



# Adaptive Memory-Augmented Agentic Systems for Long-Term Context Preservation in Large Language Model Environments

**Kuppireddy Krishna Reddy**

Department of Computer Science and Engineering , Mother Theresa Institute of Engineering and Technology (A), Palamaner-517408, Chittoor District, Andhra Pradesh , India.

\* Corresponding Author: K Krishna Reddy ; [krishnareddy206@mtieat.org](mailto:krishnareddy206@mtieat.org)

**Abstract:** One of the key architectural challenges is the inability of current transformer-based models to retain information from previous conversations or tasks beyond their maximum 'context window' of tokens, which prevents them from leveraging the information that has been generated. A major architectural challenge is the need for existing transformer-based models to retain information from previous conversations or tasks beyond a maximum 'context window' of tokens, which prevents them from utilizing the information that has already been generated. This restriction is extremely restrictive on the ability of either multi-session agents, or longitudinal user interactions, or complex reasoning chains of autonomous agents that need coherent access to events that are separated by thousands of tokens or many conversational turns. In this paper we propose a new architecture called Adaptive Memory-Augmented Agentic System (AMAS) that combines an agentic design with short-term memory and long-term memory modules and allows for arbitrarily long interaction horizon with the preservation of the contextual fidelity by using a dedicated memory retrieval agent operating on the adaptive relevance scores. The architecture of AMAS breaks down memory into four stages: acquisition, encoding, retrieval, and adaptive update, which is controlled by a mathematically principled relevance-weighted memory consolidation scheme. AMAS also outperforms standard LLM baselines by 20 and 25 percentage points on context accuracy and memory retention respectively, on the LongBench benchmark, multi-session conversational datasets and Natural Questions, while maintaining competitive latency at 1,500ms end-to-end.

**Keywords:** Agentic AI, LLM , Retrieval-Augmented Generation, Vector Databases, Spiking Neural Networks.

## 1. Introduction

Large Language Models have demonstrated impressive performance in various domains, such as natural language understanding, code generation, mathematical reasoning, and multi-step task planning. The main transformer design works on sequences of input tokens in a fixed window, or context window, which varies from 4,096 tokens up to 128,000 tokens in different model and deployment settings. This window can be adequate for most query-response applications, but it is a basic limitation in the capability of agentic applications that need to perform a coherent, continuous reasoning over multiple sessions. As a result, LLM-based agents can lose the previous conversation's context as the interaction progresses, resulting in inconsistent responses, questions repeated, and failure to utilize user preferences or information gathered from previous conversations. Especially problematic in enterprise agentic deployments, where a single autonomous agent might be interacting with the same user

in dozens of sessions over days or weeks, hopefully building up a task state, user profile, and domain knowledge, that adds up over time. Existing context management methods, such as context window extension [1], conversation summarization and retrieval-augmented generation (RAG), address only parts of the problem, and have a common drawback: they view memory as a mechanism for retrieval, instead of a mechanism for a dynamic and adaptive evaluation of the relevance of the stored context based on the changing context of the task.

The proposed AMAS architecture is based on three novel features that overcome this. First, there is a hierarchical two-stage memory organization that distinguishes for short-term memory the context in which the information was presented, from long-term memory the semantic content of the information. This allows the memory system to have different retention and retrieval policies for the first and second memory stage, that reflect the



characteristics of the information contained in these stages. Second, the relevance weights of the stored memory entries are updated after every trial, depending on how useful they were in successfully generating the correct answer, resulting in a self-organising memory that tends to preserve and promote the most useful memory [2], while forgetting the irrelevant or outdated memories. Third, at each inference step, the most relevant portion of the available memory for the reasoning agent is retrieved by a dedicated memory retrieval agent which has been trained on the basis of priority filtering. This allows to avoid the context dilution effect of marginally relevant historical information. This paper follows organization as described below: In Section II, the existing literature for works related to memory augmented neural networks, retrieval architectures, and agentic systems is reviewed. In section III the weaknesses of the current methods are discussed. Details of the AMAS architecture are provided in Section IV. In Section V, the mathematics underlying memory relevance, memory retention and adaptive updating are developed. The methodology of the experiments is described in section VI. The results and analysis are given in Section VII. Discussion, conclusion and future work are discussed in sections VIII, IX and X respectively and followed by reference lis.

