

Beyond Markovian Forgetfulness: Episodic Memory for Reasoning-Intensive Retrieval

Dohyeon Lee^{1,2}, Yeonseok Jeong³, Seung-won Hwang^{2,3*}

Korea Advanced Institute of Science and Technology (KAIST)¹

Computer Science and Engineering, Seoul National University²,

Interdisciplinary Program in Artificial Intelligence, Seoul National University³

waylight3@kaist.ac.kr, jys3136@snu.ac.kr, seungwonh@snu.ac.kr

Abstract

Reasoning-intensive information retrieval uses large language models to solve complex queries via multi-step reasoning. However, existing methods have critical limitations. Chain-of-Thought (CoT) approaches suffer from inefficiency, while state-based methods, despite better token efficiency, often fall into reasoning cycles that trap the query refinement process. To address these issues, we propose Episodic Memory for Retrieval (EMR), which enhances the state-based framework with an episodic memory. This module stores the full history of prior states for a query, allowing the model to avoid repetition of such cycles. Experiments on the BRIGHT benchmark show that EMR consistently outperforms both CoT and state-based baselines. Moreover, it is highly token-efficient, reducing token usage by 72% on average. Our results show that episodic memory is an effective and token-efficient mechanism for reasoning-intensive retrieval. The gains also generalize across different base models and stay efficient in terms of end-to-end latency. The code is available in <https://github.com/ldilab/EMR>.

1 Introduction

Traditional Information Retrieval (IR) systems often struggle to answer complex queries that require multi-step reasoning (Yang et al., 2018; Feldman and El-Yaniv, 2019). Reasoning-intensive IR (Xiao et al., 2024; Das et al., 2025) leverages large language models (LLMs) to enhance retrieval performance through techniques such as iterative query refinement and document reranking.

Early approaches to reasoning-intensive IR, such as Rank1 (Weller et al., 2025) and Rank-R1 (Zhuang et al., 2025) leveraged the Chain-of-Thought (CoT) capabilities of LLMs. While effective in producing detailed reasoning traces, these

* Corresponding Authors

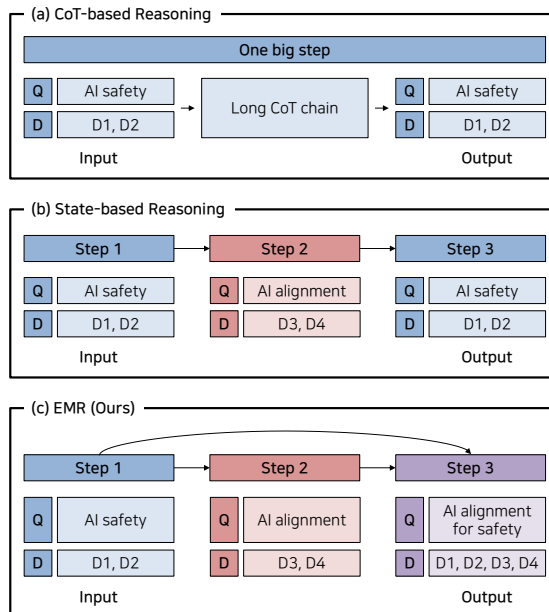


Figure 1: An illustration of different reasoning approaches for query refinement. (a) **CoT-based Reasoning** generates a single and long reasoning chain, which can be inefficient and may not result in a meaningful query update. (b) **State-based Reasoning** refines the query based solely on the immediate previous state, making it vulnerable to reasoning cycles. (c) **EMR** enhances state-based reasoning with an **episodic memory** that retains the history of all past states. This prevents cycles and enables progressive query refinement.

approaches sacrifice token efficiency. As shown in Figure 1(a), the model generates a long and redundant reasoning chain within a single step. This process can lead to a *reasoning cycle* (e.g., AI safety → AI safety) that consumes substantial computation without yielding progress.

To improve efficiency, state-based methods such as SMR (Lee et al., 2025) decompose reasoning into iterative transitions. By design, each step is prompted to produce a new state, mitigating immediate self-loops and reducing token usage. However, this framework imposes a Markovian assump-

tion, conditioning each new state only on the immediate predecessor. As a result, multi-step cycles can still emerge (e.g., AI safety \rightarrow AI alignment \rightarrow AI safety), as shown in Figure 1(b). In our experiments, we observed that such loops occur in up to 65% of queries, depending on the dataset.

