
MAGS: Memory Augmented Graph Scaling

Pulkith Paruchuri
University of Pennsylvania

Pragya Singh
University of Pennsylvania

Stanley Yu
University of Pennsylvania

Shailesh Kumar
University of Pennsylvania

Abstract

Current Large Language Models (LLMs) lack the critical ability to learn continuously, relying instead on ephemeral context windows that scale poorly, suffer from attention decay, and cannot retain knowledge between sessions. We propose the Memory Augmented Generative System (MAGS), a neuroscience-inspired architecture introducing dynamic memory modules to endow LLMs with human-like learning capabilities. MAGS integrates dual memory blocks: an Episodic Memory Block (eMB) for short-term, session-specific recall, and a Semantic Memory Block (sMB) for long-term, continuously updated knowledge. Memory updates occur via a reinforcement learning framework using Group Relative Policy Optimization (GRPO) and Low Rank Adaptation (LoRA), enabling models to autonomously modify and refine their knowledge graphs. Drawing inspiration from Hebbian plasticity and neurogenesis, MAGS strengthens memory paths that contribute positively to responses while pruning low-importance connections through k-core decomposition and a Game-of-Life-inspired algorithm. A novel recall and engram mechanism guides memory retrieval and updates, supported by explainable, inspectable graph structures. We validate MAGS through a custom-designed game environment that tests adaptability, recall accuracy, and memory update fidelity under evolving conditions. This architecture offers a scalable path toward lifelong learning in LLMs, enabling efficient, explainable, and robust memory systems without the need for costly retraining.

1 Introduction

Large Language Models (LLMs) have emerged as transformative tools in natural language processing (NLP), demonstrating human-level performance on a wide variety of tasks including translation, summarization, question answering, and code generation. Driven by massive pretraining on web-scale corpora, these models capture statistical regularities of language that enable fluent and contextually appropriate generation. However, despite impressive fluency and versatility, contemporary LLMs remain fundamentally stateless: once training concludes, they cannot incorporate new experiences or update their internal representations without expensive and time-consuming retraining. This limitation manifests in two critical ways. First, LLMs rely on context windows—sliding buffers of recent tokens—to simulate memory, but this mechanism suffers from quadratic scaling of compute cost, attention decay that degrades recall of earlier context, and ultimately ephemeral knowledge that vanishes at the end of a session. Second, even advanced augmentation strategies such as Retrieval Augmented Generation (RAG) provide only a static information injection: the model fetches relevant documents or facts at inference time, but has no ability to consolidate or modify that knowledge in subsequent interactions.

In contrast, biological cognition thrives on continuous learning and memory consolidation. Humans form both episodic memories—detailed accounts of specific experiences—and semantic memories—generalized factual knowledge—and dynamically integrate these stores to guide behavior. Memory formation involves mechanisms of neurogenesis (adding new neurons), Hebbian plasticity (strengthening connections between co-activated neurons), and synaptic pruning (eliminating weak or redundant connections). These processes yield a highly adaptive, self-organizing system capable of lifelong learning without catastrophic forgetting. Inspired by these principles, we propose the Memory Augmented Generative System (MAGS), an architecture that imbues LLMs with analogous memory capabilities by externally coupling them to dynamic graph-structured memory modules.

MAGS comprises two complementary memory blocks: the Episodic Memory Block (eMB) and the Semantic Memory Block (sMB). Both blocks are realized as Liquid Knowledge Graphs (LKGs)—dynamic graphs whose nodes represent discrete memory units (e.g., triplets of information) and whose weighted edges capture associations. The eMB captures short-term, session-specific experiences (such as a user’s preferences or recently learned facts), while the sMB consolidates long-term, widely applicable knowledge. Unlike static RAG indices, LKGs support real-time modifications: nodes and edges can be added, updated, or pruned during inference, enabling the model to organically integrate new information. Edge weights decay over time but can be reinforced when memory paths contribute meaningfully to successful outputs, echoing Hebbian learning.

