, \dots, Y_{j-1}), \quad (2)$$

where  $I = \{I_1, \dots, I_m\}$  and  $Y = \{Y_1, \dots, Y_n\}$  are the instruction and output sequence respectively. In essence, the model is trained to predict each token in  $Y$  given  $I$  and  $Y$  so far. A common loss function for this training objective is [87]:

$$\mathcal{L} = -\log P(Y_1, \dots, Y_n | I_1, \dots, I_m) = -\sum_{j=1}^n \log P(Y_j | I_1, \dots, I_m, Y_1, \dots, Y_{j-1}). \quad (3)$$

For LTM data, this phase presents a prime opportunity for memory integration. Unlike pre-training, only a moderate amount of data (i.e. pairs of  $(I, Y)$ ) is required in this phase to significantly influence the model, leading to general performance improvements [70]. As a direct comparison, humans receive a vast amount of real-time events, and they filter these signals to distill valuable information. Only a small portion of this information (i.e.  $I$ ) is stored in human’s LTM through repetitive activation and memorization, which can be retrieved and output (i.e.  $Y$ ) later for downstream tasks. To mimic this process, we must curate a refined set of instruction data and train the model iteratively on this smaller set. The quality of this refined information (or data) is crucial in influencing both human and LLM decision-making and performance. Moreover, instruction tuning can also enables LLMs to generalize to unseen tasks [88]. Additionally, like pre-training, domain-specific knowledge can also be infused into the model, turning it into a domain expert. These qualities are highly desirable for effective LTM data utilization, as the final system should be capable of performing tailored and complex actions with self-evolved model. Understanding what data accurately reflects LTM is crucial when incorporating them into LLMs. The answer varies depending on the application. For example, in an emotional support chatbot, LTM includes past sessions between users and the chatbot [89]. In a medical assistant scenario, LTM involves patients’ past medical records [90].

**Alignment tuning.** Despite extensive data cleaning, the large volume of pre-training data may still vary in quality, potentially leading to undesired responses [91]. To address this, several techniques have been proposed and implemented after the model’s fine-tuning process. These techniques, collectively referred to as alignment tuning, aim to align the model with human preferences. RLHF [86] and RLAIIF [92] are two popular alignment methods, each with its own use case and limitations. RLHF tunes a reward model that rates different outputs based on human feedback, which is then used to further fine-tune the LLM [93]. In contrast, RLAIIF directly links the LLM to a larger, more aligned model for learning. Due to the complexity and overall performance of RLHF, other methods such as Direct Preference Optimization (DPO) [94], Knowledge Transfer Optimization (KTO) [95], Optimized Reward Preference Optimization (ORPO) [96], and Simplified Preference Optimization (SimPO) [97] are often employed in different scenarios for alignment tuning. However, although it is technically feasible to adjust model preferences for specific applications or even incorporate domain knowledge in reward models to impart biases [98; 99], curating a high-quality preference dataset that aligns with the desired LTM remains challenging.

In summary, model parameterization occurs across all three phases of LLM training, which theoretically allows for LTM integration at any stage of model self-evolution. However, due to constraints such as data availability, training stability, and computational resources, we believe that supervised fine-tuning (SFT) remains the predominant technique for direct LTM data integration. Achieving true model self-evolution, however, may require a more sophisticated approach, particularly in obtaining high-quality feedback signals for effective self-guidance.

### 5.2.2 Advantage of Model Parameterization

Compared to the LTM integration strategies based on RAG and ICL, parameterized encoding method has significant advantages in terms of input format and inference efficiency. As most current LTM data are in the textual form due to its clearer interpretability, better context control, and ease of implementation [100], supporting textual LTM requires LLMs to have a sufficiently large context

window to incorporate all necessary information during prompt construction. However, decoder-based LLMs are typically constrained by a fixed, limited context window. To address this, significant efforts have been made to extend the context window effectively [101; 102; 103; 104; 105; 106; 107] or to process information in parallel [108; 109]. Despite these advancements, the information stored in LTM can still exceed the capacity of state-of-the-art context windows. Additionally, due to the architecture of transformers [110], which serve as the foundation for most LLMs, computational cost increases quadratically with input length. This makes extremely long inputs less practical for some applications. A limited context window necessitates the use of a retriever to choose relevant information from a large, manually curated LTM resource, such as a database or vector store [111]. This shifts the challenge from the LLM to the design of an effective retriever. Injecting LTM data directly into the model addresses this limitation, as it eliminates the need to include such information in the input prompt. This allows the LLM to focus on the task itself, with memory retrieval integrated into the inference process. This integration leverages the LLM’s advanced language understanding and reasoning capabilities, avoiding dependence on a (usually) less capable external retriever.

If we choose to use training methods to encode LTM data into LLMs, a vital question is how to decide the data structures and types for training. As we known, in LLM-based agents, short-term memory manages contextual information, while long-term memory stores past experiences, reflections, and profiles. Profiles are arguably one of the most critical components of an agent, directly influencing its functions and interactions [112]. Typically, profiles contain basic details such as age, gender, and occupation [113], and may also include traits, interests, and behavioral patterns [114]. While it is possible to fit such descriptions into a prompt, this approach only portrays an average person playing a specific role, lacking the depth needed to reflect an individual’s rich history and experiences [115]. To integrate such data into the model in a unified way, one can use SFT with specially curated datasets. In Character-LLM [115], training data is constructed not only from profiles but also from experiences and scenarios, enabling LLMs to display more human-like characteristics. Similar works [116; 117] also fine-tune models using character-specific data derived from past memories and experiences, proving effective in aligning models more closely with their intended personas.

In addition to in-character memory, one of the most crucial aspects of LTM is its associated domain knowledge. It has been well established that LLMs can be significantly enhanced in specific domains through fine-tuning, enabling smaller models to potentially surpass larger models. This is especially evident in domains that require extensive knowledge for complex reasoning, such as healthcare and finance [70]. For instance, the Med-PaLM series [118; 119], fine-tuned from the PaLM model [120], became the first model to pass the U.S. Medical Licensing Examination (USMLE), demonstrating human-expert level proficiency in the medical domain. In legal services, Legal Artificial Intelligence (LegalAI) [121] has been continuously evolving and has reached new heights with the advent of LLMs. Fine-tuned models can now perform in-depth legal question-answering, simulating the role of a lawyer [122; 123], and efforts are underway to develop intelligent legal service systems [124]. In finance, domain-specific fine-tuning has also shown significant promise, as seen in models like BloombergGPT [125] and FinGPT [126]. Domain knowledge plays a vital role in model self-evolution, and these examples demonstrate the feasibility of integrating specialized knowledge, which elevates models beyond their general-purpose capabilities.

### 5.3 Incorporating LTM Data with Mixed Strategy

Since using LTM data solely as an external knowledge base or as additional training data each has its own advantages and limitations, a simple question is whether we can integrate these two methods (RAG with ICL and SFT) to leverage their strengths and achieve more effective model self-evolution. Recent works on RAG also see an increase in fine-tuning models at each stage of RAG: retrieval, augmentation, and generation. By focusing on fine-tuning, we can significantly enhance an RAG system’s efficiency and effectiveness in handling LTM data.

#### 5.3.1 Fine-Tuning in the Retrieval Stage

The retrieval stage in RAG systems involves selecting relevant documents or knowledge chunks from external sources to boost the model’s intrinsic knowledge. Fine-tuning the retriever model at this stage optimizes its ability to retrieve more relevant, domain-specific information.

By fine-tuning the retriever with LTM data, the precision of retrieving highly relevant content improves. In fields such as law or medicine, this ensures that the retriever pulls accurate and

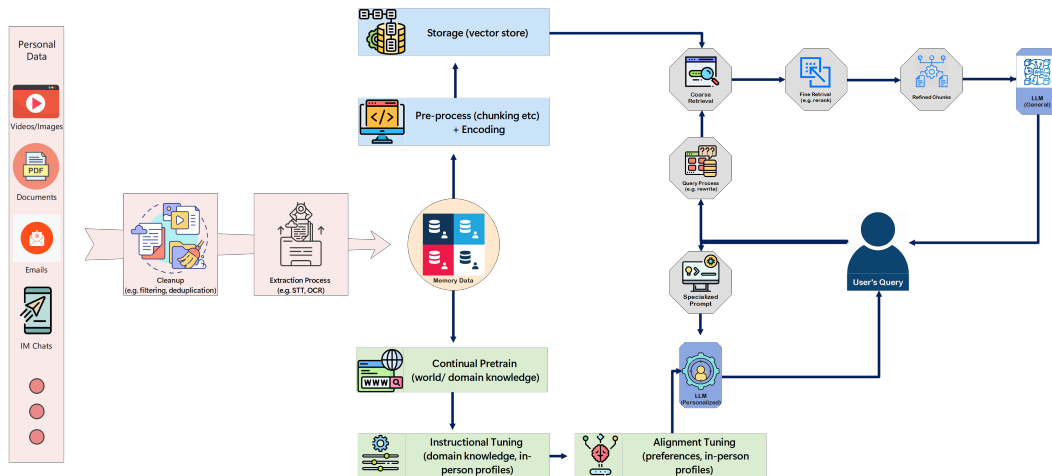


Figure 5: How personal data can be utilised through textual and parametric memory retrieval

specialized content, increasing relevance and reducing noise during the retrieval process. This addresses the limitations of generic models, which may retrieve unrelated or less useful data [127]. This process is often applied to embedding models to align the representations of both the query and external documents with the domain’s vocabulary and semantics. This alignment enhances retrieval effectiveness by improving semantic similarity matching. Moreover, fine-tuning has been used to align the retriever with the generator. One such method is the LM-Supervised Retriever (LSR), where the output of LLMs is used to train the retriever in a feedback loop [127]. For instance, REPLUG LSR [128] trains the retriever to select documents that can reduce the perplexity of the LLM’s output sequences.

### 5.3.2 Fine-Tuning in the Augmentation Stage

The augmentation stage in RAG systems integrates retrieved LTM data with the model’s existing knowledge to create a comprehensive context for generating an informed response. Fine-tuning can optimize the integration and selection of relevant information at this stage. Since the augmentation stage acts as an intermediary between retrieval and generation, fine-tuning can be applied to the retriever and/or the generator (or as an additional module in between).

For instance, fine-tuning enhances contextual filtering and selection by eliminating redundant or irrelevant information from the retrieved documents. A notable example is Self-RAG [67], which uses reflection tokens to indicate whether retrieved documents are relevant to the query, whether the generated responses are supported by the documents, or whether the response effectively addresses the query. This approach significantly improves and streamlines the augmentation process. Additionally, models can optimize their augmentation strategies through document grounding, particularly when handling large volumes of documents. Fine-tuned models are better equipped to identify relevant and critical information from the retrieved content, resulting in more accurate response generation and reducing the likelihood of hallucination [129].

### 5.3.3 Fine-Tuning in the Generation Stage

The generation stage is where LLM, enhanced by the retrieved context, produces the final response. Fine-tuning at this stage improves the quality of the generated text and ensures alignment with the desired style, tone, and factual accuracy. For instance, fine-tuning enables the LLM to conform to specific formats and enforce a particular response style. This technique is especially useful for equipping LLMs with function-calling capabilities and tool usage [130].

Fine-tuning also allows the model to be customized for specific tasks or domains, ensuring that the generated content adheres to domain-specific norms and conventions. For example, in medical applications, fine-tuning helps the model generate clinically appropriate responses, while in legal tasks, it ensures adherence to legal terminology and reasoning. Additionally, the Chain-of-Thought

process [131] can be integrated into model generation, guiding the model through complex reasoning and facilitating connections between relevant knowledge in the retrieved documents [132].

In summary, the integrated LTM utilization method combines the strengths of both strategies, optimizing data use and supporting model self-evolution. However, practical challenges remain, such as determining the optimal frequency of SFT to balance efficiency and effectiveness.

#### 5.4 Utilizing LTM with Multi-Agents

Except for the primary encoding, storage, and retrieval studies about how to utilize LTM in specific LLMs, we propose that LTM can also be adopted by using RAG-based multi-agent systems to achieve model self-evolution. Several interconnected strategies can be employed to create a cohesive and efficient memory framework that mimics human cognition, including 1) dynamic memory storage, 2) context-based retrieval, 3) memory consolidation, and 4) self-reflection, all within a multi-agent framework. In this framework, different agents can understand LTM patterns of user behavior and improve their decision-making and action planning, and we will introduce them one by one. Note that each module can be implemented by a single LLM-powered agent.

1. **Dynamic Memory Storage** is a key component in improving long-term memory is the dynamic management of stored information. Instead of passively accumulating data, memory systems must prioritize the retention of relevant information. This involves continuously updating the memory bank based on interaction context, user preferences, and relevance to the task. By applying a vectorized storage system, memories can be efficiently retrieved and updated, ensuring that the most pertinent information remains easily accessible. This approach mirrors the brain’s ability to filter and prioritize experiences, focusing on valuable memories while allowing less critical information to fade.
2. **Context-Based Retrieval** allows memory systems to dynamically recall relevant information based on current needs. Rather than relying solely on explicit keywords, retrieval can be enhanced by leveraging semantic similarity, temporal context, and retrieval frequency. This method improves both the precision and relevance of the information recalled, much like how humans retrieve memories by associating past experiences with present situations.
3. **Memory consolidation** is also crucial for maintaining LTM for further utilization. Systems must periodically assess the relevance and usage of stored memories, reinforcing essential data through consolidation while gradually deprioritizing less frequently accessed information. This helps avoid issues such as catastrophic forgetting by ensuring critical memories remain intact while the system avoids becoming cluttered with outdated or irrelevant data.
4. **Self-reflection** enables models to refine their LTM by periodically reviewing stored information. Through this reflective process, general principles can be abstracted from specific experiences, leading to more efficient memory management and retrieval strategies in the future. This self-reflective capacity not only optimizes memory storage but also enhances decision-making and adaptability over time.

Comparing with previous LTM utilization strategies, the proposed multi-agent based LTM management framework are more general to corporate with distinct LLMs. The mentioned four modules/functions can be distributed across different agents with specific tasks. By analyzing user behavior patterns over time, agents can better plan actions and make more informed decisions, resulting in a more adaptive and personalized experience. The integration of these memory-focused processes allows agents to operate cohesively, leveraging their shared knowledge to create a more efficient and human-like system for managing and utilizing long-term memory.