We address these limitations with **Episodic Memory for Retrieval (EMR)**, which augments the state-based paradigm with an explicit memory module to eliminate reasoning cycles. Inspired by *episodic memory* in agentic approaches (DeChant, 2025; Nuxoll and Laird, 2012), which refers to the ability to recall past experiences, our method records the full history of states generated for a query. The model is then prompted to condition on this trajectory when producing the next state. This trajectory-aware mechanism prevents the revisiting of prior states and encourages the generation of progressively refined queries.

Figure 1(c) shows how a query progresses through refinement: For instance, by referencing its memory of having already generated AI safety and AI alignment, the model is guided to produce a more advanced query like AI alignment for safety, rather than repeating previous steps. This episodic memory is implemented as a text-based log within the model’s system prompt, ensuring constant accessibility during inference. Furthermore, through prompt compression, EMR achieves superior token efficiency compared to the state-based method, despite the additional tokens to store past states.

To validate EMR, we conducted extensive experiments on BRIGHT, a widely-used benchmark for reasoning-intensive IR. Across the evaluated settings, EMR consistently outperforms both CoT-based and state-based baselines. Furthermore, we show that our approach is token-efficient, reducing token consumption by 72% on average compared to the state-based approach. This result suggests that the gains in token efficiency from episodic memory are larger than the cost introduced by maintaining it.

2 Related Work

2.1 CoT-based Reasoning for IR

The advent of LLMs has enabled new approaches for addressing complex information retrieval tasks that demand multi-step reasoning. Early approaches in this domain, such as Rank1 (Weller et al., 2025) and Rank-R1 (Zhuang et al., 2025),

leveraged CoT prompting to deconstruct intricate information needs.

However, the lengthy reasoning chains they generated are often redundant, which could misguide the query refinement process. Subsequent CoT compression approaches, like O1-Pruner (Luo et al., 2025), aimed to mitigate this by shortening the CoT chain. Despite these efforts, the fundamental limitation of the CoT paradigm is its single-pass generation process. This approach is vulnerable to reasoning cycles, motivating a shift toward more granular, state-based methods.

2.2 State-based Reasoning for IR

To address the inefficiencies of CoT, methods like SMR (Lee et al., 2025) reformulated the task as a sequence of state transitions. This approach significantly improved token efficiency by breaking down the reasoning process into discrete steps.

However, its reliance on a Markovian assumption makes it susceptible to multi-step reasoning cycles. This motivates a mechanism that can retain and leverage the history of past states to ensure consistent forward progress.

2.3 Memory in Agentic Approaches

The need for a history-aware mechanism has been explored in the field of agentic AI. In this domain, episodic memory (DeChant, 2025; Nuxoll and Laird, 2012) refers to an agent’s ability to record and recall a sequence of past experiences, allowing it to engage in more complex planning.

Inspired by this approach, our primary contribution is the integration of this episodic memory concept into the state-based IR framework. Unlike previous methods conditioned only on the prior state, our approach conditions the policy on the entire history of previously generated queries. This allows our model to systematically prevent reasoning cycles while preserving the high token efficiency of the state-based paradigm, allowing more consistent forward progress during query refinement.

3 Method

Our approach enhances the state-based reasoning paradigm as a sequential state-transition problem, defining its core components: states, actions, and the policy model (§3.1). Building upon this foundation, we then introduce our core contribution, EMR, and describe how its episodic memory is constructed and applied to guide the reasoning process (§3.2).

3.1 Problem Setup: State-based Reasoning

Following the recent work (Lee et al., 2025), we formulate the reasoning-intensive IR task as a sequential decision-making process. This process is characterized by three fundamental components: a state representing the current context, an action that modifies it, and a policy that selects the best action. This formulation can be viewed as a Markov Decision Process (MDP), where the policy guides the transitions between states to progressively refine the initial state.