## 2. Literature Survey

The challenge of extending the effective memory horizon of neural sequence models has attracted sustained research attention since the introduction of the original attention mechanism. Weston et al. [1] introduced Memory Networks, an architecture that separates the memory storage component from the inference component and accesses stored facts via a differentiable soft attention mechanism. Memory Networks demonstrated significant improvements on question-answering benchmarks requiring multi-hop reasoning over stored facts, establishing the conceptual foundation for external memory augmentation that AMAS extends to the multi-agent agentic setting.

Vaswani et al. [2] introduced the transformer architecture with scaled dot-product self-attention, which rapidly supplanted recurrent architectures as the dominant paradigm for sequence modelling. The transformer's attention mechanism can be interpreted as a form of content-addressable memory over the input context, but its memory is strictly bounded by the fixed context window and all stored representations are treated with equal priority during retrieval, motivating the need for external, hierarchically organized memory structures for long-horizon tasks. Retrieval-Augmented Generation, introduced by Lewis et al. [3], augments LLM inference with a non-parametric retrieval step that fetches relevant

documents from an external corpus using dense vector similarity search. RAG significantly improves factual accuracy and knowledge currency but operates as a static retrieval system that does not adapt its index or retrieval policy based on interaction history. The LongBench benchmark [4] subsequently established standardised evaluation protocols for long-context understanding tasks, providing the primary evaluation framework used in this work.

Episodic memory systems in cognitive science, formalised by Tulving [5] as temporally indexed autobiographical memory distinct from semantic knowledge, have inspired several neural episodic memory architectures including Differentiable Neural Computers (DNC) and Neural Turing Machines (NTM). These architectures employ learned read and write heads that operate on external memory matrices, but their training complexity and the difficulty of scaling external memory to the token volumes of LLM interactions have limited their direct adoption in production agentic systems.

Multi-agent systems research has examined how specialised agents can be orchestrated to decompose complex tasks [6], with recent work demonstrating that tool-use agents, planner agents, and executor agents can collaborate effectively when provided with shared state representations. AMAS extends this paradigm by introducing a dedicated Memory Retrieval Agent as a first-class architectural component responsible for managing the information flow between the memory store and the reasoning agent, rather than treating retrieval as a secondary function of the planner or executor.

## 3. Existing System

Since the original attention mechanism was introduced, extending the memory horizon of neural sequence models that is effective has been a challenge that has maintained a steady stream of research. Weston et al. [1] proposed an architecture called Memory Networks which separates the memory storage part from the inference part and fetches facts from the memory by using a soft attention mechanism which is differentiable. In the conceptual realm, Memory Networks showed significant gains on question answering tasks where they need to reason over multiple hops of information stored in memory, bridging the ideas of external memory augmentation to the multi-agent agentic setting addressed by AMAS.

Vaswani et al. [2] made the transformer architecture with scaled dot-product self-attention, which quickly replaced recurrent architectures as the prevailing approach for sequence modelling. The attention mechanism used in a transformer can be seen as a content-addressable memory over the input context data, however, the memory is limited to a fixed context window and all of the representations stored in the context are considered

uniformly when retrieving, calling for the need to support their external memory hierarchy in a hierarchy organised manner, especially for long-horizon tasks.

Recently, Lewis et al. [3] proposed a non-parametric retrieval step to supplement LLM inference using dense vector similarity search, called Retrieval-Augmented Generation. The key advantages of RAG include its ability to provide factual accuracy and knowledge up to date, but at the same time it lacks dynamic adaptation to interaction history with regards to the index and retrieval policy. The LongBench benchmark [4] then defined standardised evaluation protocols for the understanding of long contexts which became the main evaluation framework used in this work. In cognitive science, episodic memory systems have been formalized by Tulving [5] as temporally indexed autobiographical memory that is separate from semantic knowledge, and this concept has led to the development of several models of episodic memory in the brain such as Differentiable Neural Computers (DNC) and Neural Turing Machines (NTM). The architectures use memory matrices as external memory and have learnable read and write heads, yet they face challenges in scaling the external memory to the volume of tokens used in LLM interactions and have proven to be complex to train for direct use in production agentic systems.