#### 5.5 Utilizing LTM to Support Agent Self-Evolution

In previous subsections, we have reviewed the existing research and the multi-agent integration method we proposed to use RAG and ICL for model self-evolution with LTM. Furthermore, take the Agent Hospital study conducted by the Tsinghua University team as an example [133], we will show how the use of LTM helps improve the model’s self-capabilities in a specific simulacrum medical scenario. The proposed LTM accumulated and utilization method MedAgent-Zero can be mainly divided into three modules: medical records accumulation, medical experience reflection, and RAG-based LTM utilization.



1. **Medical Records Accumulation.** In this work, the doctor agent first accumulates LTM from successful medical cases to handle future problems. Taking disease diagnosis as an example, after learning the symptoms of a virtual patient, if the doctor agent provides a correct response, this diagnostic process will be saved as a medical record, just like a real-world doctor does, forming text-based LTM data. It is worth noting that even in the real world, the accumulation of medical records is very beneficial for enhancing a doctor’s medical capabilities, as it can provide strong decision support when encountering similar issues in the future.
2. **Medical Experience Reflection.** In the second module, we focus on how the doctor agent accumulates LTM through self-reflection from failed medical cases. Therefore, unlike storing successful medical cases in the memory base directly, we use self-reflection techniques in this module to allow the doctor agent to reflect and summarize after answering the diagnostic problem. The original question, its own response, and the correct response are adopted to generate experience-based LTM data. These data are also text-based, and to ensure their accuracy, we verify if the doctor agent can correctly answer the original question after utilizing the summarized experience (LTM will be dropped once cannot correctly answer). In other words, we designed the construction method for LTM data based on the agent’s interaction with different outcomes.
3. **RAG-based LTM Utilization.** Based on the two types of accumulated LTM data mentioned above, we propose using RAG technology to first retrieve stored cases and experiences when addressing new problems. It is important to note that the retrieval is based on the questions that generated the LTM data. Besides, since no parameter updates are required throughout the process, experience accumulation and application can be conducted online by ICL, making efficient use of the LTM data. As the doctor agent continually acquires valid LTM data during inference, its ability to handle medical questions also improves. Experimental results show that after self-evolving through diagnosing tens of thousands of patients, the doctor agent’s diagnostic capabilities in the respiratory department even surpassed the current SOTA model trained on real-world data.

To summarize, similar to the accumulation and inheritance of human knowledge that can help the development of culture, with appropriate algorithmic support, LTM can also effectively help models in achieving self-capability enhancement, which is a vital step in model self-evolution. Although the handling and application types of LTM in this work still have limitations, it fully demonstrates the potential of this direction. In Section 7, we will provide future optimization directions.

## 6 The Practice of model self-evolution based on LTM

### 6.1 Practice on LTM Data Acquisition

In our quest to enhance LTM data retention and accessibility, we have embarked on a comprehensive investigation of various methodologies. This exploration encompasses a spectrum of approaches, each designed to augment the robustness and reliability of LTM. Firstly, we have delved into the realm of real-world data collection. This involves the systematic gathering of data from actual medical environments and scenarios, ensuring that the information is as authentic and representative as possible. Secondly, we have experimented with real synthetic data generation. This innovative approach [134] amalgamates actual data with artificially generated data to create a more comprehensive dataset. These synthetic data are crafted to fill in gaps where real data may be scarce or insufficient, thereby providing a more balanced and extensive pool of information for analysis and learning. Lastly, we have ventured into fully synthetic data generation. This method relies entirely on artificially created data with techniques like Chain-of-Thought[135], which can be tailored to meet specific requirements or simulate particular conditions that may be difficult to replicate in real-world scenarios. The flexibility of synthetic data allows for exploring a wide array of possibilities, pushing the boundaries of what can be achieved with LTM data. Through these diverse tryouts, we strive to unlock the full potential of LTM data, paving the way for more efficient and effective model evolution.

### 6.1.1 Real-world LTM Data Collection

Developing robust data collection and labeling protocols is crucial for ensuring diversity and high quality in long-term memory (LTM) datasets, especially in specialized fields like healthcare and mental health. To meet this need, we collaborated with the Tianqiao and Chrissy Chen Institute (TCCI) and Shanghai Mental Health Center (SMHC) to create a sophisticated data collection framework. This framework is used in psychological consultations at two SMHC branches in Xuhui and Minhang districts, collecting real-world outpatient data between psychiatrists and incoming patients.

Our framework design prioritizes data diversity and quality. The collection process combines oral (recorded consultations) and written (consent forms, clinical notes) interactions, resulting in a comprehensive and heterogeneous dataset essential for building strong LTM. Adherence to standard operating procedures (SOPs) during patient interactions ensures consistency and quality across the dataset. Clinicians follow specific protocols, such as capturing vital signs and recording spontaneous speech, to maintain consistency across sessions.

Ethical and privacy considerations represent foundational pillars of this process. Informed consent was meticulously secured, with nurses providing detailed explanations to ensure that patients understood the recording procedures. Written consents were securely tracked and managed. Post-collection, transcription services (e.g., iFlytek) were employed to convert audio data into text, followed by a rigorous anonymization process overseen by research assistants to safeguard privacy and ensure ethical compliance.

In addition, a well-structured operational framework with clearly delineated roles was implemented to enhance both accountability and efficiency. Nurses introduced the project and gathered consent, while research assistants supervised the process. Clinicians executed the SOP during consultations. A coordinated protocol for data management and device usage ensured timely collection, transfer, and calibration of both data and equipment, with responsibilities distributed between research assistants and physicians. This refined approach integrates diverse data sources while maintaining high-quality standards, which is crucial for collecting LTM data and developing robust AI systems in specialized fields like mental health.

Specifically, the LTM dataset collection encompassed over 1,000 participants from SMHC, accumulating more than 30,000 minutes of high-quality doctor-patient audio recordings. Participants ranged in age from 12 to 80 years old, were required to be fluent in Mandarin Chinese, capable of providing informed consent, and free from significant physical illnesses. Following all consultations, professional physicians conducted diagnoses using ICD-10 and DSM-5, meticulously documenting the present illness history, chief complaints, and prescribed medications, among other details. As previously mentioned, all recordings and associated metadata underwent careful anonymization to protect participant privacy.

From these recordings, we filtered out the participants who came for the follow-up consultation and verified 1,160 sample data points to create an LTM dataset in psychiatry (referred to as the SMHC dataset). The dataset comprises 553 participants diagnosed with Major Depressive Disorder (MDD), 426 diagnosed with Anxiety Disorder (AD), and 181 individuals classified as "Other" (neither MDD nor AD), with a total number of 25 different diagnoses. Diagnoses were made using the Chinese version of ICD-10. The demographic characteristics of the sample are presented in Table 1.

Table 1: Demographics of all participants in SMHC study.

	<b>MDD (N = 553)</b>	<b>AD (N = 426)</b>	<b>Others (N = 181)</b>
<b>Age (years)</b>	29.2 ± 9.95	34.04 ± 12.03	27.46 ± 11.11
<b>Gender (%)</b>			
- Male	176 (31.8)	138 (32.4)	77 (42.5)
- Female	377 (68.2)	288 (67.6)	104 (57.5)

The applicability of this thorough, ethically compliant, and technologically integrated workflow extends beyond health management. For instance, in diverse business contexts such as office collaboration, deploying similarly tailored data systems has shown equally promising results. This approach not only ensures data diversity and quality but also addresses ethical considerations and operational efficiency, ultimately laying a strong foundation for constructing personalized LTM data.

### 6.1.2 Synthetic LTM Data Acquisition

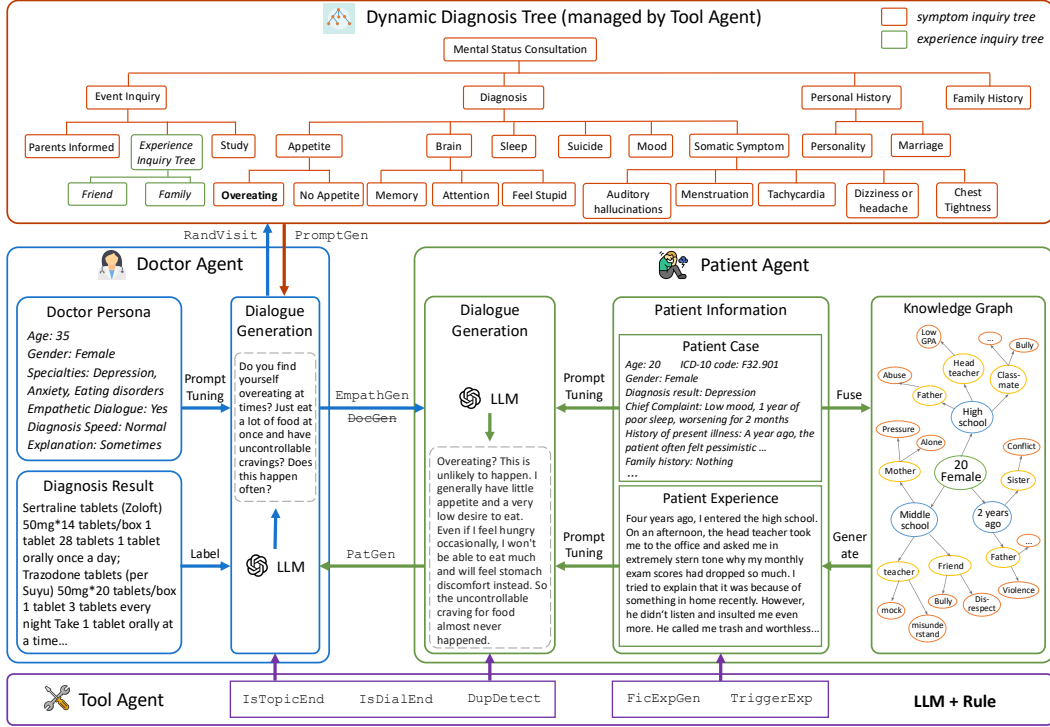


Figure 6: The neuro-symbolic multi-agent framework for synthesizing diagnostic conversation of mental disorders.

### Real Data Enhanced Synthetic LTM Generation

While real-world data collection produces high-quality data, it requires significant time, effort, and other resources. To alleviate this challenge, real-data-enhanced synthetic data generation can be a viable option. Notably, the data collected as described in is promising to boost the AI mental healthcare community. However, it cannot be directly used publicly due to stringent privacy and ethical considerations. To address this issue, we propose synthesizing diagnostic conversations by utilizing anonymized patient cases collected in this study that are more accessible. We introduce a neuro-symbolic multi-agent framework that takes patient cases as input to generate synthetic diagnostic conversations of mental disorders. As illustrated in Figure 6, this framework involves three types of large language model agents: a doctor agent, a patient agent, and a symbolic tool agent responsible for managing diagnostic topic shifts. This framework features two major innovations:

(i) One-to-many patient case-to-dialogue generation that maximizes the utilization of precious real patient cases. Unlike previous studies [136; 137] that generate one conversation with one patient case. Our proposed framework is capable of generating multiple diverse diagnostic conversations with one single patient case. Specifically, three methods ensure the diversity and correctness of the diagnostic process.

- First, doctor agents with different diagnosis habits are designed and randomly selected for each conversation.
- Second, we use LLM with the knowledge graph to generate multiple fictitious patient experiences based on one patient case. The patient experiences serve as background information for patient agents during generation. Since the diagnosis of mental disorders mainly relies on symptoms rather than concrete events, integrating the fictitious patient experiences enhances the diversity of synthesized conversation while maintaining the accuracy of the diagnostic process.
- Third, the sequence of diagnostic topics is randomly determined for each conversation.

Dataset	Professionalism	Communication (i)	Communication (ii)	Fluency (i)	Fluency (ii)	Similarity	Safety
D <sup>4</sup>	6.6	7.9	7.8	<b>8.6</b>	<b>8.2</b>	7.2	<b>0</b>
CPsyCounD	5.2	5.4	5.6	8.4	8.0	4.4	<b>0</b>
Role-playing	6.8	6.6	7.2	6.9	5.5	6.4	<b>0</b>
MDD-5k	<b>8.6</b>	<b>8.3</b>	<b>8.4</b>	<b>8.6</b>	7.6	<b>8.8</b>	<b>0</b>

Table 2: Human evaluation of different datasets.

(ii) Another significant innovation lies in text generation under symbolic control via a dynamic diagnosis tree. This tree consists of a fixed symptom inquiry tree and a dynamic experience inquiry tree. Current clinical diagnosis of mental disorders strictly follows standards from ICD-11 [138] or DSM-5 [139]. To simulate this process, we design the fixed symptom inquiry tree based on Structured Clinical Interview for DSM-5 (SCID-5) [140], covering all the diagnostic topics for important symptoms inquiry. The experience inquiry tree is constructed by extracting possible topics from the patient’s response to past experiences. It’s designed to establish deeper engagements with the patient.

By applying the proposed framework and leveraging the SMHC real data, we further utilized LLM to generate synthesized data and release the largest Chinese **Mental Disorder Diagnosis** dataset MDD-5k. It is also the first labeled mental disorders diagnostic conversation dataset with diagnosis results from professional psychiatrists.

MDD-5k contains 5000 high-quality diagnostic conversations, averaging 26.8 turns and 6906.8 Chinese words per dialogue. It is built upon the SMHC dataset mentioned in Section 6.1.1, covering more than 25 different diseases (with a majority of MDD and AD). All cases have been cleaned and filtered by global standards to ensure the complete protection of patient private information. To reinforce the clinical setting, professional psychiatrists and psychotherapists supervise the whole process and filter out unqualified dialogues. In addition, they provide diagnosis summaries based on the history of the portrait and dialogue. As shown in Table 2, the synthesized MDD-5k dataset outperforms other diagnostic conversation datasets in human evaluation.

### COT Enhanced Synthetic LTM Generation

Although we can acquire LTM data by collecting personalized interaction histories and human feedback from real-world scenarios, this strategy is often limited by high collection costs and low efficiency, making it difficult to achieve large-scale data collection quickly. In particular, real-world data often lacks key steps in the reasoning process. Therefore, we propose using LLM generation to supplement the reasoning process, thereby constructing more comprehensive LTM data for the study of model self-evolution.