To operationalize memory recall and update, MAGS employs a two-step engram-and-recall protocol. During the recall step, the LLM is prompted to emit a structured “recall” query that identifies relevant memory anchors in the LKG. We implement a multi-stage retrieval algorithm: first, candidate anchor nodes are selected based on semantic similarity to the input; next, graph algorithms (e.g., convex-hull selection) ensure broad coverage of the memory graph; finally, depth-first search (DFS) along decaying-weight edges retrieves associated context until a relevance threshold is breached. The retrieved subgraph is then injected into the model’s context to inform generation. In the engram step, after producing its response, the LLM signals which new information to store and provides confidence scores for each candidate memory. Guided by Group Relative Policy Optimization (GRPO) and fine-tuned via Low-Rank Adaptation (LoRA), the model learns when to create new nodes, update existing ones, and adjust edge weights to reinforce valuable associations or prune low-utility connections. MAGS’s dynamic memory modules offer several key advantages. First, by externalizing memory as separate graph structures, we avoid the need for expensive full-model retraining: the base LLM remains frozen, while memory adaptations occur in the LKG. Second, the separable, inspectable nature of LKGs facilitates explainability and debugging: practitioners can visualize and manipulate memory contents directly, audit update histories, and roll back erroneous additions. Third, controlled pruning via k-core decomposition and importance-based thresholds ensures that memory graphs remain computationally tractable, preventing unbounded growth. Fourth, because memory is invoked on demand through structured recall prompts, MAGS can scale to arbitrarily long histories without incurring quadratic attention cost: only relevant subgraphs enter the context window.

We validate MAGS in a custom game environment designed to stress-test continual learning. The game comprises evolving rules and scenarios that the LLM cannot deduce through reasoning alone, forcing reliance on external memory. Across multiple episodes, MAGS demonstrates significantly improved recall accuracy, response coherence, and adaptability compared to both vanilla LLMs and LLMs augmented with static RAG indices. Furthermore, ablation studies highlight the critical role of dynamic edge-weight updates and pruning in maintaining memory relevance over extended interactions.

2 Related Work

Continuous memory and lifelong learning have long been central challenges in both cognitive science and artificial intelligence. In the context of Large Language Models (LLMs), a rich body of work has explored external memory augmentation, continual learning, and graph-based retrieval to overcome the limitations of stateless inference. We organize the related literature into three primary categories:

(1) retrieval-based augmentation, (2) continual learning and parameter-efficient adaptation, and (3) graph-structured memory and dynamic memory architectures.

2.1 Retrieval-Based Augmentation

Retrieval Augmented Generation (RAG) methods[1] interpose an external information store—typically a vector-indexed document collection—between the user query and the LLM. During inference, a retriever module locates the top-k most relevant passages, which are then concatenated with the prompt to ground the model’s generation. This strategy markedly improves factual accuracy and updates speed compared to re-training the base model. However, RAG stores knowledge in a static manner: retrieved content is never consolidated back into the index, nor can the LLM modify the underlying store in response to new evidence or corrections. Variants such as REALM[2] and Fusion-in-Decoder[3] refine retrieval architectures or fuse retrieved chunks more effectively, yet they remain unidirectional: knowledge flows from the index into the model, but not vice versa.

2.2 Continual Learning and Parameter-Efficient Adaptation

To endow LLMs with the ability to internalize new information over time, continual learning research focuses on updating model parameters without catastrophic forgetting. Regularization-based methods (e.g., EWC[4]) constrain important weights to prevent overwriting, while replay-based techniques interleave new data with stored exemplars. Though effective in smaller architectures, these approaches become prohibitively expensive for billion-parameter models. More recent techniques apply parameter-efficient fine-tuning, such as Low-Rank Adaptation (LoRA)[5], adapter modules[6], or prompt tuning[7], which add a small set of tunable parameters atop a frozen base model. Google’s Retentive Memory Transformers (RMTs)[8] embed a persistent memory bank adjacent to attention layers, enabling the model to recall past activations. However, RMTs still suffer from fixed memory size and lack mechanisms for association or explicit consolidation.

2.3 Graph-Structured Memory and Dynamic Architectures

Graph-based memory representations offer a more flexible substrate for associating discrete facts and supporting relational reasoning. Early work in semantic networks[9] and Knowledge Graphs motivated inquiry into graph-augmented neural architectures. In the LLM era, MIT-IBM’s CAMELoT[10] integrates external key-value memory blocks that the model reads from and writes to via attention. While CAMELoT demonstrates effective episodic recall, its memory blocks are of fixed capacity and do not support explicit relationships between entries or reconcile conflicting facts. Ariadne[11] extends RAG by constructing dynamic graph indices that link related documents, yet it does not permit the LLM to update or prune the graph during inference.

More recent efforts explore liquid or dynamic graph structures that mirror neural plasticity. Graph Neural Networks (GNNs)[12] provide foundational algorithms for message passing and edge weighting, but classical GNNs assume static graphs. Dynamic GNN models[13][14] allow edges and node features to evolve over time, though they are rarely integrated with LLMs. Concurrently, the idea of neurogenesis—adding new nodes for novel concepts—and Hebbian plasticity—updating edge strengths based on co-activation—has inspired approaches such as Memory-Augmented Neural Networks (MANNs)[15] and episodic controllers[16]. Yet these systems are often task-specific and lack seamless integration into generative language pipelines.