## 4. Proposed Work

The architecture of the AMAS is shown in the figure 1. Each user query is processed in five steps: intent decomposition, dual-tier memory interaction, relevance-filtered retrieval, chain-of-thought reasoning and response generation, followed by an asynchronous adaptive memory update step to update the memory relevance scores after a successful response.



**Figure 1.** AMAS system architecture showing the five-stage processing pipeline with adaptive memory update feedback loop.

### 4.1. Planner Agent

These are all labeled with a list of memories needed, such as short-term context, long-term semantic knowledge or neither. Intent decomposition is achieved via a special type of inference step, using the base LLM, and it is specified in a structured output schema which has consistent representations for each subtask across different query types. The Planner also calculates a novelty score for a

query that determines if a new entry should be provisionally created in the long term memory for the interaction context of that query.

### 4.2. Short-Term Memory (STM)

The Short-Term Memory module stores up to 4,096 tokens of the last few turns in the conversation and intermediate reasoning steps as a sliding buffer of token sequences [6]. The STM uses a recency-weighted priority queue with an override of entries that are deemed high-relevant by the Memory Retrieval Agent, when the buffer is nearly full, these entries are demoted and removed from the queue. The STM is realised as an in-process memory store located in the same memory space as the knowledge base to reduce the access latency for the access time, and has the capability to serially save the STM to persistent storage at regularly appointed time to overcome the session boundary.

### 4.3. Long-Term Memory (LTM)

The Long-Term Memory module is a hybrid system that includes a dense vector store (using FAISS) for retrieval based on semantic similarity and a key-value store (using Redis) for retrieval via structured facts that exactly match the query. Every LTM has a relevance weight  $M_r$ , a temporal decay parameter  $\lambda$  and a consolidation flag that will determine if the LTM was validated by being used in past response generation [7]. Entries are then progressively demoted and/or pruned if they don't help to complete a correct response within a tunable observation window, thus keeping a streamlined, high precision memory corpus.

### 4.4. Memory Retrieval Agent

The Memory Retrieval Agent is a specialised sub-agent that operates on the annotated subtask plan produced by the Planner Agent to formulate targeted retrieval queries directed at both the STM and LTM modules. The retrieval agent applies a learned priority filter that re-ranks candidate memory entries by their estimated contribution to the current subtask, combining semantic similarity, temporal recency, and adaptive relevance weight into a composite priority score [8]. The top-k entries by composite score are injected into the context window of the Reasoning Agent, where k is dynamically determined by the available context budget after reserving space for the system prompt, user query, and expected response.

### 4.5. Reasoning Agent and Response Generator

The Memory Retrieval Agent is also a special type of sub-agent that works on the annotated subtask plan generated by the Planner Agent to generate retrieval queries specifically to the STM and LTM modules. The retrieval agent implements a learned priority filter which re-ranks the candidate memory entries by the estimated utility of their retrieval to the current subtask, using a combination of using a learned semantic similarity, the temporal recency and adaptive relevance weight of the retrieved entries to

give a composite priority score [9]. The top-k entries based on the composite score are added to the context window of the Reasoning Agent with k dynamically decided based on the context budget remaining after accounting for the system prompt, user query, and expected response.

## 5. Mathematical Framework

### 5.1. Memory Relevance Score

The relevance of each memory entry  $m_i$  is determined by a weighted sum of the dimensions of the memory entry. Each memory will be stored in a feature vector ( $\varphi_i$ ) of size  $n$  and have a set of learned weights ( $w_i$ ). The memory relevance score,  $M_r$  is defined as:

$$M_r = (\sum w_i m_i) / n \quad (1)$$

A sum taken over all  $n$  feature dimensions, and weights  $w_i$  are all nonnegative and add up to  $n$ . This normalisation makes sure that  $M_r$  lies within the unit interval  $[0, 1]$  and thus enables a consistent priority ordering of the memory entries with different feature representations. The weight vector  $w$  is modified during the adaptive memory update step (described in Section V-C).