State. We define the state s_t at step t as a tuple consisting of the current query and a set of retrieved documents. This represents a snapshot of the reasoning process at a given moment. Formally, a state is defined as:

$$s_t = (q_t, D_t) \quad (1)$$

where $q_t \in \mathcal{Q}$ is the query at step t , with \mathcal{Q} being the list of all possible queries. $D_t \subset \mathcal{D}$ is the list of documents retrieved using the query q_t , where \mathcal{D} represents the entire document corpus. The goal is to transition from an initial state s_0 to a final state s_n where D_n is optimized for the user’s information need.

Action. An action a_t represents a strategic operation chosen by the policy to advance the search process. Our action space \mathcal{A} consists of three distinct operations:

$$\mathcal{A} = \{ \text{REFINE}, \text{RERANK}, \text{STOP} \} \quad (2)$$

The REFINE action is deployed to broaden the search space. It uses the LLM to generate a new query q_{t+1} , retrieves a new set of documents, and expands the current document list D_t by appending these new results. Conversely, the RERANK action is used to exploit the currently held information. It leverages the LLM to perform a listwise reranking of documents in D_t , prioritizing the most relevant items. For computational feasibility, this list is then truncated to the top- k entries, creating a more focused document set D_{t+1} while the query q_t remains stable. Finally, the process terminates with the STOP action when the policy determines that the gathered documents in D_t are sufficient. This terminates the loop and returns D_t as the output.

Policy. The Policy π serves as the decision-making engine of our reasoning agent, strategically guiding the search process toward a res-

olution. Embodied by an LLM, the policy assesses the current state $s_t = (q_t, D_t)$ at each step t to determine the most promising action $a_t \in \{ \text{REFINE}, \text{RERANK}, \text{STOP} \}$ to execute. For efficiency, the LLM generates both the chosen action and its direct result within a single forward pass. For instance, if the policy chooses REFINE, it outputs the action itself along with the newly refined query. This one-call-per-step architecture minimizes latency and computational cost. To ensure the reasoning process remains bounded, the policy is also constrained by a predefined maximum number of steps, T_{\max} . If the current step count reaches this threshold, the policy issues a STOP action, preventing overly long or costly reasoning chains.

Reasoning Cycle. A significant drawback of the state-based formulation is its vulnerability to reasoning cycles, where the agent becomes trapped in a loop by revisiting previously explored states. This issue is a direct consequence of the policy’s Markovian nature: its decisions are conditioned solely on the current state s_t , with no memory of the path taken to reach it. Formally, the reasoning cycle is defined as follows:

$$q_t = q_{t'} \quad (t' < t) \quad (3)$$

where equality denotes exact lexical identity. In practice, we flag a cycle when the newly generated query string exactly matches any previously generated query in memory. Such loops prevent meaningful progress, leading to redundant computation and suboptimal retrieval performance. This limitation highlights the need for a mechanism that can break the Markovian chain by incorporating a memory of its history into the agent.

3.2 Proposed: EMR

This section discusses how EMR introduces an episodic memory, by first describing the conceptual role of episodic memory in state-based reasoning (§3.2.1) and then detailing its implementation for reading and writing states in textual form (§3.2.2). Finally, we introduce the memory compression mechanism designed to maintain high token efficiency, even with the overhead of storing memory states (§3.2.3).

3.2.1 Role of Episodic Memory

The role of the episodic memory in EMR is to break the Markovian dependency inherent in the