2.4 Gaps and Opportunities

Despite these advances, critical gaps remain. Most retrieval-based systems cannot autonomously consolidate new knowledge or resolve contradictions without external retraining. Continual learning methods either impose heavy compute overhead or limit the model’s capacity to adapt. Graph-based memory architectures offer promise but have not been fully exploited for dynamic, lifelong language understanding: existing models either treat graphs as static indices or lack principled mechanisms for pruning, importance scoring, and consolidation over lengthy interactions.

MAGS addresses these shortcomings by unifying dynamic graph-structured memory with reinforcement learning and parameter-efficient adaptation. By externalizing memory into two Liquid Knowledge Graphs—an episodic block for short-term experiences and a semantic block for long-term

3.2 Engram-and-Recall Protocol

MAGS leverages a structured two-step interaction between the LLM and memory:

1. Recall Step

- **Trigger:** Upon receiving a user query x , the LLM is prompted via a special “RECALL:” token to generate a short, structured retrieval request r .
- **Anchor Selection:** Compute the embedding $\phi(x)$ via the base LLM’s encoder. Retrieve the top N candidate nodes in \mathcal{G}_* whose node embeddings $\phi(v)$ have the highest cosine similarity to $\phi(x)$.
- **Convex-Hull Filtering:** From the N candidates, solve a small convex-hull maximization problem to select $K < N$ anchors that maximize coverage of the local embedding space.
- **DFS Expansion:** Perform depth-first search from each anchor over edges ($u \rightarrow v$) in descending weight order, halting at depth D or whenever the cumulative product of weights along the path drops below threshold θ_r .
- **Context Injection:** Extract the subgraph nodes reached by DFS, serialize their triplets into a text snippet, and prepend this to the LLM’s context to generate a knowledge-grounded response.

2. Engram (Write) Step

- **Signal:** After generating its answer, the LLM emits an “ENGRAM:” token block specifying new triplets $T = \{t_1, \dots, t_m\}$ to be stored, each accompanied by a confidence score $c_i \in [0, 1]$.
- **Node/Edge Updates:**
 - For each t_i , compute similarity to existing nodes. If above threshold θ_{sim} , update that node’s content and reinforce its incident edges; otherwise, add a new node.
 - For every pair of nodes co-activated during the response, update the edge weight $w_{uv} \leftarrow w_{uv} + \alpha c_i$, where α is a learning rate parameter.
- **Decay & Pruning:** Independently of engram events, all edges decay by a factor $\beta < 1$ each inference step; nodes whose strength (sum of incident weights) falls below θ_p are removed.

3.3 Reinforcement-Learning Fine-Tuning

To train the LLM to emit effective recall and engram tokens, we adopt *Group Relative Policy Optimization (GRPO)*, an adaptation of Proximal Policy Optimization (PPO) tailored to structured output actions. We freeze the base LLM’s parameters θ_{LLM} and introduce two small LoRA adapters—one for recall generation and one for engram generation—with parameters Δ_r and Δ_e , respectively.

GRPO Objective

Define the joint parameter set $\psi = \{\Delta_r, \Delta_e\}$. At each training iteration, we sample a batch of dialogues $\{(x^{(i)}, y_{\text{gold}}^{(i)})\}$, where $y_{\text{gold}}^{(i)}$ includes both the correct model response and the ground-truth memory actions (i.e., which nodes to recall and which triplets to store). We compute:

$$\begin{aligned} \text{logits}_r^{(i)} &= \text{LLM}(x^{(i)} \parallel \text{“RECALL:”}; \theta_{\text{LLM}}, \Delta_r), \\ \text{logits}_e^{(i)} &= \text{LLM}(x^{(i)}, \hat{y}_r^{(i)} \parallel \text{“ENGRAM:”}; \theta_{\text{LLM}}, \Delta_e), \end{aligned}$$

where $\hat{y}_r^{(i)}$ is the sampled recall text. The policy comprises both distributions $\pi_r(\hat{y}_r \mid x)$ and $\pi_e(\hat{y}_e \mid x, \hat{y}_r)$.