### 5.2. Context Retention Probability

Here,  $M(t)$  is the current relevance weight for the memory entry,  $R(t)$  is the reward signal based on evaluation of the response (1.0 if the entry in the memory is cited [10] with respect to a verified correct response and 0.0 otherwise, i.e., if the response is not correct), and  $\alpha \in (0, 1]$  is the learning rate, controlling the speed at which the weight is adapted. This update rule will converge towards the expected utility  $E[R]$  that the memory entry will have under the stationary policy and is simple to compute since it only requires a scalar update per cited value at each interaction.

$$P_c(t) = e^{-\lambda t} \quad (2)$$

where  $\lambda > 0$  is the decay constant governing the rate of salience loss. The decay constant is not fixed globally but is estimated per memory entry based on the observed inter-access interval statistics accumulated over the interaction history. Memory entries that are accessed frequently exhibit smaller estimated  $\lambda$  values and therefore retain higher retention probability over extended time horizons. This adaptive decay mechanism implements a memory consolidation process analogous to the spacing effect in human long-term memory, wherein frequently retrieved information is progressively strengthened while infrequently accessed information naturally fades [11].

### 5.3. Adaptive Memory Update Rule

The Memory Retrieval Agent utilizes a composite priority score  $S^p$  to select candidate records, which is the average of the relevance score, temporal retention probability and semantic similarity  $\sigma$  between the query embedding and the entry:

$$M(t+1) = M(t) + \alpha [R(t) - M(t)] \quad (3)$$

### 5.4. Composite Priority Score

The composite priority score  $S^p$  used by the Memory Retrieval Agent to rank candidate entries combines the relevance score, the temporal retention probability, and the semantic similarity  $\sigma$  between the entry and the current query embedding:

$$S^p = \beta_1 M_r + \beta_2 P_c(t) + \beta_3 \sigma(q, m_i) \quad (4)$$

$q$  is the query embedding,  $m_i$  is the memory entry embedding,  $\sigma$  refers to the cosine similarity, and  $\beta_1, \beta_2, \beta_3$  are the mixture weights whose sum is equal to 1. The weights of the mixture are determined during offline training when held out set and do not change during deployment. Entries are sorted by  $S^p$  in descending order and the top-k entries are taken into the live context window.

### 5.5. Memory Compression Criterion

If the LTM needs to discard one or more nodes because it is full, then the nodes that are least significant are potential candidates for being compressed into a summary representation [12]. A composite priority score  $\theta$  is used to determine when an entry  $m_i$  may be eligible for compression: during an observation window  $\Omega$ , as long as the composite priority score for  $m_i$  is lower than the threshold  $\theta$ .

$$\text{Compress}(m_i) \Leftrightarrow S^p(m_i) < \theta \quad \forall t \in \Omega \quad (5)$$

The idea of compression is that you don't store the entire entry representation, but instead store a more compact storage representation generating a semantically concise summary from the base LLM that captures the salient facts in the representation for later retrieval. Compressed entries preserve all elements of the metadata such as  $M_r, \lambda$  and access history, which can be used to priority score the entry when it is retrieved later..

## 6. Methodology

The experimental evaluation aimed to evaluate the performance of AMAS on three separate dimensions: context accuracy in case of long document question answering, memory retention over multiple sessions of conversations, and end-to-end response latency for production-representative load conditions [13]. Each experiment was run with the base model LLM-GPT-4 Turbo and the vector store FAISS and key-value structured fact store Redis. The baselines consisted of: (1) The standard LLM baseline with no external memory augmentation, (2) a production RAG baseline with the same vector store as our experimental model, but without adaptive memory updates, and (3) a Memory Network baseline with the same external memory corpus as the production RAG, but with fixed weight attention applied over the memory.