To address the issue of missing intermediate reasoning steps in the thought process, we introduce a technique called Retrieval-Thinking-Generation (RTG) for LTM data synthesis, designed to enhance memory recall. The core idea is to leverage Chain-of-Thought (CoT) reasoning to strengthen Retrieval-Augmented Generation (RAG) capabilities.

The generated LTM data should contribute to training model, especially for RAG-based strategies. Therefore, the generated data needs to include not only the LTM positive data that should be utilized but also negative samples (distractors) that should not be used for RAG models. The key to constructing helpful training data is to construct high-quality distractors, so our designed RTG complete process is as follows, and an example of an RTG prompt is provided in Appendix A.

- First, we can select high-quality question-answer pairs from existing datasets as the foundational model input. These sources can include Wikipedia-based Q&A, web searches, or domain-specific datasets.
- Second, we focus on constructing high-quality distractors for the existing Q&A pairs. These two types of data should have a high degree of similarity (to train the model’s discriminative capability), but the distractors should not be able to answer the question. We can complete the selection process by using a context similarity-based strategy to filter related paragraphs from the original corpus used to construct the Q&A pairs and then verify them.
- Finally, based on the Q&A pairs along with distractors (which can range from several to dozens), we will use CoT (Chain of Thought) Reasoning to divide the data utilization

into several explicit reasoning processes, including: 1) Determining whether the current paragraph is related to the question; 2) If related, does it help answer the question? 3) Combining all paragraphs that can help answer the question to generate a response and list all cited documents. Such a reasoning and utilization process will help improve the use details of LTM data.

In addition, since LLMs may generate hallucinations, we ensure the constructed data meets correctness through the following validations: 1) The index of the ground-truth paragraph should be cited; 2) The paragraphs marked as cited should appear in the citation list; 3) All cited documents should be in the given candidate set.

To verify the effectiveness of the RTG strategy, we constructed training data for LTM utilization based on this approach and annotated a small-scale test dataset to evaluate its performance.

### Training Data Generation

To validate LLMs’s memory retrieval capability, we constructed a custom training dataset based on our RTG method. The two primary sources for this dataset are the SQuADv2 dataset [141] and WebGLM[142], which focus on Wikipedia and web data, respectively. A key distinction between our dataset and the original datasets is the introduction of a curated list of distractors, built on top of the provided context or references.

For Wikipedia data, we calculated embeddings of the context and retrieved the most similar paragraphs from Wikipedia<sup>2</sup> to generate the distractors. For web-based data, we used SerpApi<sup>3</sup> to search for relevant questions and built distractors from the most similar content chunks. Notably, we applied the Levenshtein distance to filter out paragraphs or chunks that are highly similar to the original context of reference. We utilized the GPT-4o API to perform reasoning and generate results with correct citations.

During data construction, we observed that not all questions in the original dataset are self-contained, meaning additional context is required for the question to be fully understood. For example, in SQuADv2, the question "In what R&B group was she the lead singer?" is ambiguous, as "she" should refer to Beyonce. Such ambiguity can potentially lead to the generation of random distractors, thereby lowering the overall quality of the dataset. To address this issue, techniques such as coreference resolution [143] were employed to filter out incomplete questions by identifying unresolved coreferences.

Finally, we curated a training dataset comprising 40,000 data instances. Each instance contains a question, 20 passages (including an average of 1-3 gold passages), the document ID of the relevant passage, a CoT process with citations and quotes, and the answers. The ratio of Wikipedia to web-based data is approximately equal.

### Evaluation Dataset Construction

To address the limitations of existing Q&A evaluation datasets, where candidate documents vary in length, lack targeted design for distractors, and use single-dimensional metrics, making them unsuitable for evaluating the application capabilities of LTM data after training with the generated dataset, we have constructed a dataset called LTM-COT-1.

LTM-COT-1 consists of 158 instances with high-quality preprocessing and annotations. Each instance is composed of two parts: (1) a question-answer pair with corresponding contexts and supporting evidence, and (2) 19 distractor passages. The constructing strategy is similar to the training set.

The dataset is derived from the SQuADv2 validation set [141], where five annotators provided answers to the questions. To minimize potential ambiguity, we selected only those instances where the annotators’ responses were fully consistent. To address edge cases where distractor passages might inadvertently contain the correct answer, we employed the APIs from OpenAI<sup>4</sup> and Anthropic<sup>5</sup> to conduct verification. Finally, five human annotators conducted an additional review to confirm data quality and eliminate any inconsistencies.

<sup>2</sup>Cohere Multilingual Embeddings for Wikipedia in 300+ Languages: <https://huggingface.co/datasets/Cohere/wikipedia-2023-11-embed-multilingual-v3>

<sup>3</sup><https://serpapi.com/>

<sup>4</sup><https://openai.com/>

<sup>5</sup><https://claude.ai/>

## 6.2 Practice on LTM Data Utilization

### 6.2.1 LTM Utilization with SFT and RAG

We used the private LTM-COT-1 dataset along with several publicly available RAG-related datasets to test whether the model demonstrates enhanced data utilization capability after training with the constructed LTM data. We selected the Llama-3-70B-instruct version as the backbone model and performed SFT training using the generated LTM data, the trained model is named as Homer-70B.

#### Experimental results on LTM-COT-1

We use two metrics to evaluate model performance here: 1) Answer Accuracy, calculated by dividing the number of correct answers by the total number of questions, where GPT-4 is adopted to judge the correctness. 2) Cite Score, indicating the performance of selecting citation paragraphs in terms of both accuracy and recall. Each correct earns one point and incorrect reduces one point. The final score is obtained by dividing the number of correct citations.

In Table 3, we present the evaluation results of various prominent LLMs on our private benchmark dataset, and the prompt settings are also listed in the table. We can see that: 1) After training with the proposed synthetic dataset, our Homer-70B achieved the best in both Answer Accuracy and Citation Score, showing that reasoning data generation is helpful in achieving better utilization. 2) GPT-3.5 and Command R+ often cite multiple documents to avoid missing the golden document, but only get an undesirable behavior. GPT-4o, using vanilla RAG instead of RTG prompts to align with current industry settings, performs the second showing its strong ability. 3) The base model Llama3-70B has a weak ability to select paragraphs, also showing the usefulness of the generated dataset.

Table 3: Experimental results on LTM-COT-1 benchmark

Model	Answer Accuracy	Citation Score
GPT-3.5 <sup>1</sup>	91.8%	56.3%
Command R+ <sup>2</sup>	96.8%	54.1%
Llama3-70B-instruct <sup>3</sup>	97.4%	76.2%
GPT-4o <sup>1</sup>	98.1%	86.7%
<b>Homer-70B</b>	<b>98.7%</b>	<b>91.2%</b>

<sup>1</sup> Vanilla RAG.

<sup>2</sup> Official API.

<sup>3</sup> RTG prompts w/o training.

#### Experimental results on public benchmarks.

Furthermore, we evaluated our model on a diverse set of standard academic benchmarks: ARC, HellaSwag, MMLU, TruthfulQA, Winogrande, and GSM8K in Table 4. The results show that our model achieves the highest average score across all compared models, indicating that despite being tested on various tasks, our model retains strong generalization performance learned from the generated RTG dataset.

Table 4: Experimental results on distinct public benchmarks.

Model	DBRX Instruct	CommandR Plus	Llama3-70B-instruct	Homer-70B
Model Size	A36B/132B	104B	70B	70B
ARC	68.90%	70.99%	71.42%	<b>72.70%</b>
HellaSwag	89.00%	<b>88.60%</b>	85.69%	88.13%
MMLU	73.70%	75.70%	<b>80.06%</b>	79.81%
TruthfulQA	<b>66.90%</b>	56.30%	61.81%	64.46%
Winogrande	<b>81.80%</b>	85.40%	82.87%	81.61%
GSM8K	66.90%	70.70%	<b>85.44%</b>	83.78%
<b>Averaged</b>	74.50%	74.60%	77.88%	<b>78.42%</b>
<b>Delta</b>	-3.91%	-3.82%	-0.53%	

Table 5: Experimental results of MedAgent-Zero and LTM-enhanced MedAgent-Zero in the MedQA subset [133]. The best performance is in bold and % is omitted for all results.

#Interactions	5,000	10,000	15,000	20,000	25,000	30,000	35,000	40,000	45,000	50,000
Original	93.06	94.44	93.06	94.44	91.67	94.44	93.06	93.06	91.67	91.67
LTM-enhanced	91.67	94.44	91.67	93.06	93.06	91.67	93.06	<b>95.83</b>	94.44	93.06

### 6.2.2 LTM Utilization for Agent Self-Evolution in Medical Domains

Currently, some studies have proposed distinct strategies to support agent self-evolution, but they mainly rely on the storage of interaction raw data and limited self-reflection. We believe that if the raw data can be preprocessed based on the concept of LTM, even when using the same evolution strategies, better performance should be achievable. Therefore, we adopted the MedAgent-Zero method from Tsinghua University’s study [133], but combined it with the LTM concept for improvement, and conducted the following experiment.

We first briefly review the core strategy of the MedAgent-Zero algorithm. By allowing the medical agent to interact with generated patients during diagnosis and treatment, the agent accumulates medical cases from correct interactions and reflects on experiences from incorrect interactions. After such evolution, its performance is evaluated on the respiratory subset of the MedQA dataset. Here, we focus on how to further improve and optimize the experiences gained from incorrect interactions, which is helpful in achieving better medical performance. Based on the idea of LTM construction, we believe that the experiences gained from single-interaction self-reflection should be further optimized to achieve better utilization.

Inspired by the necessary preprocessing of LTM, we propose that the experiences gained can be further refined through a few-shot prompting approach, making them more aligned with the task’s specific needs and characteristics. In other words, raw experience will be rewritten into task-specific LTM. As part of the specific strategy, we first selected three high-quality LTM data samples through manual annotation and optimization in a preliminary experiment. The three LTM samples were then used as few-shot examples to rewrite all experiences gained during interactions, resulting in a new version of LTM experiences.

For the evaluation, we followed the existing work’s approach by testing on the MedQA respiratory subset [133] and expanded the training scale to include 50,000 virtual patients. Experimental results are shown in Table 5, where we used a voting mechanism based on three repeated experiments to enhance stability. From the table, we can see that the LTM-rewritten strategy achieves better continuous evolution, avoiding the significant performance decline trend that can occur during interactions. While the original version saw a faster improvement in the initial phase, the unprocessed raw experience often led to conflicts as the number of interactions increased, eventually causing a decline in the model’s performance. In contrast, the LTM-rewritten model achieved the best results in this subset (95.83% accuracy), demonstrating the effectiveness of LTM in enhancing performance.

Although we only used a simple LTM rewriting method for this experiment, the results validated its effectiveness. We also believe that the introduction of LTM strategies will provide even greater benefits for self-evolution in broader scenarios.

### 6.2.3 LTM Utilization with memory system design

LLMs and LLM-powered agents demonstrate advanced language comprehension and can engage in high-quality interactions, making them promising tools for mental health diagnostics, particularly for depression [133]. However, several challenges are encountered in real-world application, including insufficient domain-specific knowledge, limited empathetic abilities, and an inability to learn and adapt based on prior experiences. To address these issues, the integration of LTM is crucial. Inspired by the training process of psychiatrists, we propose an innovative three-level hierarchical LTM structure—*Conversation Records*, *Electronic Medical Records (EMR)*, and *Diagnostic Skills*. This architecture enhances agents’ ability to accumulate and apply knowledge, improving diagnostic accuracy. To evaluate its efficacy, we developed a novel conversational agent simulation system to mimic diagnostic sessions between patient agents and psychiatrist agents. This system generates realistic depression diagnostic conversations, providing a valuable tool for training intern psychiatrists and conducting preliminary depression risk assessments in individuals exhibiting depressive symptoms.

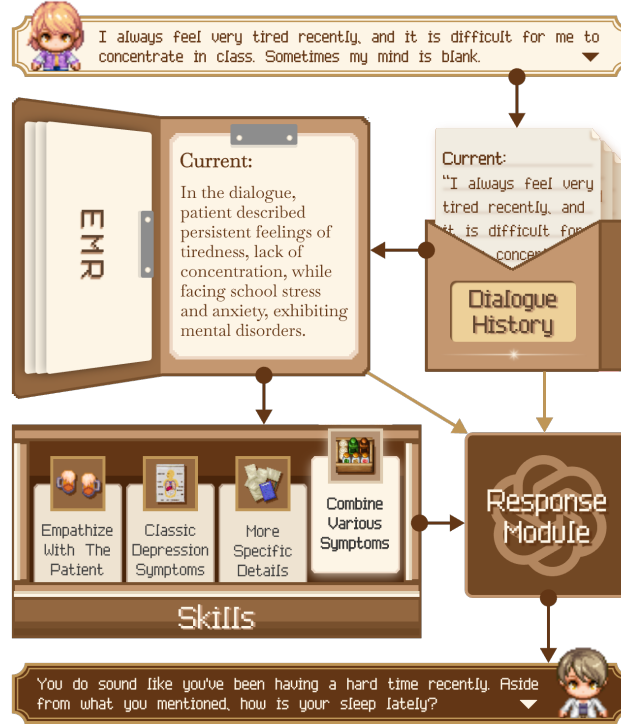


Figure 7: **The Tertiary Memory Structure.** The utterance of the diagnosis conversation will be stored in **Dialogue History**. The whole dialogue history in the session will be summarized into electronic medical records (**EMR**). **Skills** are generated by the supervisor plugin. All the memory will contribute to the dialogue generation.

Our experimental results, based on real-life depression diagnosis conversations D<sup>4</sup> dataset [49], demonstrate an average improvement of 6.05% in depression diagnosis accuracy and 1.8% in suicide risk prediction.

### Tertiary LTM Mechanism

The memory structure reflects the process psychiatrists follow: patient information is recorded through conversation, summarized into EMR, and used to refine diagnostic skills.