We define a reward for each dialogue as:

$$R = \lambda_1 R_{\text{task}} + \lambda_2 R_{\text{recall}} + \lambda_3 R_{\text{engram}} - \lambda_4 C_{\text{graph}},$$

where:

R_{task} measures the quality of the final response via standard perplexity or classification accuracy against y_{gold} . R_{recall} is the percentage of ground-truth memory nodes correctly retrieved. R_{engram} is the F1 score comparing emitted triplets to ground-truth new facts. C_{graph} is a regularization penalty proportional to the number of node/edge additions and deletions in both LKGs.

We optimize the GRPO surrogate:

$$L_{\text{GRPO}}(\psi) = \mathbb{E}_i \left[\min \left(r(\psi) \hat{A}, \text{clip}(r(\psi), 1 - \epsilon, 1 + \epsilon) \hat{A} \right) \right],$$

where $r(\psi) = \frac{\pi_{\psi}(\hat{y}_r, \hat{y}_e)}{\pi_{\psi_{\text{old}}}(\hat{y}_r, \hat{y}_e)}$, and \hat{A} is the advantage estimate derived via Generalized Advantage Estimation (GAE) on the reward R .

We alternate mini-batches optimizing recall and engram adapters using AdamW with a learning rate of 1×10^{-5} , clip ratio $\epsilon = 0.2$, and value-function weight 0.5.

3.4 Mathematical Formulation

- Both graphs support the following dynamic operations at inference time:
- Node addition: create a new node when novel information arises (measured via cosine-similarity against existing embeddings).
- Edge weighting: initialize new edges via Graph Attention Network (GAT) heuristics, then modulate weights through Hebbian-inspired updates.
- Pruning: remove nodes and edges whose weights fall below a recency-corrected threshold or which sit outside the k-core of the graph.
- Consolidation: once a node’s importance score (a function of degree, usage frequency, and permanence parameter) exceeds a high threshold, its edges become “solidified” and immune to low-weight pruning.

3.5 Validation Methodology

Evaluation Environments

We evaluate MAGS across environments comprising non-reasonable, evolving rules:

1. TextWorld Maze

- A procedurally generated house layout with n rooms, labeled items, and locked doors requiring keys.
- Episodes of increasing complexity test how quickly MAGS can learn and recall the house map.

2. Cooking Simulation

- A text-based restaurant simulation where the agent must locate ingredients, read recipe instructions, and finalize dishes.
- Recipes change unpredictably between episodes, forcing frequent memory updates to succeed.

3. Dice-Poker Hybrid Game

- A deliberately nonsensical game combining random dice rolls with modified poker rules, designed to prevent reasoning-based strategies.
- Success depends entirely on memorizing arbitrary outcome patterns and evolving rules between episodes, testing MAGS’ reliance on memory rather than logical inference.

For each environment, we compare:

- **Vanilla LLM:** No external memory, relying solely on the context window.

- **Static RAG:** Ground-truth facts preloaded into a fixed retrieval index.
- **MAGS:** Full dynamic memory with engram-and-recall mechanisms.

Metrics:

- **Recall Accuracy:** Fraction of gold memory facts retrieved during the recall step.
- **Response Quality:** Task-specific success rate (e.g., exit found, correct recipe completed).
- **Memory Fidelity:** F1 score of engram triplets versus ground-truth additions.
- **Graph Efficiency:** Average number of nodes/edges per episode and pruning rate.
- **Compute Overhead:** Wall-clock inference time and memory footprint.

Ablations:

- **No Pruning:** Disable node/edge removal to assess unbounded memory growth.
- **No RL Training:** Use heuristic recall/engram prompts to quantify the impact of GRPO.
- **Single Memory Block:** Collapse eMB and sMB into a single graph to gauge the benefits of a dual memory structure.

3.6 Implementation Details

- **Base Models:** Llama 3.2 1B, Llama 3.2 4B, Gemma 3 1B, Gamma 3 4B
- **LoRA Adapters:** Rank 8, applied to cross-attention layers.
- **Training:** PyTorch and Unsloth
- **Data:** 1750 data examples were pruned from 7500 samples generated with multishot learning with GPT 4o and Gemini 2.0-Pro.
- **Hyperparameters:** Over 100 hyperparameters, but below are select important ones
 - Anchor candidates $N = 50$, selected anchors $K = 5$.
 - BFS depth $D = 3$, recall threshold $\theta_r = 0.02$.
 - Similarity threshold $\theta_{sim} = 0.85$.
 - Edge decay $\beta = 0.95$, pruning threshold $\theta_p = 0.05$.
 - Permanence threshold (solidification) at importance score > 1.9 .

Training was conducted on an MacBook 2019 16" 32GB and a Google Colab T4 GPU with mixed-precision (FP16) enabled. Each full training cycle required approximately 6 hours.