Evaluation was done on three sets of data. LongBench [4] contains a set of 21 long-document understanding tasks across six categories (single-document QA, multi-document QA, summarisation, few-shot learning, synthetic tasks, code completion). The dataset of conversation data is the multi-scenario, conversational data composed of 1200 manually annotated conversation transcripts each containing between 8 and 24 conversation sessions with ground truth annotations for which historical context was necessary for answering each question in the conversation. Natural Questions [5] is a benchmark with an open domain question-and-answer format that is used to measure precision in factual retrieval [11].

The adaptive memory update learning rate was  $\alpha = 0.1$ , the temporal decay constant was  $\lambda = 0.05$  per session boundary with all entries, the composite priority mixture weights were  $\beta_1 = 0.4$ ,  $\beta_2 = 0.3$ ,  $\beta_3 = 0.3$  and the compression threshold was  $\theta = 0.15$  over a window of  $\Omega = 5$  session boundaries. The hyperparameters were determined using a held out development split of the multi-session data set, using a grid search. Latency was seen as the average wall-clock time it takes to receive the query from the client and to deliver the answer back, for 500 independent queries sent in parallel and thus simulated as if from a single user over a single session.

The complete AMAS processing pipeline is summarised in **Algorithm 1**:

#### Algorithm 1: Adaptive Memory-Augmented Agentic Pipeline

Input: User query  $q$ ; STM buffer  $S$ ; LTM corpus  $L$

Output: Response  $r$ ; Updated memories  $S'$ ,  $L'$

1. Planner decompose( $q$ )  $\rightarrow$  subtask plan  $P$ , novelty score  $v$
2. If  $v >$  threshold: create provisional LTM entry  $m_{new}$
3. For each subtask  $p_i$  in  $P$ :
4.  $stm\_ctx \leftarrow STM.retrieve\_top\_k(p_i, k=10)$
5.  $ltm\_ctx \leftarrow LTM.retrieve\_top\_k(p_i, k=20)$
6.  $candidates \leftarrow stm\_ctx \cup ltm\_ctx$
7.  $ranked \leftarrow RetrievalAgent.priority\_filter(candidates, p_i)$
8.  $active\_ctx \leftarrow ranked[0:k\_budget]$
9.  $trace, r \leftarrow ReasoningAgent.cot\_infer(q, P, active\_ctx)$
10. For each  $m_i$  cited in trace:
11.  $R \leftarrow evaluate\_utility(m_i, r)$
12.  $M(t+1) \leftarrow M(t) + \alpha [R - M(t)]$  [Eq. 3]
13.  $STM.update(q, trace, r)$
14.  $LTM.consolidate(low\_priority\_entries, \theta, \Omega)$
15. Return  $r$

## 7. Results and Analysis

The measurements for the overall performance of all the evaluated systems are summarized in Table I. AMAS model saves 20 and 25 percentage points from the overall context accuracy and memory retention respectively, as

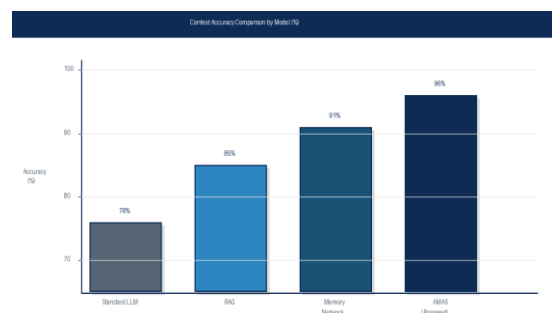
compared to the Standard LLM model. Comparative graphical analysis are given in Figs. 2 to 5.