- **Conversation Records** represent the transcripts of the diagnosis conversation between patient agents and psychiatrist agents. It helps the psychiatrist agent to conclude the electronic medical records of the patient agents.
- **Electronic Medical Records (EMR)** summarize patient information, including demographics, chief complaints, and symptomatology. They serve as a concise reference for diagnosing similar cases and for improving diagnostic skills.
- **Diagnostic Skills** act as the optimizer in training psychiatrist agents. Diagnostic skills are generated by comparing the psychiatrist agent's diagnoses to ground truth diagnoses provided by professional psychiatrists in real life in D<sup>4</sup>. These skills help refine the agents' diagnostic accuracy, guiding them to align more closely with professional psychiatrists by highlighting overestimated and underestimated diagnosed symptoms.

As shown in Figure 7, this three-tiered memory structure facilitates the abstraction and refinement of raw diagnostic information. Each layer serves as a distillation process, transforming detailed data into increasingly concise and enduring representations. This structure allows the system to store large amounts of data efficiently while ensuring that the most critical information is preserved for future retrieval, thereby enhancing diagnostic performance over time [144].



We conducted series of experiment to assess the impact of these memory types on depression diagnosis by manipulating the memory module’s activation: 1) No memory module, 2) Memory retrieval from EMR, 3) Memory retrieval based on diagnostic skills, and 4) Memory retrieval from both EMR and diagnostic skills. The results are depicted in Figure 8.

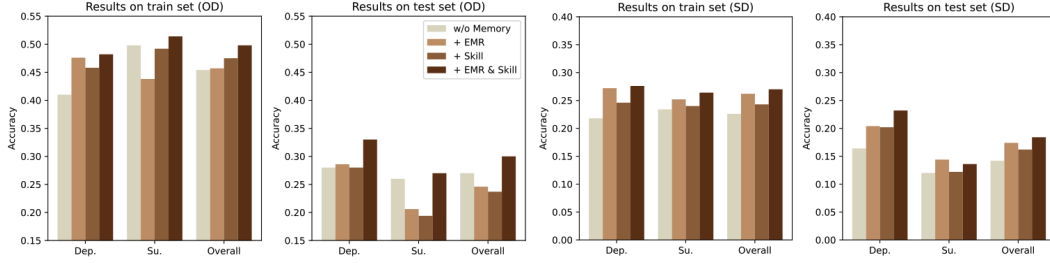


Figure 8: **The Results of Ablation Study on Memory Layers.** The first 2 pictures indicates the results(depression diagnosis and suicide risk prediction) based on the original dialogue history setting(OD), while the last 2 pictures implies the results based on the simulated dialogue setting(SD).

The findings suggest that diagnostic skills are particularly effective in Original Dialogue (OD) settings, likely due to their alignment capabilities with LLM models. In Simulated Dialogue (SD) settings, EMR proves more beneficial, as the simulations often lack the precision of real-world conversations, and EMR can assist by referencing similar patient cases. Combining EMR and diagnostic skills yielded the most consistent and accurate results, outperforming LLMs without memory integration [144].

### Effects of supervisor agent

Another key innovation in this study is the introduction of a supervisor agent, which monitors the patient’s symptoms, updates the symptom list, and generates targeted questions to assist the psychiatrist agent, streamlining the diagnostic process. It also manages the dialogue progression across predefined stages, improving diagnostic efficiency. This approach prevents the psychiatrist agent from asking redundant questions, instead guiding the conversation to focus on the unknown symptoms, thereby enhancing the efficiency and effectiveness of the diagnostic dialogue. Additionally, the supervisor agent manages the progression of the dialogue through three predefined stages of the diagnostic process [144].

We conducted an ablation experiment by disabling the reflection and question-generation, illustrating the fact that the supervisor agent significantly enhances risk prediction accuracy when its reflection and question-generation features are active. Such a supervisory agent could be replaced by professional psychiatrists to monitor the data quality of both the doctor and patient agents. When the goal is to assist real patients in obtaining easier diagnoses, the efficiency and accuracy can be improved. When the aim is to train actual doctors using patient agents, the realism of the agents can be enhanced. When both the doctor and patient are agents, this can serve as a platform for generating simulated data under the supervision of real doctors.

### 6.2.4 LTM Utilization with Real-time Weight Update

As we explore model self-evolution and the pivotal role of Long-Term Memory (LTM) in enhancing AI’s cognitive abilities, it is crucial to explore innovations in neural architectures that could better facilitate these advancements. While Transformer architectures excel in reasoning and knowledge compression—typically through methods like supervised fine-tuning (SFT)—this knowledge is stored in the model’s weights, which are usually frozen during inference. As a result, real-time weight updates are not feasible in standard Transformer models, limiting the model’s ability to learn and adapt dynamically during inference. To overcome this, exploring new architectures, potentially integrating sequence-based structures like RNNs, could enable breakthroughs in LTM utilization and support more effective model self-evolution.

One notable study, [57], proposes a new class of sequence modeling layers known as Test-Time Training (TTT) layers. These layers are designed to overcome the limitations of existing RNN layers by transforming the hidden state into a machine learning model itself, updated through self-supervised

learning steps. This innovative approach allows the model to continue learning even during test time, enhancing its performance on long-context tasks where traditional RNNs usually struggle. Thus, we propose to adopt TTT as the backbone structure to verify its performance in making use of LTM data, so some primary experiments are designed and conducted here.

## Experimental Settings

At the beginning of Section 4.2, we mentioned that the first challenge in using LTM data is adapting to continuously updated LTM data. So our experiment aims to verify the adaptability of TTT under distinct data distribution.

We choose multilingual datasets to conduct these experiments as distinct languages have various token distributions, where French, Chinese, and English are adopted. By performing inference and updating model weights in different languages, we assess two questions: 1) whether the model can effectively learn new distribution patterns and improve its performance on the corresponding test sets. 2) whether TTT leads to catastrophic forgetting, i.e., whether the model’s performance on the original test dataset deteriorates after adapting to the new distribution.

Specifically, we adopted the Book3 dataset [145] (a subset of The Pile) as the training dataset and trained a 1.3B TTT-linear model, denoted as  $Model_{En}$ . The Book3 dataset comprises a vast collection of English digitized books that range from classical literature to contemporary works, which has been widely used to train LLMs in long context. We selected two books for our experiments: "LE RÊVE DE SUZY" [146] as the French book and "The Smiling, Proud Wanderer (Xiào ào jiānghú)" [147] as the Chinese book. After tokenization, we extracted 32,000 tokens from each book, with a fixed 2,000 tokens reserved as the test set. The remaining tokens were used as the training set, and we evaluated the impact of different training set lengths on model inference performance.

First, we performed inference on  $Model_{En}$  using the training set and retained the inner-loop weights  $W_{train}$  updated after inference. In this context, the process of parametric learning can be viewed as compressing a large amount of training data into the model’s weights, where the model effectively learns and stores underlying patterns from its training data. Specifically, each step of inference updates the weights  $W_t$  by applying gradient descent on a self-supervised loss  $\ell$ . This process is akin to dynamically compressing the training data into a hidden state, which is then used to predict the next token. The update rule for the weights is defined as:

$$W_t = W_{t-1} - \eta \nabla \ell(W_{t-1}; x_t), \quad (4)$$

where  $\eta$  is the learning rate and  $\ell$  is the self-supervised loss, which measures the difference between the model’s predicted token and the true token. The gradient of the loss with respect to  $W$  drives the weight updates, enabling the model to adapt to new data distributions. This continual update mechanism ensures that the model captures significant input signals that produce large gradients, optimizing the weights accordingly.

Subsequently, we replaced  $Model_{En}$ ’s initial weights with  $W_{train}$ , the updated weights after inference, resulting in the domain-specific models  $Model'_{FR}$  and  $Model'_{CN}$ . These models are tailored to the respective French and Chinese datasets, having undergone inference and subsequent inner-loop weight updates on these new language-specific training sets.

Two experiments are designed to answer the two questions: 1) Experiment 1 is to compare the performance of  $Model_{En}$  and  $Model'_{FR}$  /  $Model'_{CN}$  on the test set, where the training loss and perplexity (PPL) are adopted as metrics. 2) Experiment 2 tests the performance of  $Model_{En}$  and  $Model'_{FR}$  /  $Model'_{CN}$  on the original Book3 test set to assess whether catastrophic forgetting occurs. To comprehensively evaluate the effect of TTT, we tried different lengths of train tokens and observed their impact on the performance of both the test tokens split from the Chinese/French books and the test tokens from the original Book3 pre-trained data.

## Experimental Results

### Experiments 1: Learn on new distribution pattern.

This experiment aimed to assess whether the TTT mechanism could effectively adapt to new data distributions and enhance model performance during inference.

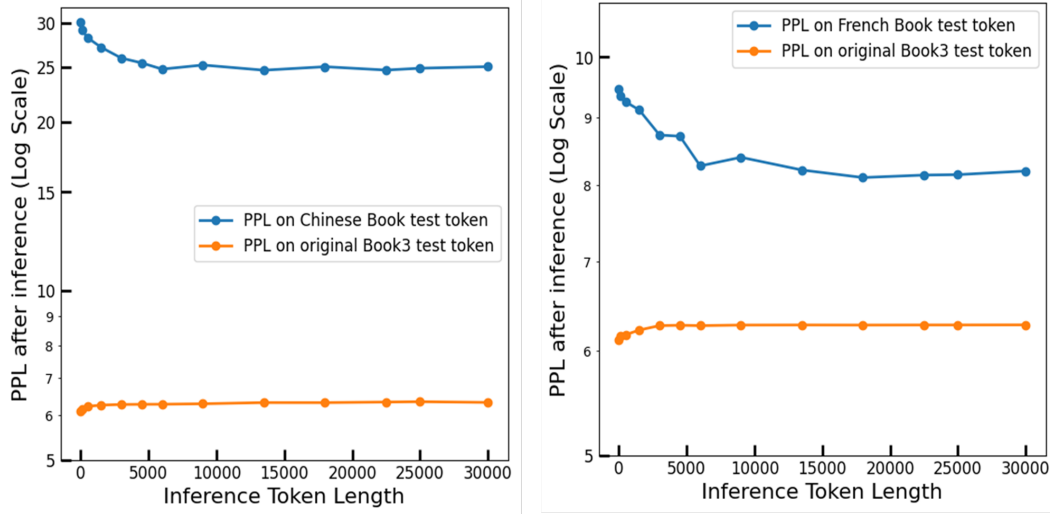


Figure 9: **Left: PPL on test tokens after inference on Chinese books. Right: PPL on test tokens after inference on French books.** Both figures show the experiment results on learning new distributions and catastrophic forgetting tests. The figures illustrate the differences in perplexity (PPL) between using the initial weights  $W_0$  (*inference token length = 0*) and the updated weights  $W_{train}$  after inference on Chinese and French books, measured on the corresponding test tokens. The results highlight how the model’s performance shifts after inference on each language tokens.

As illustrated in Figure 9, the results demonstrate that by updating the model’s weights ( $W_0 \rightarrow W_{train}$ ) during inference on both the French and Chinese datasets, the model successfully adapted to the new distribution patterns. This adaptation led to a significant reduction in perplexity on the corresponding Chinese and French book test tokens. The plots show that as the inference token length increased, model performance improved consistently before stabilizing beyond a certain point. These findings confirm that TTT effectively learns from new language distributions, improving performance on both the French and Chinese test datasets, which directly addresses our first research question above.

### Experiment 2: Catastrophic Forgetting Test.

The second experiment aimed to determine whether updating model weights through TTT inference leads to catastrophic forgetting, which would manifest as a significant performance drop on the original test set.

After conducting inference on the French and Chinese datasets and updating the weights, the model exhibited a minor increase in loss and perplexity when tested on the original Book3 dataset in English. The performance degradation was more pronounced after French inference compared to Chinese inference, though the impact was relatively limited overall, especially when shorter train token lengths were used. These findings suggest that while TTT does introduce some degree of catastrophic forgetting, the effect is minor, and the model largely retains its ability to generalize to the original data distribution.

To summarize, the experimental results demonstrate that TTT can effectively help the model learn new language distributions and improve performance on the test sets. While there is some slight increase in PPL during the catastrophic forgetting test, as seen in the results, this impact is minimal compared to the substantial gains in learning the new language distributions. This indicates that the model retains its generalization ability even after adapting to new data.

### Explore on the Sequence Modeling Structure

In our current experiments, the modeling function  $f$  used in the TTT layers is primarily instantiated as either a MLP or a linear model. During the adaptation process to new language distributions, such as French and Chinese, we observed that across various model sizes, MLP consistently outperformed the linear model in learning new distribution patterns. This suggests that a more expressive  $f$  is better equipped to capture the complexities of evolving new distributions.

Future work could focus on developing more sophisticated instantiations of  $f$ , particularly in the context of models that integrate LTM. As context length increases, such as in tasks across millions of tokens,  $f$  may need to scale in complexity. One promising direction is to explore hierarchical models, where  $f$  could incorporate recurrent layers (e.g., RNN or LSTM) or convolutional neural networks that are specifically designed to accumulate and compress information over extended sequences. This would enable the model to continuously learn and adapt, even in environments with extremely long contexts.

Moreover, another interesting research field is multi-level learning to learn. If  $f$  itself is a self-attention layer, it could act as a nested inner loop, where each level refines the understanding of the context in relation to LTM data. In this framework, each successive layer of learning could address more abstract or temporally distant features, effectively creating a multi-stage adaptation process that optimizes both short-term and long-term dependencies. This structure would be especially relevant for models aiming to leverage LTM in dynamic or changing environments.

Except for the study of TTT, the evolution towards new model structures aligns well with the goals of model self-evolution and the effective utilization of LTM. By incorporating advanced RNNs with expressive hidden states, there are also some other efficient architectures like Mamba [148], RWKV [149], and biologically-inspired neuron diversity [150; 151; 152], so we can build more sophisticated and efficient models capable of better simulating human-like cognitive processes.

### 6.3 Development of LTM-Based Multi-Agent Framework

In Section 5.3, we employed a hybrid strategy of RAG and training to utilize the generated LTM data, achieving promising performances. Here, we will introduce a LTM-Based Multi-Agent Framework called Omne. Omne is a deeply customized development framework based on the AutoGen Multi-Agent Framework, specifically designed to address the practical application challenges of LTM in AI systems. It extends a range of memory-related infrastructure, including a unified memory model, a multimodal message processing system, and flexible memory storage and manipulation mechanisms. The core goal of Omne is to provide a comprehensive solution that enables the effective deployment of LTM in real-world engineering projects, thereby enhancing the long-term memory capabilities and task-handling efficiency of AI systems.