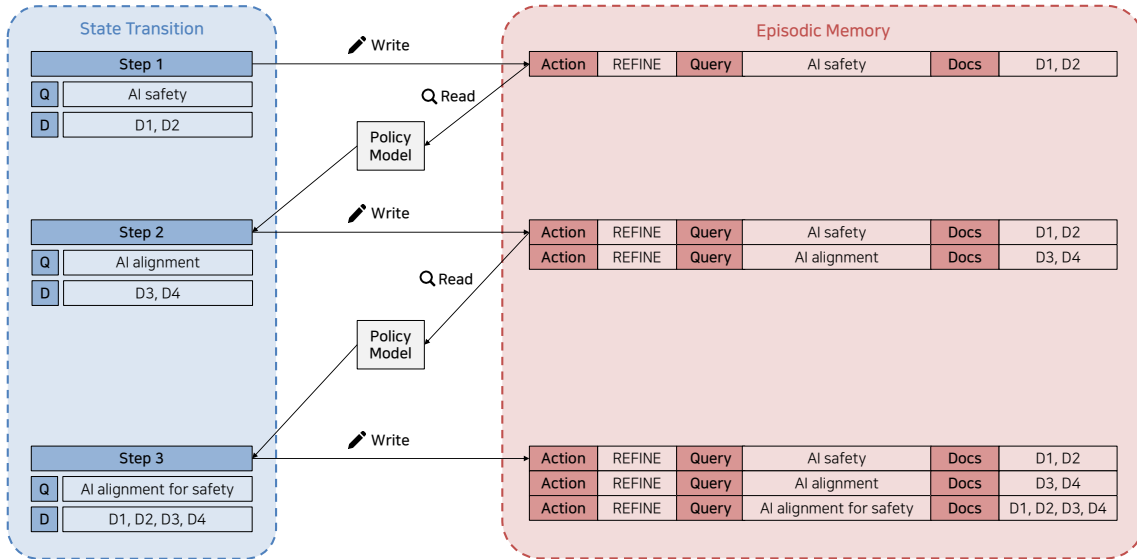


Figure 2: An illustration of the EMR architecture and the role of episodic memory. The figure shows the interplay between the iterative **State Transition** process (left) and the **Episodic Memory** (right). At each step, the policy model generates the next state based on the memory’s history, and then updates the memory with that newly generated state.

standard state-based approach. As defined in §3.1, the vanilla policy π makes decisions based solely on the current state s_t , rendering it blind to the path taken to arrive there.

To address this, we introduce an episodic memory, denoted as M_t , which is the complete history of all states visited up to step t . Formally, the memory is a sequence of past states:

$$M_t = (s_0, s_1, \dots, s_{t-1}) \quad (4)$$

where $s_i = (q_i, D_i)$ is the state at step i .

By incorporating this memory, we redefine the policy function. The original policy $\pi(s_t)$ is transformed into a memory-augmented policy, π' , which conditions its decisions on both the current state and the entire history stored in the memory. The action selection process is thus updated as follows:

$$a_t = \pi'(s_t, M_t) \quad (5)$$

This seemingly simple modification has two profound implications for the reasoning process, as illustrated in Figure 2.

First, the memory-augmented policy can explicitly prevent cycles. By having access to the set of all past queries $\{q_0, q_1, \dots, q_{t-1}\}$ within M_t , the policy π' can be instructed to avoid generating any query that has been previously explored. This directly resolves the cyclic behavior defined in Equation 3.

Second, beyond cycle prevention, the memory provides a rich historical context to generate more targeted queries, rather than just reacting to the immediate present. The following sections will detail the practical implementation of how this conceptual memory is read and written as a textual log within the LLM’s prompt.

3.2.2 Memory Read and Write Operations

To utilize the conceptual memory M_t described in the previous section, we designed a token-efficient textual format that can be seamlessly integrated into the policy model’s prompt. This process involves two key operations: writing the historical states into a textual log and reading this log as context for the next decision. The entire memory structure is prepended to the main task prompt at each step t , ensuring the model has constant access to its history.

The memory is formatted as plain text and organized into two distinct sections as shown in Table 1: a history of actions and a repository of document contents.

The History of Recent Actions section serves as a log of the reasoning trajectory. A crucial design choice for token efficiency is that each entry stores only the identifiers (IDs) of the documents retrieved at that step, rather than their full contents. This significantly reduces the token count

```

## History of Recent Actions
[1] Action: {a_1} Query: {q_1} Ranks: {IDs of D_1}
...
[t-1] Action: {a_{t-1}} Query: {q_{t-1}} Ranks: {IDs
of D_{t-1}}

## Memory of Documents
[ID of d_1] {d_1}
[ID of d_2] {d_2}
...
[ID of d_n] {d_n}

```

Table 1: Text format of the history of previous states in the system prompt of EMR.

that grows with each step, as document contents can be hundreds of tokens long while their IDs are just single tokens or short strings.