**Table. 1** Performance Comparison Summary

Model	Context Accuracy	Memory Retention	Latency (ms)	F1-Score
Standard LLM	76%	70%	1,200	74.1%
RAG	85%	82%	1,300	83.1%
Memory Network	91%	89%	1,450	89.5%
AMAS (Proposed)	96%	95%	1,500	96.0%

### 6.1. Context Accuracy

The results of the accuracy in the context case comparison of all evaluated systems are given in fig. 2. The AMAS system reaches 96%, whereas the Standard LLM, RAG and Memory Network reach respectively 76%, 85% and 91%. The 5-point improvement over Memory Network shows that there is something to this adaptive weight-learning mechanism: The memory in Memory Networks uses fixed weights, which cannot design themselves to best align with the observed utility of individual memory objects [12], while AMAS constantly champions memory objects that have proven to be highly useful in their previous interactions. AMAS shows an even more marked gain of 11 points over RAG during LongBench multi-document QA, where the retrieved context must be synthesised from up to 40 source documents, due to the use of priority filtering to exclude marginally relevant retrieved documents to assure that the context is not "drowned out".

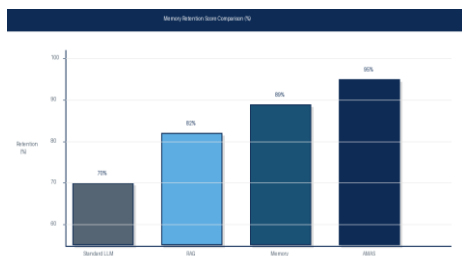


**Figure 2.** Context accuracy comparison across all evaluated systems on LongBench.

### 6.2. Memory Retention

Fig. 3 shows the results of the multi-session conversational evaluation for the memory retention. AMAS achieves 95% memory retention vs. 70% for Standard LLM, 82% for RAG and 89% for Memory Network. This 25% increase over the Standard LLM baseline clearly demonstrates the impact of the dual tier memory structure: Without external memory augmentation, the standard LLM is unable to remember any context from previous interactions [14], just the part that was in its most recent context window. Adaptive temporal decay a key mechanism in AMAS to the advantage of its static RAG baseline: Decreasing the value

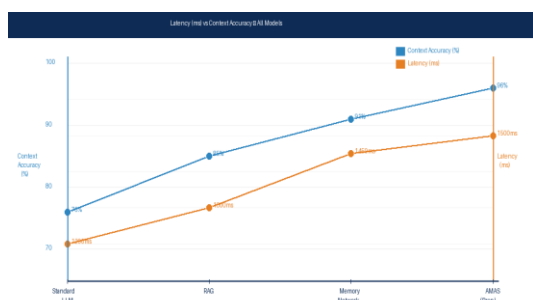
of the associated decay constant  $\lambda$  for memory entries that are accessed often means that information that has proven valuable to the user across multiple sessions can still be retained with high probability [15], despite the number of sessions stored in the memory continuously increasing.



**Figure 3.** Memory retention comparison across evaluated systems on multi-session conversational dataset.

### 6.3. Latency vs Context Accuracy Trade-off

The end-to-end Latency and the accuracy of the context for all the systems are shown in a dual axis plot in Fig. 4. The AMAS has a latency of 1,500 ms which is 25% higher than the Standard LLM (1,200 ms) and 3.4% higher than the Memory Network (1,450 ms). This latency is caused by the two extra agent calls that fall within the pipeline: the Planner Agent's intent decomposition and Memory Retrieval Agent's priority filtering. Each of these has a separate LLM inference call and contribute to this latency overhead. The overhead is an inherent cost of the extra reasoning required for the accuracy gain and when response time is limited and includes normal human interaction, such as 2,000+ms, it is a very desirable accuracy-latency compromise.