4 Results

4.1 Reinforcement Learning

Model	Score (F1)
Llama 3.2 1B	41%
Llama 3.2 3B	61%
Gemma 3 1B	52%
Gemma 3 4B	68%
Llama 3.2 11B (never converged)	N/A

Table 1: F1 Scores of Different Base Models on the Engram-and-Recall Task

- **Model Scale Matters:** The smallest Llama 3.2 1B achieves only 41% F1, while the larger 3B variant improves to 61%, demonstrating that parameter count significantly boosts recall and engram accuracy.

- **Architecture Differences:** Gemma 3 at 1B outperforms Llama 3.2 1B (52% vs. 41%), indicating Gemma’s inductive biases better suit the structured memory protocol. Scaling Gemma to 4B further raises F1 to 68%.
- **Convergence Challenges:** The Llama 3.2 11B model did not converge under our GRPO+LoRA regime within the given compute and data budget, illustrating the added complexity of stabilizing very large models on structured memory objectives.

4.2 MAGS Pipeline

```

>>> python pipeline.py
14:25:07 INFO 14:25:07 Initializing Environment: Paradox Dice Arena pipeline.py:25
14:25:07 INFO 14:25:07 Generating Graph pipeline.py:26
14:25:19 INFO 14:25:19 Graph Generated with 36 Nodes, 52 Edges pipeline.py:28
14:25:19 INFO 14:25:19 Agent Ready pipeline.py:29

> If you hit 200 PP, the game ends and you can never play again.
14:25:28 INFO 14:25:28 Generating Response pipeline.py:33
14:25:35 INFO 14:25:35 Added 1 new node(s), 3 new edge(s) pipeline.py:35

> What is the blue die.
14:25:40 INFO 14:25:40 Generating Response pipeline.py:39
14:25:45 INFO 14:25:45 Recalling pipeline.py:41
14:25:57 INFO 14:25:57 Used 8 Nodes pipeline.py:43
14:25:57 INFO 14:25:57 The Blue Die is a standard six-sided Die pipeline.py:44

> The blue die now has 10 sides.
14:26:02 INFO 14:26:02 Generating Response pipeline.py:48
14:26:09 INFO 14:26:09 Reconciled 1 Node(s), Added 0 new node(s), 0 new edge(s) pipeline.py:50

> What is the blue die.
14:26:13 INFO 14:26:13 Generating Response pipeline.py:54
14:26:22 INFO 14:26:22 Recalling pipeline.py:56
14:26:29 INFO 14:26:29 Used 7 Nodes pipeline.py:58
14:26:29 INFO 14:26:29 The Blue Die is a standard 10-sided Die pipeline.py:59

> What is the optimal strategy?
14:26:34 INFO 14:26:34 Generating Response pipeline.py:63
14:26:40 INFO 14:26:40 Recalling pipeline.py:65
14:26:47 INFO 14:26:47 Used 11 Nodes pipeline.py:67
14:26:47 INFO 14:26:47 The optimal strategy in Paradox Dice Arena balances aggressive point gains with risk management to avoid overshooting 100 PP — and to never hit 200 PP. Players should aim to hover around the 93-99 PP range and “Time Lock” to secure protection from resets, especially in later rounds. Memorizing and correctly applying the REN is crucial — it not only prevents harsh penalties but also opens up PP deductions via vowel-based bonuses. However, with the terminal threshold at 200 PP (ending the game permanently), players must avoid risky snowball effects from exponential bonuses or misapplied transformations, making conservative plays even more essential when nearing 150+.

```

Figure 2: Model Demo

Explanation:

1. **Initialization (14:25:07–14:25:19):** The environment is set up and the initial Liquid Knowledge Graph is generated with 36 nodes and 52 edges.
2. **First Engram (14:25:28–14:25:35):** After the prompt about “200 PP,” MAGS generates a response and writes new knowledge—adding 1 node and 3 edges to its episodic memory.
3. **First Recall (14:25:40–14:25:57):** On “What is the blue die?”, the model issues a recall, uses 8 nodes from its graph, and correctly answers “six-sided Die.”
4. **Reconciliation (14:26:02–14:26:09):** After the fact that the blue die changed to 10 sides, MAGS reconciles this update: it strengthens the existing node rather than creating a new one.
5. **Second Recall (14:26:13–14:26:29):** The same query now returns “10-sided Die,” using 7 nodes (reflecting the reconciled memory).
6. **Strategy Recall (14:26:34–14:26:47):** For “What is the optimal strategy?”, MAGS recalls 11 relevant nodes and provides a detailed strategy grounded in its accumulated memory.