#### 6.3.1 Core Modules of Omne

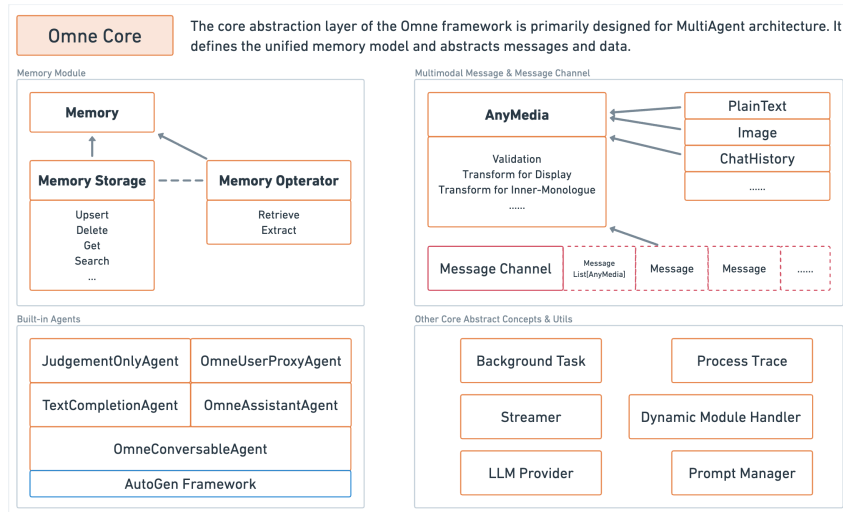


Figure 10: Illustration of Omne Core, which constitutes of a memory module, multimedia messaging channel, built-in agents and other related tasks.

Omne Core serves as the primary abstraction layer for the framework, focusing on the design of a multi-agent architecture. It defines a unified memory model and abstracts the processing of messages

and data. By offering a standardized memory abstraction mechanism, the Omne framework allows developers to manage memory with consistent operations at any level of the framework.

Coupled with the message channel mechanism provided by the framework layer, Omne enables more complex memory operations, such as asynchronous memory retrieval and on-demand memory extraction and reprocessing across different business scenarios. This design allows Omne to efficiently handle various complex memory and business integration scenarios, offering a robust infrastructure for developing highly intelligent AI applications.

Additionally, whether using traditional Retrieval-Augmented Generation (RAG) techniques or in-context learning approaches, the memory abstraction provided by the Omne framework can be leveraged. This flexibility allows developers to mix and match different technologies to meet specific business requirements and application scenarios.

In terms of multimodal message processing, Omne offers a complete set of message storage and manipulation abstractions. Developers can utilize built-in media types like plain text and images, or extend them based on specific business needs. The framework also supports flexible message visibility definitions to accommodate different business contexts.

For example, when handling chat history that includes both text and images, if using an LLM provider that supports image processing (e.g., GPT-4), the message channel will automatically preserve the images and convert them into a compatible format. If switching to a model that does not support image input (e.g., Llama 3.1), the system will automatically convert images into descriptive text. This design ensures data continuity and usability when transitioning between models with varying capabilities.

### 6.3.2 Omne Assistant

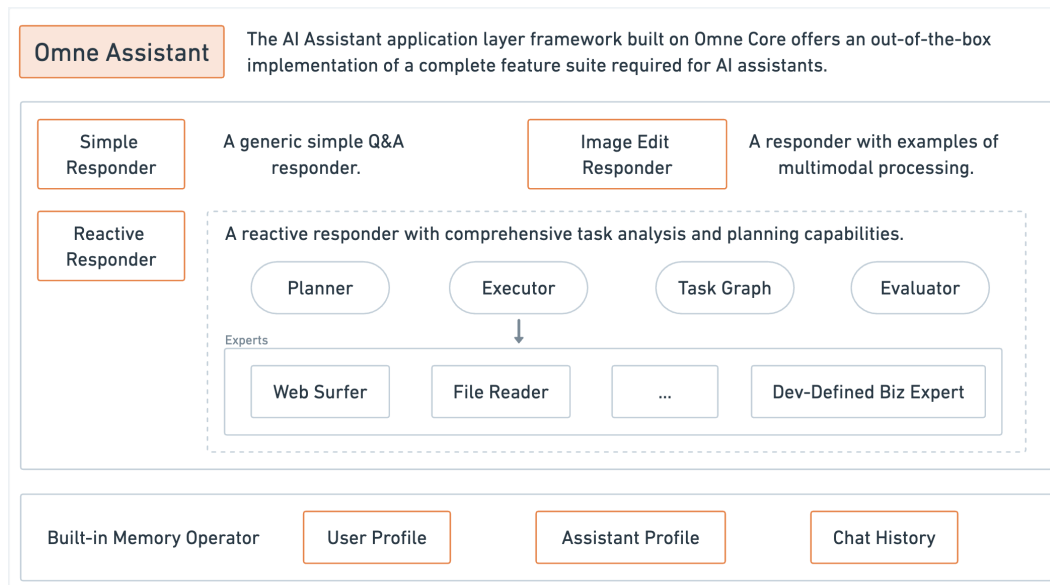


Figure 11: Omne Assistant

Built on Omne Core, Omne Assistant provides a ready-to-use application layer framework specifically designed for the development of AI assistants in chat scenarios. It includes the essential functionalities required for AI assistants, enabling developers to quickly build fully functional chatbots without the need to design foundational components from scratch.

Omne Assistant comes with a **Simple Responder**, a general-purpose Q&A responder that can handle basic user chat interactions for instant communication. Additionally, the framework provides a **Reactive Responder**, which has advanced task analysis and planning capabilities, allowing it to manage more complex user requests that require multi-step reasoning and task orchestration.

With these built-in components, Omne Assistant allows developers to focus on implementing custom functionalities, speeding up the development and deployment of AI assistant applications equipped with long-term memory capabilities.

### 6.3.3 Case Study

In practical applications for planning complex tasks, it is often necessary to balance task quality with generation speed. Therefore, we explore the possibility of allowing sufficiently powerful models to "pre-plan" various tasks and store the results as "memories." During actual online inference, similar tasks can be "recalled," and a faster model can use these pre-planned results to carry out the real task planning.

Using the tools provided by the Omne framework, we can quickly implement In-context Learning (ICL) capabilities in business applications. Below is a simple process outline. *Usage Example: Enhancing Planning Effectiveness for Highly Complex Tasks Using In-context Learning.*

- Prepare a set of task objectives and pre-constructed task contexts, generating different plans by selecting various models and prompting strategies.[Offline Phase]
- Automatically filter effective experiences through simulation of the execution of the planning results.[Offline Phase]
- Involve human intervention to assess the validity of the experiences and make adjustments to the results if necessary.[Offline Phase][Optional]
- Save the effective experiences as memories (the simplest method is to combine task objectives and contexts into task descriptions, using these descriptions as vector retrieval indices).[Offline Phase]
- Identify user objectives and construct the context needed for task planning.[Online Phase]
- Retrieve "experience memories" to find similar planning experiences and reference these experiences for actual task planning.[Online Phase]

At this point, we have implemented In-context Learning for complex task planning using the Memory Operator and other related suites provided by Omne. It is noteworthy that due to the well-designed framework, the entire process is iterative and minimally invasive to business operations.

The Omne framework is designed and implemented with the goal of applying "memory in a Multi-Agent architecture." It consistently treats "memory" as a primary element and core concept, aiming to build a technical framework that can grow alongside long-term memory technologies and facilitate engineering deployment.

### 6.3.4 GAIA benchmark

We evaluated the Omne framework on the GAIA benchmark [153], a rigorous test for general AI assistants comprising over 400 question-answering tasks. These tasks involve complex logical deduction and reasoning processes, requiring capabilities such as numerical calculations, web browsing, video and speech processing, and file manipulation.

To explore the boundaries of AI models, we utilized GPT-4o and o1-preview[154] within Omne. Omne was equipped with four tools: web surfer, bing search engine, file reader based on llamaparse[155] (for various file manipulation), and a logic expert built with o1-preview. Additionally, we implemented a scene router based on o1-preview to assess whether a question could be resolved locally, without external web information and resources. If one question can be solved locally, then we used logic expert to solve it directly (and file reader if needed). Otherwise we turn to a gpt-4o based multi-agent system with all 4 tools mentioned above. Our final results were submitted to official evaluation server of GAIA in huggingface.

Leveraging OpenAI's GPT-4o and o1-preview as base models, Omne achieved first place (40.53%) on the test set and second place (46.06%) on the validation test, establishing a new state-of-the-art on this milestone dataset in AI research. Notably, Omne attained 26.53% accuracy on the most sophisticated and demanding level 3 questions, demonstrating its potential to solve real-world problems by leveraging powerful foundational models, particularly those with strong reasoning and logical capabilities.

Results: Test		Results: Validation						
Agent name	Average score (%)	Level 1 score (%)	Level 2 score (%)	Level 3 score (%)	organisation	Model family	Submission date	
Trase Agent v0.2	47.27	58.49	46.51	26.92	Trase Systems	Multi-Agent - Gemini	2024-10-16	
Magentic-1 (o1)	46.06	56.6	46.51	23.08	MSR AI Frontiers	o1 and GPT-4o (var: 100k)	2024-10-15	
omne	46.06	60.38	44.19	23.08		o1-preview, gpt-4o	2024-10-20	
Hugging Face Agents + GPT-4o	44.24	58.49	43.02	19.23	Hugging Face 🤖	GPT-4o	2024-06-26	
AgentIM v1.1	40	50.94	40.7	15.38		GPT-4-turbo	2024-05-28	
Trase Agent	40	47.17	40.7	23.08	Trase Systems	GPT-4o	2024-08-07	
Multi-Agent Experiment v0.1 (pove)	39.39	54.72	38.37	11.54	MSR AI Frontiers	GPT-4-turbo	2024-03-01	

Figure 12: GAIA leaderboard in validation set

Results: Test		Results: Validation						
Agent name	Average score (%)	Level 1 score (%)	Level 2 score (%)	Level 3 score (%)	organisation	Model family	Submission date	
omne_v0.1	40.53	53.76	37.11	26.53		o1-preview, gpt-4o	2024-10-20	
Trase Agent v0.2	39.53	55.91	37.74	14.29	Trase Systems	Multi-Agent - Gemini	2024-10-11	
Multi Agent	38.87	53.76	37.74	14.29			2024-10-10	
das_agent_v0.4	38.21	51.61	36.48	18.37	NA	GPT-4o	2024-10-04	
Magentic-1 (o1)	38	54.84	32.7	22.92	MSR AI Frontiers	GPT-4o and o1-preview	2024-10-19	
Trase Agent v0.1	35.55	50.54	33.33	14.29	Trase Systems	Fine-tuned GPT-4o	2024-09-04	
Sibyl System v0.2	34.55	47.31	32.7	16.33	Baichuan Inc.	GPT-4o	2023-11-03	
Hugging Face Agents + GPT-4o	33.33	49.46	28.3	18.75	Hugging Face 🤖	GPT-4o	2024-06-27	

Figure 13: GAIA leaderboard in test set

## 7 Our Future Plans

As discussed in Section 2, we believe that model personalization could be one of the key pathways to a second emergence of intelligence. Here, we will briefly describe the key research directions and areas we will focus on in the future, and we welcome more researchers to join us in our efforts to use LTM to improve model personalization. By accumulating data, such as LTM data, different LLM agents could develop diverse and differentiated abilities. The combination of diverse and differentiated individuals is more likely to spark new intelligence. We are also confident that LTM and model personalization could play a pivotal role in achieving this ambitious goal. In the future, we aim to answer the following challenging questions:

1. How to better construct LTM data?
2. How to Design New Model Architectures for LTM?
3. How Can LTM Help Users Ask Better Questions?
4. How to Integrate LTM with Inference-Time Search?
5. How to Use LTM for Agent Self-Evolution in Complex Scenarios?
6. How to Use LTM in Multi-Agent Scenarios?

### 7.1 How to better construct LTM data?

To better collect longitudinal data, it is essential to establish a system capable of tracking individual data over time. This can be achieved by incentivizing users to continuously engage with data platforms while ensuring privacy and security, encouraging long-term data collection. Additionally, the use of wearable devices and IoT sensors can provide continuous multimodal data collection, which is crucial for building long-term memory models. To balance diversity and data consistency, it is important to curate datasets that represent a wide range of user characteristics while maintaining uniformity in data labeling and collection methods.

We are also advancing our efforts in data synthesis and developing an end-to-end data synthesis system. This system incorporates a continuous learning and adaptive data synthesis framework, allowing it to evaluate the performance of synthetic data and adjust the generation process accordingly. The system includes a feedback loop, enabling it to detect data deficiencies in model personalization and automatically optimize and refine data synthesis methods.

## 7.2 How to Design New Model Architectures for LTM?

One perspective is that as models progress, the increasing length of the context window may provide a pathway toward achieving long-term memory. However, given that the context window is currently limited to the KV cache of short-term sessions, it cannot genuinely achieve cross-task, cross-session long-term memory. We believe that it is necessary to redesign the underlying model architecture, transitioning LLMs from merely relying on "context windows" to a deeper, more structured long-term memory mechanism. The key challenge in this direction is how to enable the model not only to store and retrieve historical information but also to dynamically manage this memory through adaptive adjustments.

Our vision is to conceptualize LLMs as world models, drawing inspiration from the concept of cortical columns, where each layer of the model can be seen as an independent entity. In this way, each layer of the model can dynamically adjust its weights based on new data distributions. This could significantly enhance computational efficiency and provide the model with more flexibility, enabling it to handle more complex and diverse data distributions. Similar to neurons in the brain, different model layers can activate distinct memory fragments based on specific tasks, forming an adaptive reasoning mechanism. We have made preliminary attempts on this idea with TTT, and future research will explore this direction in more depth.