The Memory of Documents section acts as a deduplicated repository. It maps document IDs to their full textual content. This structure decouples the log from the voluminous document content, preventing redundant storage of the same document if it is retrieved multiple times across different steps.

We deliberately chose a simple plain text format over structured formats like JSON or XML. Our preliminary experiments indicated that while structured formats offered no significant performance improvement, they consistently incurred a higher token overhead due to syntactic characters (e.g., brackets, quotes, and tags). The plain text approach provides the best balance of performance and token economy for this task.

Read Operation. The read operation is performed implicitly at the beginning of each step t . The entire two-part memory block, representing M_t , is formatted as a single string and placed within the LLM’s prompt, typically following the system message. This provides the full historical context required for the memory-augmented policy $\pi'(s_t, M_t)$ to make its next decision.

Write Operation. The write operation occurs after the policy executes an action a_t and a new state $s_{t+1} = (q_{t+1}, D_{t+1})$ is generated. The process involves two updates:

- A new formatted string is appended to the History of Recent Actions.
- Each document in the retrieved set D_t is checked. If a document’s ID is not already present in the Memory of Documents, its ID and full content are appended to that section.

This read-write cycle continues until the STOP action is issued, ensuring that the memory is always a complete and up-to-date record of the entire reasoning process.

3.2.3 Memory Compression

Despite deduplication, the cumulative length of unique document contents still poses a significant challenge to token efficiency. As a result, as the reasoning process explores more documents, the memory’s token footprint can grow excessively, leading to increased computational costs and potential context window limitations. To mitigate this issue, we introduce a memory compression mechanism that reduces each document to only its most query-relevant sentences.

This compression is performed via a three-step extractive summarization process whenever a new set of documents D_t is retrieved by a query q_t .

Sentence Segmentation. All documents in D_t are split into sentences using spaCy¹, creating a large pool of candidate sentences.

Relevance Scoring. Each sentence is scored for relevance to q_t using a pre-trained MiniLM cross-encoder².

This compression only affects the current prompt construction, as the backend preserves the original text of every document. Because no data is permanently lost, the system can re-summarize the same content differently for later requests.

Global Filtering and Composition. Unlike per-document selection, we rank every sentence from the entire pool and keep only the top- k overall. This global approach ensures that only the most relevant sentences enter the Memory of Documents, effectively pruning irrelevant documents. We select top- k filtering specifically to ensure predictable inference costs. While threshold-based methods might offer better adaptivity, they often result in inconsistent prompt sizes that are harder to manage. This selection process reduces token overhead and focuses the model on information essential for reasoning. Furthermore, while stronger encoders could enhance the quality of sentence selection, such improvements are orthogonal to the core mechanism of EMR.

¹<https://spacy.io/>

²<https://huggingface.co/cross-encoder/ms-marco-MiniLM-L6-v2>

4 Results and Analysis

4.1 Experimental Setup

Datasets and Metrics. We evaluate on the BRIGHT benchmark (Su et al., 2024), which contains 12 reasoning-intensive retrieval datasets spanning domains such as mathematics, code, and science. We use nDCG@10 as the primary metric and report additional results on Recall and MAP in the Appendix A.3.

Baselines. We compare EMR with three CoT-based baselines, Rank1 (Weller et al., 2025), RankR1 (Zhuang et al., 2025), and O1-Pruner (Luo et al., 2025), and with the state-based baseline SMR (Lee et al., 2025). We report results with two retrieval backbones, BM25 (Robertson et al., 2009) and ReasonIR (Shao et al., 2025). More details are provided in Appendix A.1.

Implementation Details. Unless otherwise noted, we follow the original implementations and settings of prior methods. We use Qwen2.5-32B and Qwen3-32B as the main backbones. Full implementation details are provided in Appendix A.2.

4.2 Analysis

We structure our analysis around three research questions to validate the effectiveness of EMR.