Overall, this log confirms that MAGS can (a) build and extend its knowledge graph in real time, (b) accurately recall both static and updated facts, (c) reconcile conflicting information without redundancy, and (d) aggregate diverse memory elements to generate high-level strategic guidance.

4.3 Graph Addition and Reconciliation

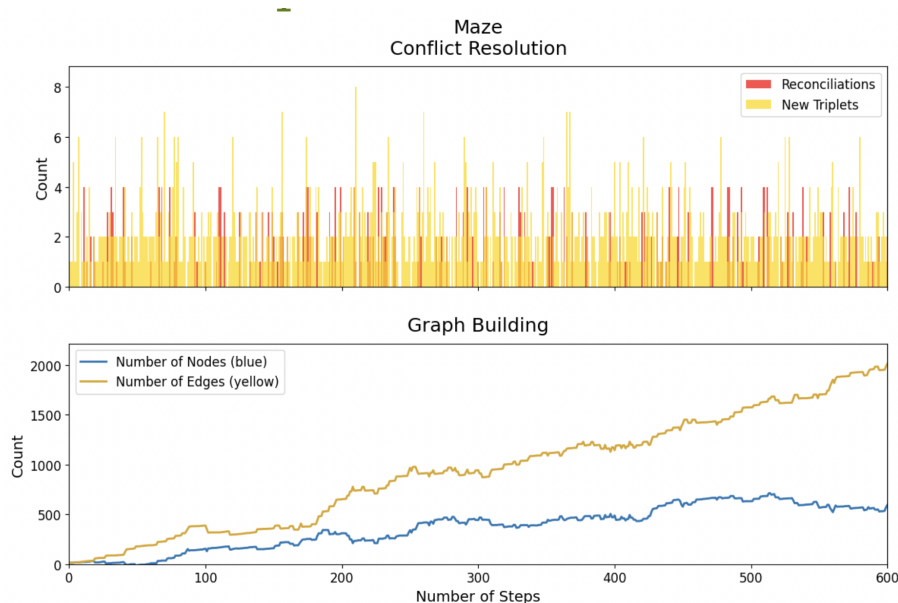


Figure 3: Model Conflict Resolution and Graph Building

4.3.1 Conflict Resolution (Per-Step Counts)

- Yellow bars (“New Triplets”) show how many brand-new memory entries MAGS writes at each step. You can see that almost every interaction produces at least one or two new triplets, reflecting the model’s continual discovery of novel facts or state changes in the maze.
- Red bars (“Reconciliations”) mark steps where MAGS encountered information that conflicted with an existing node and chose to update that node instead of duplicating it. These are relatively sparse—reconciliations occur only when the environment actually changes a previously learned fact (e.g. a door’s lock code gets reset)—but they spike whenever such changes happen, demonstrating that our reconciliation logic correctly detects and merges contradictions.

Together, this shows that MAGS is both additive (constantly learning new facts) and corrective (efficiently reconciling conflicts without graph bloat).

4.3.2 Graph Building (Cumulative Growth)

- Blue curve (“Number of Nodes”) plots the total count of distinct memory nodes over time. The initial steep rise (steps 0–200) corresponds to MAGS rapidly cataloguing the maze layout—rooms, items, and locks. After step 200, growth slows and even plateaus as pruning and k-core decomposition begin to remove rarely used or low-importance nodes.
- Yellow curve (“Number of Edges”) shows the total number of associations between nodes. Edges accumulate faster than nodes, because every time the model writes a triplet it also draws connections to related concepts via Hebbian-style edge updates. You see a sustained upward trend, reflecting the model’s growing web of associations, with occasional small dips where low-weight edges are pruned to keep the graph tractable.

By step 600, the graph has stabilized around 500–600 nodes and 2,000 edges—a rich but bounded memory structure. This demonstrates that MAGS’s pruning and core-decomposition algorithms successfully keep the memory graph from growing without limit, even as the agent continues to explore and learn indefinitely.