More importantly, the concept of allowing models to autonomously adjust long-term memory based on environmental sensing opens up new avenues for the development of more intelligent AI. Unlike traditional pre-set memory management, models under this paradigm can perceive changes in external environments in real-time and adjust their memory structure accordingly. This adjustment goes beyond simple recall and storage of data and includes the perception and optimization of data distribution changes, ensuring that the model can efficiently and personally adapt to various tasks.

This form of learning not only enhances the model's adaptability but also enables it to provide more personalized reasoning when facing diverse data. This flexible memory management and optimization mechanism will empower future LLMs to not only better understand user needs but also handle complex reasoning tasks more efficiently. The core of this vision is for the model to optimize its memory management as the environment changes, gradually evolving into an intelligent agent capable of self-learning and optimization. This agent will demonstrate flexibility in handling dynamic and uncertain environments. The integration of long-term memory and model personalization will pave the way for advancements in future AI development.

## 7.3 How Can LTM Help Users Ask Better Questions?

We believe that while LLMs are already capable of providing good responses to user questions, users may struggle to formulate optimal questions, thereby limiting the effectiveness and potential applications of LLMs.

For example, in a news reading context where a model is expected to generate in-depth questions that stimulate critical thinking and enhance readers' engagement. However, obtaining training data that enables a model to ask "high-quality" questions poses a significant challenge. One potential approach to generating such data is to extract and summarize questions from news interview programs. These questions are typically crafted by experienced hosts and journalists, who are adept at posing thought-provoking and insightful queries. By analyzing and summarizing these questions alongside the relevant news content, we can construct a dataset that embodies the key characteristics of high-quality inquiry, such as probing assumptions, exploring diverse perspectives, and addressing complex issues.

Exploration is a key method by which agents discover new possibilities, while pruning accelerates convergence and avoids ineffective paths. Striking a balance between exploration and pruning is crucial to improving an agent's ability to ask better questions.

How can an agent find effective strategies during initial exploration and then quickly prune ineffective paths in the later stages? Could we introduce a dynamic exploration-pruning strategy, where the agent leans more toward exploration during the early stages of high uncertainty and accelerates pruning as certainty increases? When facing new problems, how can an agent leverage long-term memory to transfer past strategies from similar problems to the current task, thus reducing exploration time and resource consumption?



How can long-term memory guide pruning? Can agents identify inefficient or unfeasible paths based on past exploration results, allowing them to skip these paths in new tasks? For instance, in a medical system, an agent could quickly prune ineffective treatment plans based on past data, focusing instead on more effective strategies. How can we balance local exploration with global pruning? In multi-agent systems, one agent's exploration may influence the strategies and outcomes of others. How can we ensure that the pruning process does not hinder other agents' exploration or the evolution of the overall system?

## **7.4 How to Integrate LTM with Inference-Time Search?**

A major advance in OpenAI's O1 model is the integration of reinforcement learning (RL) to enhance LLMs' multi-step reasoning and planning capabilities. Specifically, RL introduces a powerful inference-time search mechanism, enabling the pretrained model to effectively "think" during the inference phase. By generating multiple reasoning trajectories and selecting the optimal one, the model can refine its reasoning process in real-time. This combination of "memorization" (long-term memory, LTM) and "thinking" (dynamic reasoning) will jointly define the future intelligence level of AI systems.

We pose the following questions for further exploration:

How can the dynamic evolution of long-term memory (LTM) enhance LLM's reasoning and search abilities? The dynamic evolution of LTM could play a critical role in supporting more robust and adaptive reasoning. LTM allows a model to not only store but also update knowledge across interactions, which could directly improve reasoning by providing more relevant historical context and insights. LTM could act as a repository for key information, reducing the need for the model to "relearn" common knowledge with each inference. This could allow the LLM to avoid repeating mistakes or inefficient paths it had previously explored. By evolving this memory over time, the LLM can avoid redundant or inefficient search paths, guiding it toward more relevant reasoning trajectories.

Can LTM be useful for representing LLMs' thought processes as a sequence of state abstractions, so that reinforcement learning can efficiently explore the large space of reasoning paths? We believe LTM could provide as a structured representation of the LLM's thought processes, capturing key intermediate states during multi-step reasoning. By encoding these states into meaningful abstractions, LTM can create a more tractable search space for RL to navigate. This would allow the reinforcement learning agent to focus on important or novel reasoning paths, making exploration more efficient and targeted. Essentially, LTM could help map the vast reasoning space into manageable, hierarchical abstractions that RL can use to prune less optimal trajectories and focus on those with higher potential outcomes.

How can reinforcement learning-based search methods enhance LTM? RL can improve LTM by augmenting it through real-time exploration during the inference phase. For example, if the model's memory lacks certain knowledge or if its reasoning process encounters ambiguities, an RL-driven search method could generate targeted questions or reasoning trajectories to actively fill the gaps in memory.

Consider a scenario where an AI system encounters a new task or topic that challenges its current knowledge base. The RL-based approach could iteratively search for relevant questions that would retrieve or generate knowledge to augment LTM. For another example, RL could evaluate different strategies for updating LTM by comparing how changes in memory influence the success of reasoning over time. By using feedback signals from multiple reasoning trajectories, the RL agent can continuously refine how the memory is updated and leveraged, ensuring that the memory evolves in a way that enhances future reasoning tasks.

## **7.5 How to Use LTM for Agent Self-Evolution in Complex Scenarios?**

In the previous discussions, we mentioned that LTM plays a crucial role in the capability evolution of agents, particularly with the additional LTM gained from interaction history and reflective reasoning, which can directly enhance task performance. However, these improvements come with a prerequisite: we must have ground-truth or a well-established evaluation mechanism to provide feedback, guiding the accumulation of the agent's LTM.

Nevertheless, in complex task scenarios, on one hand, collecting a large amount of high-quality ground-truth is very challenging, often requiring high-quality annotations from experts and significant participation. On the other hand, building evaluators for individual tasks is also very difficult, especially in generative tasks, where effective evaluation methods are currently lacking. Therefore, how to use LTM for agent self-evolution in complex systems remains a vital challenge.

We believe a feasible research approach is to achieve self-evolution through environmental feedback, thereby enabling effective accumulation and use of LTM. Specifically, environmental feedback may come from the physical world, the Simulacrum world, or a combination of both. The evolution in the physical world is closely related to embodied AI research, but the evolution efficiency is relatively low; constructing a Simulacrum world is very challenging, but it offers higher efficiency and lower cost. Moving forward, we will conduct in-depth exploration on Agent Self-Evolution combined with LTM.

## 7.6 How to Use LTM in Multi-Agent Scenarios?

We believe that multi-agent cooperation may lead to the second time of intelligence emergence, which is the ultimate goal we pursue for agent self-evolution. In multi-agent systems, individual agents' actions are typically discrete and local, but the system requires them to co-evolve to achieve global optimization. Agents must use long-term memory for task planning while continuously exploring and pruning to enable overall evolution. However, achieving this goal presents a series of intermediate challenges, especially in the accumulation and utilization of LTM during multi-agent collaboration. Specifically, the following points are crucial.

First, how can we accumulate individual LTMs for collaborating agents? In collaborative scenarios, since multiple agents are responsible for different parts of a complex task, the final completion of the task cannot be simply used as direct feedback for the accumulation of each agent's LTM. Instead, an additional evaluation and decomposition strategy is required to convert it into feedback for each stage, similar to a step-grained reward rather than an overall reward (especially when this reward may be binary). This is critical for the co-evolution of multiple agents.

Second, should there be a certain mechanism for sharing and utilizing LTM during multi-agent collaboration? Currently, there are two forms of multi-agent collaboration: one where each agent completes a sub-module of the task, and another where multiple agents collaborate through discussion to make the final decision. Regardless of the specific form of cooperation, if we can introduce some form of shared memory or communication mechanism to help agents understand the decisions and outcomes of other agents, allowing them to consider global objectives more effectively when making local decisions, we believe it will directly and effectively contribute to enhancing overall capabilities.

We believe that addressing the above issues will significantly advance the evolution of multi-agent collaboration.

## 8 Conclusion

In this report, we propose that long-term memory mechanisms and the evolving nature of knowledge will be crucial. Current models treat all data uniformly, from ancient to modern, without capturing the gradual progression of knowledge. Human cognition, however, evolves—children learn from simple to complex concepts, and knowledge builds upon itself over time. This temporal structuring of data, where models learn through a progression of difficulty or sequence, could help them grasp not just static facts but the relationships and evolution of knowledge itself.

By leveraging longer-term memory architectures, models could begin to capture this evolutionary aspect of learning. Instead of simply memorizing information, they could learn the dynamics of knowledge development, allowing them to operate over extended temporal scales. This would also involve recursive learning, where feedback from real-world environments shapes the model's growth, leading to a self-boosting mechanism. As neural networks evolve, they may focus more on fast pruning of search spaces during gradient descent and benefit from enhanced parallelism in recursive computations. Ultimately, scaling laws will persist, but model architectures will need to adapt to handle increasing complexity while maintaining efficiency.

Looking forward, model personalization, inspired by human diversity, could drive a new level of intelligence. In multi-agent systems, diverse, LTM-powered agents could collaborate more effectively, balancing exploration and pruning to achieve global optimization. This personalization and adaptability may be key to fostering a second emergence of intelligence, where agents evolve together to tackle increasingly complex problems. This exploration into long-term memory and evolutionary learning will drive advancements in AI, particularly for models that need to continuously adapt, learn, and personalize over time. In this paper, we aim to articulate our vision and roadmap for future research directions and critical focus areas. We invite fellow researchers to engage with our findings and collaborate with us in the pursuit of utilizing LTM to enhance model personalization.

## References

- [1] Min Xu, Jeanne M. David, and Suk Hi Kim. The fourth industrial revolution: Opportunities and challenges. *International Journal of Financial Research*, 9(2):90–95, April 2018.
- [2] Sumanth Doddapaneni, Krishna Sayana, Ambarish Jash, Sukhdeep Sodhi, and Dima Kuzmin. User embedding model for personalized language prompting, 2024.
- [3] I. Lehr Brisbin, Charles T. Collins, Gary C. White, and D. Alasdair McCallum. A new paradigm for the analysis and interpretation of growth data: The shape of things to come. *The Auk*, 104(3):552–554, 1987.
- [4] Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model?, 2020.
- [5] Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 21702–21720. Curran Associates, Inc., 2023.
- [6] Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Frozen pretrained transformers as universal computation engines. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(7):7628–7636, 2022.
- [7] Fabricio Ballarini, Diego Moncada, Maria Cecilia Martinez, Nadia Alen, and Haydée Viola. Behavioral tagging is a general mechanism of long-term memory formation. *Proceedings of the National Academy of Sciences*, 106(34):14599–14604, 2009.
- [8] Jeff Hawkins. *A thousand brains: A new theory of intelligence*. Basic Books, 2021.
- [9] Brian R Johnson and Timothy A Linksvayer. Deconstructing the superorganism: social physiology, groundplans, and sociogenomics. *The Quarterly Review of Biology*, 85(1):57–79, 2010.
- [10] N.J. Hagens. Economics for the future – beyond the superorganism. *Ecological Economics*, 169:106520, 2020.
- [11] Huao Li, Yu Chong, Simon Stepputtis, Joseph Campbell, Dana Hughes, Charles Lewis, and Katia Sycara. Theory of mind for multi-agent collaboration via large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2023.
- [12] Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. Are emergent abilities of large language models a mirage? In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 55565–55581. Curran Associates, Inc., 2023.
- [13] Christopher J. Hewer and Robert Roberts. History, culture and cognition: Towards a dynamic model of social memory. *Culture & Psychology*, 18(2):167–183, 2012.
- [14] Ali Oroojlooy and Darius Adam Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *Applied Intelligence*, 53:13677–13722, 2023.
- [15] Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, and Xia Hu. Data-centric ai: Perspectives and challenges. In *Proceedings of the 2023 SIAM International Conference on Data Mining (SDM)*, pages 945–948. Society for Industrial and Applied Mathematics, 2023.
- [16] C. Delpierre and T. Lefèvre. Precision and personalized medicine: What their current definition says and silences about the model of health they promote. implication for the development of personalized health. *Frontiers in Sociology*, 8:1112159, 2023.
- [17] Ajay Bandi, Pydi Venkata Satya Ramesh Adapa, and Yudu Eswar Vinay Pratap Kumar Kuchi. The power of generative ai: A review of requirements, models, input-output formats, evaluation metrics, and challenges. *Future Internet*, 15:260, 2023.

- [18] E. Talvitie. Self-correcting models for model-based reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31. AAAI Press, 2017.
- [19] J. C. Schön. Structure prediction in low dimensions: concepts, issues and examples. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 381(2250), Jul 2023.
- [20] Marcel Nonnenmacher, Srinivas C Turaga, and Jakob H Macke. Extracting low-dimensional dynamics from multiple large-scale neural population recordings by learning to predict correlations. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [21] Leonard A Smith. Identification and prediction of low dimensional dynamics. *Physica D: Nonlinear Phenomena*, 58(1):50–76, 1992.
- [22] Zhengwei Tao, Ting-En Lin, Xiancai Chen, Hangyu Li, Yuchuan Wu, Yongbin Li, Zhi Jin, Fei Huang, Dacheng Tao, and Jingren Zhou. A Survey on Self-Evolution of Large Language Models, June 2024.
- [23] S Hochreiter. Long short-term memory. *Neural Computation MIT-Press*, 1997.
- [24] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, volume 159, pages 1877–1901, 2020.
- [25] James L. McGaugh. Memory—a century of consolidation. *Science*, 287(5451):248–251, 2000.
- [26] Larry R. Squire. Memory systems of the brain: A brief history and current perspective. *Neurobiology of Learning and Memory*, 82(3):171–177, 2004.
- [27] S Sridhar, A Khamaj, and MK Asthana. Cognitive neuroscience perspective on memory: overview and summary. *Frontiers in Human Neuroscience*, 17:1217093, 2023.
- [28] Susana A. Merlo, Mariano A. Belluscio, María E. Pedreira, et al. Memory persistence: From fundamental mechanisms to translational opportunities. *Translational Psychiatry*, 14:98, 2024.
- [29] Benjamin Straube. An overview of the neuro-cognitive processes involved in the encoding, consolidation, and retrieval of true and false memories. *Behavioral and Brain Functions*, 8:35, 2012.
- [30] Yadin Dudai. The neurobiology of consolidations, or, how stable is the engram? *Annual Review of Psychology*, 55:51–86, 2004.
- [31] Larry R. Squire and John T. Wixted. The cognitive neuroscience of human memory since h.m. *Annual Review of Neuroscience*, 34:259–288, 2011.
- [32] Angelina R. Sutin, Justin Brown, Martina Luchetti, Damaris Aschwanden, Yannick Stephan, and Antonio Terracciano. Five-factor model personality traits and the trajectories of episodic memory: A meta-analysis of individual participant data from 120,640 participants and 471,821 memory assessments. *The Journals of Gerontology: Series B*, 78(3):421–433, March 2023.
- [33] Y. N. Kenett, S. Humphries, and A. Chatterjee. A thirst for knowledge: Grounding curiosity, creativity, and aesthetics in memory and reward neural systems. *Creativity Research Journal*, 35(3):412–426, 2023.
- [34] D. Z. Hambrick and R. W. Engle. The role of working memory in problem solving. In J. E. Davidson and R. J. Sternberg, editors, *The Psychology of Problem Solving*, pages 176–206. Cambridge University Press, 2003.