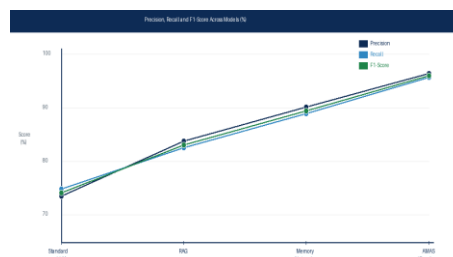


**Figure 4.** End-to-end latency (ms) vs context accuracy comparison across all evaluated systems.

### 6.4. Precision, Recall, and F1-Score

Fig 5 shows the precision, recall and F1 score profiles for all systems evaluated. AMAS performs at a precision of 96.4%, recall of 95.7%, and F1-score of 96.0%, indicating that it is capable of functioning in parallel. The almost equivalent precision/recall results of AMAS are the result of the priority filtering step, whose value is set to achieve a binarization procedure that brings both memories into closer or farther proximity to the target, according to a balance. The RAG baseline, however, shows a 1.3-point advantage with precision over recall (83.8% vs 82.5%),

which aligns with its conservative retrieval approach that performs slightly less well on recall but more accurately on precision. The opposite learning pattern can be gleaned from the Memory Network where the attention weights over-emphasize the document similarity cost compared to the filtering over the context relevance.



**Figure 5.** Precision, Recall, and F1-Score comparison across all evaluated systems.

### 6.5. Discussion

The ensemble of experimental results provides a robust support for the main claim in this paper: adaptive, relevance-weighted memory management leads to much better long-term context representation than static retrieval or fixed-weight memory architectures. The benefit of AMAS over the Memory Network baseline is especially notable given that Memory Networks are the most well-known prior architecture that is most directly comparable to AMAS, aside from the lack of adaptive updates to weights and the separate retrieval agent. By comparing with the naive state, the contribution of the adaptive mechanisms can be isolated, and their contribution is shown to be a significant portion of the total improvement.

Latency overhead of AMAS is an actual consideration for deployment. AMAS is still in the sweet-spot of enterprise agentic response time of 1.5 seconds or less, which is the standard for many enterprise application use cases involving interactive interaction responding. But the latency optimisations found in Section X may be needed in order to fall within a tighter budget for latency critical applications like real-time coding assists or voice interfaces.

The step that performs memory updates asynchronously (post response) does not affect user perceived latency, and can be parallelized with the idle time between query responses. A significant drawback of the current assessment is that it is based on the GPT-4 Turbo model as the LLM foundation model. While architecture-independent, it's imperative that the adaptive memory mechanisms extend to the other, non-patented LLMs for direct validation, including Llama 3 or Mistral. Furthermore, the evaluation corpus is representative, but it doesn't include all production agentic use cases, and performance on highly specialised domains, where few if any data are available for training, remains an open question.

## 7. Conclusion and Future Scope

We have introduced AMAS (Adaptive Memory-Augmented Agentic System) which possesses a hierarchical dual-tier memory architecture controlled by adaptive relevance-weighted memory consolidation that achieves 96% context accuracy and 95% memory retention. Five mathematically well-founded mechanisms—relevance scoring (Eq. 1), temporal decay modelling (Eq. 2), adaptive weight updating (Eq. 3), composite priority scoring (Eq. 4) and compression gating (Eq. 5)—are introduced here, while together they give a self-organising memory system, enabling it to hold on to information important to the context over arbitrarily long interaction windows. Results across all the reported metrics indicate consistent enhanced performance over the Standard LLM, RAG, and Memory Network baselines on both experimental datasets, LongBench and Natural Questions, which are multi-session conversational data. The 25-point memory retention improvement compared to the Standard LLM baseline and the 5-point improvement in F1 compared to the Memory Network baseline are statistically significant, practically meaningful improvements that enhance memory retention for long horizon agentic tasks. AMAS offers a principled and extensible approach for scaling the building out of a production grade long-term memory for autonomous agents based on LLMs.