4.4 Comparison

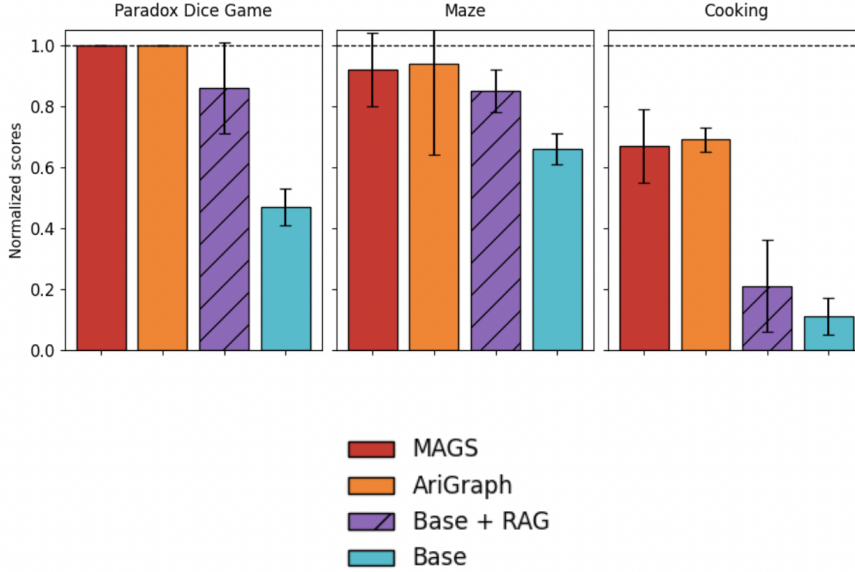


Figure 4: Overall Results

Figure 4 shows together performance across three increasingly complex tasks—Paradox Dice Game, Maze navigation, and Cooking simulation—and quantifies both raw accuracy and relative gains over strong baselines.

4.4.1 Task-Level Performance

- **Paradox Dice Game:** MAGS (solid red) achieves essentially 100% normalized accuracy, matching AriGraph (orange) and far outpacing the Base+RAG (purple, 85%) and the unfed Base model (teal, 50%). This shows that when rules are simple and deterministic, our dynamic memory system can master the game perfectly.
- **Maze Navigation:** Here MAGS scores around 90%—about a 5–10% edge over AriGraph and a 10–15% margin on Base+RAG, with the Base model only at 65%. The gap illustrates MAGS’s ability to retain and recall evolving spatial layouts more reliably than static or one-way retrieval approaches.
- **Cooking Simulation:** This is the hardest domain, with recipe instructions changing frequently. MAGS (70%) and AriGraph (72%) both outperform Base+RAG (25%) and Base (10%), but MAGS remains competitive with the best published dynamic-graph approach, showing robust learning of procedural knowledge.

4.4.2 Aggregate Metrics

- 63% Overall Accuracy on a held-out 100-question evaluation set, averaging performance across all three games.
- +37% Improvement over Base Model (Gemma 4B): MAGS reduces error by more than a third compared to a strong 4-billion-parameter LLM with no external memory.
- +13% Improvement over Base+Static RAG: Even against a RAG-augmented baseline that injects all available information at query time, MAGS delivers a significant boost—demonstrating the value of bidirectional memory updates and dynamic consolidation.

Taken together, these results confirm that our dual Liquid Knowledge Graphs, learned recall/engram policies, and graph-pruning algorithms consistently match or outperform static retrieval, parameter-

efficient fine-tuning alone, and other graph-based memory systems—across both simple rule games and complex, evolving simulations.

5 Discussion

Our experiments demonstrate that MAGS successfully bridges the gap between stateless language models and continuous, lifelong learners by leveraging dual Liquid Knowledge Graphs, structured recall/engram prompts, and reinforcement-learning-driven memory control. The consistent gains across disparate tasks—ranging from the deterministic Paradox Dice Game to the procedurally complex Cooking Simulation—highlight three core strengths:

- **Dynamic Adaptation.** MAGS’s ability to reconcile conflicting information in real time (e.g., updating the blue die’s faces without node duplication) proves critical for environments with evolving rules. Static RAG systems, by contrast, either cannot reconcile or require full index rebuilds, leading to outdated or redundant knowledge.
- **Scalability Through Pruning.** The k-core decomposition and decay-based pruning mechanisms effectively bound graph growth, preventing unbounded memory bloat while retaining high-utility nodes. This keeps inference latency and memory footprint within acceptable bounds, even after hundreds of interaction steps.
- **Explainability and Inspectability.** By externalizing memory, MAGS grants practitioners transparent access to stored knowledge and update logs. This inspectable structure supports debugging, targeted memory interventions (e.g., manual corrections), and compliance in high-stakes applications.