- [35] Kun Zhao and Ginestra Bianconi. Social interactions model and adaptability of human behavior. *Frontiers in Physiology*, 2, 2011.
- [36] M. S. Weldon and K. D. Bellinger. Collective memory: Collaborative and individual processes in remembering. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 23(5):1160–1175, 1997.
- [37] M. Wimber, A. Alink, I. Charest, et al. Retrieval induces adaptive forgetting of competing memories via cortical pattern suppression. *Nature Neuroscience*, 18:582–589, 2015.
- [38] Wanjun Zhong, Lin Guo, Qian Gao, Hao Ye, and Yijun Wang. Memorybank: Enhancing large language models with long-term memory. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(17):19724–19731, 2024.
- [39] Yiming Du, Hongru Wang, Zhengyi Zhao, Bin Liang, Baojun Wang, Wanjun Zhong, Zezhong Wang, and Kam-Fai Wong. Perltqa: A personal long-term memory dataset for memory classification, retrieval, and fusion in question answering. In *Proceedings of the 10th SIGHAN Workshop on Chinese Language Processing (SIGHAN-10)*, pages 152–164. Association for Computational Linguistics, 2024.
- [40] Y. Wang, N. L. Zhang, and T. Chen. Latent tree models and approximate inference in bayesian networks. *arXiv preprint arXiv:1401.3429*, 2014.
- [41] W. Wei and L. Liu. Trustworthy distributed ai systems: Robustness, privacy, and governance. *arXiv preprint arXiv:2402.01096*, 2024.
- [42] Stevan Milovanović, Zorica Bogdanović, Aleksandra Labus, Dušan Barać, and Marijana Despotović-Zrakić. An approach to identify user preferences based on social network analysis. *Future Generation Computer Systems*, 93:121–129, 2019.
- [43] Longbing Cao. In-depth behavior understanding and use: The behavior informatics approach. *CoRR*, abs/2007.15516, 2020.
- [44] Benediktus Rolando and Herry Mulyono. Unlocking the power of data: Effective data-driven marketing strategies to engage millennial consumers. *Transekonomika: Akuntansi, Bisnis Dan Keuangan*, 4(3):303–321, 2024.
- [45] Luka Abb and Jana-Rebecca Rehse. A reference data model for process-related user interaction logs. *Springer*, pages 57–74, 2022.
- [46] Dipen R Bhuva and Sathish Kumar. A novel continuous authentication method using biometrics for iot devices. *Internet of Things*, 24:100927, 2023.
- [47] Muhammad Umar Khan, Zainoor Ahmad Choudry, Sumair Aziz, Syed Zohaib Hassan Naqvi, Afeefa Aymin, and Muhammad Atif Imtiaz. Biometric authentication based on emg signals of speech. *IEEE*, pages 1–5, 2020.
- [48] Sonia Sabir, Inayat Ali, and Eraj Khan. H-lps: a hybrid approach for user’s location privacy in location-based services. *arXiv preprint arXiv:2212.08241*, 2022.
- [49] Bing Yao, Changying Shi, Lei Zou, et al. D4: a chinese dialogue dataset for depression-diagnosis-oriented chat. *arXiv preprint arXiv:2205.11764*, 2022.
- [50] Alexey Kurakin, Natalia Ponomareva, Umar Syed, Liam MacDermed, and Andreas Terzis. Harnessing large-language models to generate private synthetic text. *arXiv preprint arXiv:2306.01684*, 2023.
- [51] Avi Singh, John D. Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, et al. Beyond human data: Scaling self-training for problem-solving with language models. *CoRR*, abs/2312.06585, 2023.
- [52] Fabrizio Gilardi, Meysam Alizadeh, and Ma.I Kubli. Chatgpt outperforms crowd-workers for text-annotation tasks. *CoRR*, abs/2303.15056, 2023.

- [53] R. Wang, W. Zhou, and M. Sachan. Let’s synthesize step by step: Iterative dataset synthesis with large language models by extrapolating errors from small models. *arXiv preprint arXiv:2310.13671*, 2023.
- [54] C. Cui, X. Peng, and M. Riedl. Thespian: Multi-character text role-playing game agents. *arXiv preprint arXiv:2308.01872*, 2023.
- [55] Abhay Zala, Jaemin Cho, Han Lin, Jaehong Yoon, and Mohit Bansal. Envgen: Generating and adapting environments via llms for training embodied agents. *arXiv preprint arXiv:2403.12014*, 2024.
- [56] Y. Yu, Q. Zhang, J. Li, Q. Fu, and D. Ye. Affordable generative agents. *arXiv preprint arXiv:2402.02053*, 2024.
- [57] Yu Sun, Xinhao Li, Karan Dalal, Jiarui Xu, Arjun Vikram, Genghan Zhang, Yann Dubois, Xinlei Chen, Xiaolong Wang, Sanmi Koyejo, Tatsunori Hashimoto, and Carlos Guestrin. Learning to (learn at test time): Rnns with expressive hidden states. *arXiv preprint arXiv:2407.04620*, 2024.
- [58] Wenhui Chen, Hexiang Hu, Xi Chen, Pat Verga, and William W. Cohen. Murag: Multimodal retrieval-augmented generator for open question answering over images and text. In *arXiv preprint arXiv:2210.02928*, 2022.
- [59] Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. Ra-cm3: Retrieval-augmented causal masked multimodal model. In *arXiv preprint arXiv:2211.12561*, 2022.
- [60] Daixuan Cheng, Shaohan Huang, Junyu Bi, Yuefeng Zhan, Jianfeng Liu, Yujing Wang, Hao Sun, Furu Wei, Denvy Deng, and Qi Zhang. Uprise: Universal prompt retrieval for improving zero-shot evaluation. In *arXiv preprint arXiv:2303.08518*, 2023.
- [61] Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. Raptor: Recursive abstractive processing for tree-organized retrieval. In *arXiv preprint arXiv:2401.18059*, 2024.
- [62] Lucrezia Bottegoni. Memwalker: Memory tree approach for long document question answering. In *arXiv preprint arXiv:2310.0502*, 2023.
- [63] H. Yang, Z. Lin, W. Wang, H. Wu, Z. Li, B. Tang, W. Wei, J. Wang, Z. Tang, S. Song, C. Xi, Y. Yu, K. Chen, F. Xiong, L. Tang, and W. E. Memory<sup>3</sup>: Language modeling with explicit memory. *arXiv preprint arXiv:2407.01178*, 2024.
- [64] 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.
- [65] Kunal Sawarkar, Abhilasha Mangal, and Shivam Raj Solanki. Blended rag: Improving rag accuracy with semantic search and hybrid query-based retrievers. In *arXiv preprint arXiv:2404.07220*, 2024.
- [66] Leonie Monigatti. Improving retrieval performance in rag pipelines with hybrid search. In *Towards Data Science*, 2024.
- [67] Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection. *arXiv preprint arXiv:2310.11511*, 2023.
- [68] Michael Glass, Gaetano Rossiello, Md Faisal Mahbub Chowdhury, Ankita Rajaram Naik, Pengshan Cai, and Alfio Gliozzo. Re2g: Retrieve, rerank, generate, 2022.
- [69] Yunjia Xi, Weiwen Liu, Jianghao Lin, Bo Chen, Ruiming Tang, Weinan Zhang, and Yong Yu. Memocrs: Memory-enhanced sequential conversational recommender systems with large language models, 2024.

- [70] 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*, 2023.
- [71] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Canton Ferrer, Cristian, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Singh Koura, Punit, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Michael Smith, Eric, Ranjan Subramanian, Ellen Tan, Xiaoqing, Binh Tang, Ross Taylor, Adina Williams, Xiang Kuan, Jian, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [72] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Canton Ferrer, Cristian, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Lewis Anderson, Georgia, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Silveira Cabral, Ricardo, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuwei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples,



Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Medina Florez, Gabriela, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind That-tai, Grant Herman, Grigory Sizov, Guangyi (Jack)Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Hou U, Kam, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelen, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabza, Manav Avalani, Manish Bhatt, Maria Tsimploukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Jubert Hermoso, Miquel, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Pavlovich Laptev, Nikolay, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Jayesh Bondu, Sai, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Cindy Zha, Shengxin, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiao Cheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu (Sid) Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

- [73] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Le Scao, Teven, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and El Sayed, William. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

- [74] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothee Lacroix, Baptiste Roziere, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Lamplem Guillaume. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [75] Jupinder Parmar, Sanjev Satheesh, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Reuse, don’t retrain: A recipe for continued pretraining of language models. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- [76] Yujia Qin, Jiajie Zhang, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. Elle: Efficient lifelong pre-training for emerging data. *Findings of the Association for Computational Linguistics*, 2022.
- [77] Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeonghun Kim, Stanley Choi Jungkyu, and Minjoon Seo. Towards continual knowledge learning of language models. *arXiv preprint arXiv:2110.03215*, 2022.
- [78] Pierre Colombo, Telmo Pessoa Pires, Malik Boudiaf, Dominic Culver, Rui Melo, Caio Corro, Andre F. T. Martins, Fabrizio Esposito, Vera Lucia Raposo, Sofia Morgado, and Desa Michael. Saullm-7b: A pioneering large language model for law. *arXiv preprint arXiv:2403.03883*, 2024.
- [79] Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhak Kim, and Bing Liu. Continual pre-training of language models. *arXiv preprint arXiv:2302.03241*, 2023.
- [80] Suchin Gururangan, Ana Marasovic, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don’t stop pretraining: Adapt language models to domains and tasks. *arXiv preprint arXiv:2302.03241*, 2023.
- [81] Kshitij Gupta, Benjamin Therien, Adam Ibrahim, Mats L. Richter, Quentin Anthony, Belilovsky Eugene, Irina Rish, and Timothee Lesort. Continual pre-training of large language models: How to (re)warm your model? *arXiv preprint arXiv:2308.04014*, 2023.
- [82] Genta Indra Winata, Lingjue Xie, Karthik Radhakrishnan, Shijie Wu, Xisen Jin, Pengxiang Cheng, Mayank Kulkarni, and Daniel Preotiuc-Pietro. Overcoming catastrophic forgetting in massively multilingual continual learning. *arXiv preprint arXiv:2305.16252*, 2023.
- [83] Adam Ibrahim, Benjamin Therien, Kshitij Gupta, Mats L. Richter, Quentin Anthony, Timothee Lesort, Eugene Belilovsky, and Irina Rish. Simple and scalable strategies to continually pre-train large language models. *arXiv preprint arXiv:2403.08763*, 2024.
- [84] Chengyue Wu, Yukang Gan, Yixiao Ge, Zeyu Lu, Jiahao Wang, Ye Feng, Ying Shan, and Ping Luo. Llama pro: Progressive llama with block expansion. *arXiv preprint arXiv:2401.02415*, 2024.
- [85] Renze Lou, Kai Zhang, and Wenpeng Yin. Large language model instruction following: A survey of progresses and challenges. *Computational Linguistics*, 2024.
- [86] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*, volume 2011, pages 27730–27744, 2022.
- [87] Zhengyan shi, Adam X. Yang, Bin Wu, Laurence Aitchison, Emine YilmaZ, and Aldo Lipani. Instruction tuning with loss over instructions. *arXiv preprint arXiv:2405.14394*, 2024.
- [88] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25, 2024.

- [89] Hao Li, Chenghao Yang, An Zhang, Yang Deng, Xiang Wang, and Tat-Seng Chua. Hello again! llm-powered personalized agent for long-term dialogue. *arXiv preprint arXiv:2406.05925*, 2024.
- [90] Kai Zhang, Fubang Zhao, Yangyang Kang, and Xiaozhong Liu. Memory-augmented llm personalization with short-and long-term memory coordination. *arXiv preprint arXiv:2309.11696*, 2023.
- [91] Zhichao Wang, Bin Bi, Shiva Kumar Pentyala, Kiran Ramnath, Sougata Chaudhuri, Shubham Mehrotra, Xiang-Bo Mao, Sitaram Asur, et al. A comprehensive survey of llm alignment techniques: Rlhf, rlaif, ppo, dpo and more. *arXiv preprint arXiv:2407.16216*, 2024.
- [92] Harrison Lee, Samrat Phatale, Hassan Mansoor, Kellie Lu, Thomas Mesnard, Colton Bishop, Victor Carbune, and Abhinav Rastogi. Rlaif: Scaling reinforcement learning from human feedback with ai feedback. *arXiv preprint arXiv:2309.00267*, 2023.
- [93] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.
- [94] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [95] Kawin Ethayarajh, Winnie Xu, Niklas Muennighoff, Dan Jurafsky, and Douwe Kiela. Human-aware loss functions (halos). <https://github.com/ContextualAI/HALOs/blob/main/assets/report.pdf/>, 2024. [Online; accessed 30-Aug-2024].
- [96] Jiwoo Hong, Noah Lee, and James Thorne. Orpo: Monolithic preference optimization without reference model. *arXiv preprint arXiv:2403.07691*, 2(4):5, 2024.
- [97] Yu Meng, Mengzhou Xia, and Danqi Chen. Simpo: Simple preference optimization with a reference-free reward. *arXiv preprint arXiv:2405.14734*, 2024.
- [98] Swaroop Nath, Tejpalsingh Siledar, Sankara Sri Raghava Ravindra Muddu, Rupasai Rangaraju, Harshad Khadilkar, Pushpak Bhattacharyya, Suman Banerjee, Amey Patil, Sudhanshu Shekhar Singh, Muthusamy Chelliah, and Nikesh Garera. Leveraging domain knowledge for efficient reward modelling in rlhf: A case-study in e-commerce opinion summarization. *arXiv preprint arXiv:2402.15473*, 2024.
- [99] Tzu-Han Lin, Chen-An Li, Hung-yi Lee, and Yun-Nung Chen. Dogerm: Equipping reward models with domain knowledge through model merging. *arXiv preprint arXiv:2407.01470*, 2024.
- [100] 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.
- [101] Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023.
- [102] Ofir Press, Noah A. Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables input length extrapolation. *International Conference on Learning Representations*, 2022.
- [103] Yira Ding, Li Lyna Zhang, Chengruidong Zhang, Yuanyuan Xu, Ning Shang, Jiahang Xu, Fang Yang, and Yang Mao. Longrope: Extending llm context window beyond 2 million tokens. *arXiv preprint arXiv:2402.13753*, 2024.
- [104] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*, 2023.