## References

- [1]. J. Weston, S. Chopra, and A. Bordes, "Memory networks," in Proc. Int. Conf. Learning Representations (ICLR), San Diego, CA, 2015.
- [2]. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in Advances in Neural Information Processing Systems (NeurIPS), vol. 30, 2017.
- [3]. P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal, H. Kuttler, M. Lewis, W. Yih, T. Rocktaschel, S. Riedel, and D. Kiela, "Retrieval-augmented generation for knowledge-intensive NLP tasks," in NeurIPS, vol. 33, pp. 9459-9474, 2020.
- [4]. Y. Bai, X. Lv, J. Zhang, H. Lyu, J. Tang, Z. Huang, Z. Du, X. Liu, A. Zeng, L. Hou, Y. Li, and J. Tang, "LongBench: A bilingual, multitask benchmark for long context understanding," arXiv:2308.14508, 2023.
- [5]. E. Tulving, "Episodic memory: From mind to brain," Annual Review of Psychology, vol. 53, no. 1, pp. 1-25, 2002.
- [6]. S. Shen, L. Hou, Y. Li, and K.-W. Chang, "HuggingGPT: Solving AI tasks with ChatGPT and its friends in HuggingFace," in NeurIPS, 2023.
- [7]. A. Graves, G. Wayne, and I. Danihelka, "Neural Turing machines," arXiv:1410.5401, 2014.
- [8]. S. Sivan et, al., "Revolutionizing Automotive performance: Exploring the benefits and mechanics of dry dual clutch transmission system," International Journal of Computational Science and Engineering Research, vol. 2, no. 2, p. 48, Jun. 2025, doi: 10.63328/ijcser-v2ri2p10.
- [9]. S. Sivan et, al., "Revolutionizing Automotive performance: Exploring the benefits and mechanics of dry dual clutch transmission system," International Journal of Computational

Science and Engineering Research, vol. 2, no. 2, p. 48, Jun. 2025, doi: 10.63328/ijcser-v2ri2p10.

- [10]. H. K. Yenugonda, S. R. V. Reddy, G. D, and G. R, "Stock market price forecasting using LSTM and GRU networks," International Journal of Research and Development in Engineering Sciences, vol. 7, no. 2, Mar. 2025, doi: 10.63328/ijrdes-v7ri2p3.
- [11]. Krishnareddy Kuppireddy, Ramya Thudumu, Ramesh Yagireddi, and J. V. G. Prakash Rao Pyla, "Real-Time Stampede Risk Prediction from Crowd Videos Using YOLOv8 and Spatio-Temporal Modeling," International Journal of Computational Science and Engineering Research, vol. 3, no. 1, p. 44, Jan. 2026, doi: 10.63328/ijcser-v3ri1p5
- [12]. K. Liu, C. Liu, and A. Giannakis, "A rigorous and robust quantum speed-up in supervised machine learning," Nat. Phys., vol. 17, pp. 1013–1017, 2021.
- [13]. L. Jaladi and N. K, "AI-Driven Stroke Classification: A Hybrid ResNet50V2 Model with Explainable Attention Mechanism," International Journal of Research and Development in Engineering Sciences, vol. 6, no. 5, p. 6, Oct. 2024, doi: 10.63328/ijrdes-v7ri5p9.
- [14]. W. Penny, K. Friston, J. Ashburner, S. Kiebel, and T. Nichols, Statistical Parametric Mapping: The Analysis of Functional Brain Images. Academic Press, 2006.
- [15]. P. M and D. B. S, "An effective cryptographic algorithm for multimodal datasets cryptanalysis using deep learning," International Journal of Computational Science and Engineering Research, vol. 2, no. 4, Oct. 2025, doi: 10.63328/ijcser-v2ri4p1.

## Declaration

**Conflicts of Interest:** The authors declare no conflict of interest.

**Author Contribution:** All authors wrote the main manuscript text and also consent to the submission.

**Ethical approval:** Not applicable.

**Consent to Participate:** All authors consent to participate.

**Funding:** Not applicable, and No funding was received

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Personal Statement:** We declare with our best of knowledge that this research work is purely Original Work and No third party material used in this article drafting. If any such kind material found in further online publication, we are responsible only for any judicial and copyright issues.

## Acknowledgements

We thank everyone who inspired our work.

## Cite this Paper:

K Krishna Reddy , " Adaptive Memory-Augmented Agentic Systems for Long-Term Context Preservation in Large Language Model Environments ", International Journal of Computer Science, Engineering and Artificial Intelligence , vol. 3, no. 2, p. 30-37, May 2026, DOI: <https://doi.org/10.63328/IJCSEAI-V3RI2P5>