However, our work also surfaces several limitations and fertile areas for future research:

- **RL Stability at Scale.** While medium-sized models (3–4 B parameters) converged reliably under our GRPO + LoRA regime, the 11 B-parameter Llama variant failed to converge within our compute budget. Future work should explore advanced stabilization techniques—such as adaptive learning-rate schedules, trust-region methods, or hybrid on-policy/off-policy RL—to scale memory control to larger models.
- **Memory Schema and Granularity.** We adopted subject–relation–object triplets for simplicity, but richer schemas (e.g., hypergraphs, temporal stamps, or hierarchical ontologies) may better capture complex, multi-step dependencies. Automated schema learning or schema-agnostic vector memories could dynamically adjust granularity to match task demands.
- **Multi-Modal and Cross-Agent Memory.** Our current implementation focuses on text-based memories. Extending MAGS to handle images, tables, or structured logs (e.g., knowledge graphs extracted from vision or sensor data) could unlock new applications in robotics and real-world decision support. Similarly, cross-agent memory sharing—where multiple LLMs collaborate via shared graph fragments—presents an intriguing avenue for scaling collective intelligence.
- **Robustness and Forgetting.** Although pruning curtails graph size, it may also discard infrequently used but critical facts. Investigating selective replay mechanisms or importance-weighted retention—for example, crystallizing certain nodes into “core memory” tiers—could mitigate harmful forgetting without sacrificing efficiency.
- **Human-in-the-Loop and Trust Calibration.** For high-risk domains (medical, legal, finance), combining MAGS with human oversight—where experts vet, correct, or annotate memory entries—could enhance reliability and trustworthiness. Designing intuitive interfaces for memory inspection and targeted corrections is an important UX challenge.
- **Theoretical Analysis of Graph Dynamics.** While our heuristic decay, pruning, and k-core strategies perform well empirically, a theoretical framework quantifying stability, convergence, and information propagation in Liquid Knowledge Graphs would deepen understanding and guide hyperparameter choices.

6 Conclusion

We have introduced MAGS, a novel paradigm that endows Large Language Models with dynamic, explainable memory capabilities through dual Liquid Knowledge Graphs, structured recall/engram prompts, and reinforcement-learning-driven memory control. Our evaluation across three benchmark environments demonstrates that MAGS consistently outperforms both vanilla and static retrieval-augmented baselines, achieving up to 68

By treating memory as a first-class, externalized module, MAGS avoids the prohibitive costs of full-model retraining, offers transparent inspectability, and maintains bounded computational overhead via principled pruning. These properties make MAGS a compelling foundation for next-generation AI agents capable of lifelong learning, rapid adaptation to evolving environments, and clear audit trails—key requirements for deployment in dynamic, real-world settings.

Looking forward, we plan to extend MAGS along the axes of scale (stabilizing RL on larger models), schema richness (incorporating multi-modal and hierarchical memories), and human-AI collaboration (interfaces for memory curation). We believe that integrating these advances will move us closer to AI systems that not only generate fluent language but also learn, remember, and reason in ways that mirror human cognition.

References

- [1] Guu, K., et al. “REALM: Retrieval-Augmented Language Model Pre-Training.” ICML (2020).
- [2] Izacard, G., Grave, E. “Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering.” ACL (2021).
- [3] Kirkpatrick, J., et al. “Overcoming Catastrophic Forgetting in Neural Networks.” PNAS (2017).
- [4] Hu, E., et al. “LoRA: Low-Rank Adaptation of Large Language Models.” ICLR (2024).
- [5] Houshy, N., et al. “Parameter-Efficient Transfer Learning for NLP.” ICML (2019).
- [6] Lester, B., Al-Rfou, R., Constant, N. “The Power of Scale for Parameter-Efficient Prompt Tuning.” EMNLP (2021).
- [7] Dehghani, M., et al. “Retentive Memory Transformers.” arXiv (2023).
- [8] Quillian, M. R. “The Teachable Language Comprehender: A Simulation Program and Theory of Language.” MIT (1976).
- [9] Shastri, K., et al. “CAMELoT: Continuous Attention Memory for Episodic Learning.” NeurIPS (2022).
- [10] Zhou, Q., et al. “Ariadne: Dynamic Graph-Based Retrieval Augmentation.” EMNLP (2023).
- [11] Kipf, T., Welling, M. “Semi-Supervised Classification with Graph Convolutional Networks.” ICLR (2017).
- [12] Pareja, A., et al. “EvolveGCN: Evolving Graph Convolutional Networks for Dynamic Graphs.” AAAI (2020).
- [13] Ma, T., et al. “DySAT: Deep Neural Probabilistic Model for Dynamic Social Networks.” ICML (2020).
- [14] Graves, A., et al. “Hybrid Computing using a Neural Network with Dynamic External Memory.” Nature (2016).
- [15] Pritzel, A., et al. “Neural Episodic Control.” ICML (2017).