- [105] Dawei Zhu, Nan Yang, Liang Wang, Yifan Song, Wenhao Wu, Furu Wei, and Sujian Li. Pose: Efficient context window extension of llms via positional skip-wise training. *arXiv preprint arXiv:2309.10400*, 2024.
- [106] Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Longlora: Efficient fine-tuning of long-context large language models. *arXiv preprint arXiv:2309.12307*, 2024.
- [107] Wenhan Xiong, Jingyu Liu, Igor Molybog, Hejia Zhang, Prajjwal Bhargava, Rui Hou, Louis Martin, Rashi Rungta, Karthik Abinav Sankararaman, Barlas Oguz, Madian Khabsa, Han Fang, Yashar Mehdad, Sharan Narang, Kshitiz Malik, Angela Fan, Shruti Bhosale, Sergey Edunov, Mike Lewis, Sinong Wang, and Hao Ma. Effective long-context scaling of foundation models. *arXiv preprint arXiv:2309.16039*, 2023.
- [108] Yaru Hao, Yutao Sun, Li Dong, Zhixiong Han, Yuxian Gu, and Furu Wei. Structured prompting: Scaling in-context learning to 1,000 examples. *arXiv preprint arXiv:2212.06713*, 2022.
- [109] Nir Ratner, Yoav Levine, Yonatan Belinkov, Ori Ram, Inbal Magar, Omri Abend, Ehud Karpas, Ammon Shashua, Kevin Leyton-Brown, and Yoav Shoham. Parallel context windows for large language models. *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*, 1:6383—6402, 2023.
- [110] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Lllia Polosukhin. Attention is all you need. *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6000–6010, 2017.
- [111] 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, pages 19724–19731, 2024.
- [112] Lei Wang, Ma Chen, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, Xin Wayne Zhao, Zhewei Wei, and Ji-Rong Wen. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18, 2024.
- [113] 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. *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, 2023.
- [114] Lei Wang, Jingsen Zhang, Hao Yang, Zhiyuan Chen, Jiakai Tang, Zeyu Zhang, Xu Chen, Yankai Lin, Ruihua Song, Zhao Xin Wayne, Jun Xu, Zhicheng Dou, Jun Wang, and Ji-Rong Wen. User behavior simulation with large language model based agents. *arXiv preprint arXiv:2306.02552*, 2023.
- [115] Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. Character-llm: A trainable agent for role-playing. *arXiv preprint arXiv:2310.10158*, 2023.
- [116] Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen, Yi Song, Jifan Yu, Yongkang Huang, Libiao Peng, Jiaming Yang, Xiyao Xiao, Sahand Sabour, Xiaohan Zhang, Wenjing Hou, Yijia Zhang, Yuxiao Dong, Jie Tang, and Minlie Huang. Characterglm: Customizing chinese conversational ai characters with large language models. *arXiv preprint arXiv:2311.16832*, 2023.
- [117] Li Cheng, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi Mi, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun. Chatharuhi: Reviving anime character in reality via large language model. *arXiv preprint arXiv:2308.09597*, 2023.
- [118] Karan Singhal, Shekoofeh Azizi, Tao Tu, S. Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, Perry Payne, Martin Seneviratne, Paul Gamble, Chris Kelly, Nathaneal Scharli, Aakanksha Chowdhery, Philip Mansfield, Blaise Agüera y Arcas, Dale Webster, Greg S. Corrado, Yossi Matias, Katherine Chou, Juraj Gottweis, Nenad Tomasev, Yun Liu, Alvin Rajkomar, Joelle Barral, Christopher Semturs, Alan Karthikesalingam, and Vivek Natarajan. Large language models encode clinical knowledge. *arXiv preprint arXiv:2212.13138*, 2022.

- [119] Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, Mike Schaekermann, Amy Wang, Mohamed Amin, Sami Lachgar, Philip Mansfield, Sushant Prakash, Bradley Green, Ewa Dominowska, Blaise Agüera y Arcas, Nenad Tomasev, Yun Liu, Renee Wong, Christopher Semturs, S. Sara Mahdavi, Joelle Barral, Dale Webster, Greg S. Corrado, Yossi Matias, Shekoofeh Azizi, Alan Karthikesalingam, and Vivek Natarajan. Large language models encode clinical knowledge. *arXiv preprint arXiv:2305.09617*, 2023.
- [120] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling language modeling with pathways. *The Journal of Machine Learning Research*, 24(240), 2024.
- [121] Haoxi Zhong, Chaojun Xiao, CunChao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. How does nlp benefit legal system: A summary of legal artificial intelligence. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- [122] Quzhe Huang, Mingxu Tao, Chen Zhang, Zhenwei An, Cong Jiang, Zhibin Chen, Zirui Wu, and Yansong Feng. Lawyer llama technical report. *arXiv preprint arXiv:2305.15062*, 2023.
- [123] Jiayi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Yuan Li. Chatlaw: Open-source legal large language model with integrated external knowledge bases. *arXiv preprint arXiv:2306.16092v1*, 2023.
- [124] Shengbin Yue, Wei Chen, Siyuan Wang, Bingxuan Li, Chenchen Shen, Shunjun Liu, Yuxuan Zhou, Yao Xiao, Song Yun, Xuanjing Huang, and Zhongyu Wei. Disc-lawllm: Fine-tuning large language models for intelligent legal services. *arXiv preprint arXiv:2309.11325*, 2023.
- [125] Sijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhajan Kambadur, David Rosenberg, and Gideon Mann. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*, 2023.
- [126] Hongyang Yang, Xiao-Yang Liu, and Christina Dan Wang. Fingpt: Open-source financial large language models. *arXiv preprint arXiv:2306.06031*, 2023.
- [127] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Wang Haofen. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.
- [128] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652*, 2023.
- [129] Chengyu Huang, Zeqiu Wu, Yushi Hu, and Wenya Wang. Training language models to generate text with citations via fine-grained rewards. *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*, 2024.
- [130] Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzales. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023.
- [131] 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.

- [132] Yanyang Li, Shuo Liang, Michael R. Lyu, and Liwei Wang. Making long-context language models better multi-hop reasoners. *arXiv preprint arXiv:2408.03246*, 2024.
- [133] Junkai Li, Siyu Wang, Meng Zhang, Weitao Li, Yunghwei Lai, Xinhui Kang, Weizhi Ma, and Yang Liu. Agent hospital: A simulacrum of hospital with evolvable medical agents. *arXiv preprint arXiv:2405.02957*, 2024.
- [134] Hajra Murtaza, Musharif Ahmed, Naurin Farooq Khan, Ghulam Murtaza, Saad Zafar, and Ambreen Bano. Synthetic data generation: State of the art in health care domain. *Computer Science Review*, 48:100546, 2023.
- [135] Yuanyuan Liang, Jianing Wang, Hanlun Zhu, Lei Wang, Weining Qian, and Yunshi Lan. Prompting large language models with chain-of-thought for few-shot knowledge base question generation, 2023.
- [136] Chenhao Zhang, Renhao Li, Minghuan Tan, Min Yang, Jingwei Zhu, Di Yang, Jiahao Zhao, Guancheng Ye, Chengming Li, Xiping Hu, et al. Cpsycoun: A report-based multi-turn dialogue reconstruction and evaluation framework for chinese psychological counseling. *arXiv preprint arXiv:2405.16433*, 2024.
- [137] Junda Wang, Zonghai Yao, Zhichao Yang, Huixue Zhou, Rumeng Li, Xun Wang, Yucheng Xu, and Hong Yu. Notechat: A dataset of synthetic doctor-patient conversations conditioned on clinical notes, 2024.
- [138] World Health Organization et al. Icd-11 for mortality and morbidity statistics (2018). *website*, 2018.
- [139] DS American Psychiatric Association et al. *Diagnostic and statistical manual of mental disorders: DSM-5*, volume 5. American psychiatric association Washington, DC, 2013.
- [140] Michael B First. Structured clinical interview for the dsm (scid). *The encyclopedia of clinical psychology*, pages 1–6, 2014.
- [141] Pranav Rajpurkar, Robin Jia, and Percy Liang. Know what you don’t know: Unanswerable questions for squad. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Association for Computational Linguistics, 2018.
- [142] Xiao Liu, Hanyu Lai, Hao Yu, Yifan Xu, Aohan Zeng, Zhengxiao Du, Peng Zhang, Yuxiao Dong, and Jie Tang. Webglm: Towards an efficient web-enhanced question answering system with human preferences. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 4549–4560, 2023.
- [143] Kevin Clark and Christopher D. Manning. Deep reinforcement learning for mention-ranking coreference models. *arXiv preprint arXiv:1609.08667*, 2016.
- [144] K Lan, B Jin, Z Zhu, et al. Depression diagnosis dialogue simulation: Self-improving psychiatrist with tertiary memory. *arXiv preprint arXiv:2409.15084*, 2024.
- [145] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- [146] H. Ardel. *Le rêve de Suzy*. Plon, 1950.
- [147] Y. Jin. *Smiling Proud Wanderer*. Liu, 1993.
- [148] Zhili He and Yu-Hsing Wang. Mamba meets crack segmentation. *arXiv preprint arXiv:2407.15714*, 2024.
- [149] Bo Peng, Daniel Goldstein, Quentin Anthony, Alon Albalak, Eric Alcaide, Stella Biderman, Eugene Cheah, Xingjian Du, Teddy Ferdinan, Haowen Hou, Przemysław Kazienko, Kranti Kiran GV, Jan Kociń, Bartłomiej Koptyra, Satyapriya Krishna, Ronald McClelland Jr.,

- Niklas Muennighoff, Fares Obeid, Atsushi Saito, Guangyu Song, Haoqin Tu, Stanisław Woźniak, Ruichong Zhang, Bingchen Zhao, Qihang Zhao, Peng Zhou, Jian Zhu, and Rui-Jie Zhu. Eagle and finch: RwkV with matrix-valued states and dynamic recurrence. *arXiv preprint arXiv:2404.05892*, 2024.
- [150] Feng-Lei Fan, Yingxin Li, Hanchuan Peng, Tiejong Zeng, and Fei Wang. Towards neuroai: Introducing neuronal diversity into artificial neural networks. *arXiv preprint arXiv:2301.09245*, 2023.
- [151] L. He, Y. Xu, and W. et al. He. Network model with internal complexity bridges artificial intelligence and neuroscience. *Nat Comput Sci* 4, page 584–599, 2024.
- [152] Jason K. Eshraghian, Max Ward, Emre O. Neftci, Xinxin Wang, Gregor Lenz, Girish Dwivedi, Mohammed Bennamoun, Doo Seok Jeong, and Wei D. Lu. Training spiking neural networks using lessons from deep learning. *Proceedings of the IEEE*, 111(9):1016–1054, 2023.
- [153] Grégoire Mialon, Clémentine Fourier, Craig Swift, Thomas Wolf, Yann LeCun, and Thomas Scialom. Gaia: a benchmark for general ai assistants. *arXiv preprint arXiv:2311.12983*, 2023.
- [154] OpenAI. Introducing openai o1-preview. <https://openai.com/index/introducing-openai-o1-preview/>, 2024.
- [155] LlamaIndex. Llama parse. [https://docs.llamaindex.ai/en/stable/llama\\_cloud/llama\\_parse/](https://docs.llamaindex.ai/en/stable/llama_cloud/llama_parse/), 2024.

## A RTG prompt

Figure A shows the prompt we use to perform RTG.

```
{contexts}
{question}
```

Answer the question using the information given in the documents. You should respond in JSON format with four keys: 'Relevant': list[int] or 'None', 'Cited': list[int] or 'None', 'Grounded': str, and 'Answer': str. Carefully perform the following instructions, in order:

- Firstly, decide which of the documents are relevant to question. You should be thorough in curating the list to ensure all relevant information is included. You should answer with a comma-separated list of document numbers. If none are relevant, you should instead write 'None'. Use the JSON key: "Relevant".
- Secondly, decide which of the documents contain facts that should be cited in a good answer to the user's question. You should be thorough in curating the list to ensure all relevant information is included. You should answer with a comma-separated list of document numbers. If you don't want to cite any of them, you should instead write 'None'. Use the JSON key: 'Cited'.
- Thirdly, provide a comprehensive step-by-step reasoning on how to answer the question and cite all the 'Cited Documents' in your reasoning. In your reasoning, copy paste the cited text and include them in <begin\_cite: doc\_num> and <end\_cite: doc\_num> when the quotes come from a document. This would mean that things outside of <begin\_cite: doc\_num> and <end\_cite: doc\_num> are not directly copy paste from the documents. e.g. <begin\_cite: 0>my fact<end\_cite: 0> for a quote from Document 0. Use the JSON key: 'Grounded'.
- Finally, provide a detailed and complete answer to user's question in high quality natural English. Use your reasoning and cited documents to help you. Do not insert any citations or grounding markup. Use the JSON key: 'Answer'