4.2.1 RQ1. Does Reducing Cycles with Episodic Memory Boost Performance?

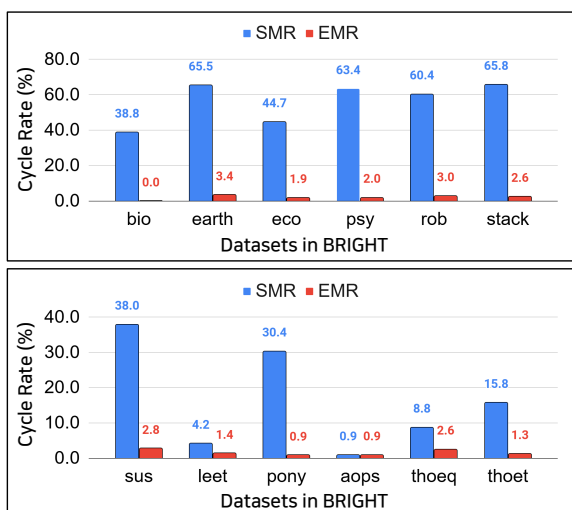


Figure 3: Comparison of the cycle ratio between EMR and SMR on the BRIGHT benchmark. The ratio indicates the percentage of queries that encounter one or more reasoning cycles during the refinement process.

Reducing Reasoning Cycles. To answer RQ1, we first measure the occurrence of reasoning cycles as defined in Eq. 3: any instance where a refined query revisits a previously generated one.

As illustrated in Figure 3, EMR (red) consistently suppresses reasoning cycles compared to SMR (blue) across all BRIGHT datasets. Aggregated over all datasets, the average cycle rate is reduced by 94%, from 38.75% for SMR to just 2.25% for EMR.

A closer analysis reveals that the baseline’s vulnerability to reasoning cycles is heterogeneous, varying with the type of the queries in each domain. On datasets like *Leetcode* and *AoPS*, SMR already exhibits few cycles. This is likely because their queries often contain code snippets or mathematical formulas, whose high lexical specificity makes exact repetition less probable. In contrast, on datasets with natural language queries such as *EarthScience*, *Psychology*, and *StackOverflow*, SMR is highly vulnerable to cyclic behavior, with cycle rates exceeding 60%. Across all datasets, EMR reduces the cycle rates to under 4%. These results show that conditioning on prior states substantially reduces reasoning cycles in state-based retrieval.

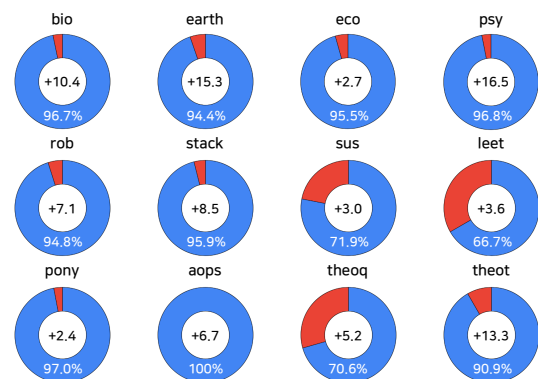


Figure 4: Analysis of queries that encounter cycles with SMR baseline. For each BRIGHT dataset, the donut chart displays the percentage of these queries for which EMR successfully resolves the cycle (blue). The central number indicates the average nDCG@10 point gain for these resolved queries.

Performance Gain in Cycle Queries. To further quantify the impact of cycle mitigation, we conducted a fine-grained analysis on queries where SMR encounters a reasoning cycle. We measure both the rate at which EMR resolves the cycle and the subsequent impact on retrieval performance.

	Bio	Earth	Econ	Psy	Rob	Stack	Sus	Leet	Pony	AoPS	TheoQ	TheoT	Avg
Retriever													
BM25	18.9	27.2	14.9	12.5	13.6	18.4	15.0	24.4	7.9	6.2	10.4	4.9	14.5
CoT-based Reasoning													
BM25 + Rank1	22.1	31.7	14.6	15.7	15.8	17.7	20.3	22.9	9.7	5.9	11.8	9.3	16.5
BM25 + Rank-R1	23.6	33.5	16.8	15.4	18.8	18.4	20.3	24.9	9.1	6.9	12.1	8.7	17.4
BM25 + O1-Pruner	23.6	31.9	17.7	18.0	17.2	20.0	19.8	24.2	8.4	6.6	10.7	7.0	17.1
State-based Reasoning													
BM25 + SMR (Qwen2.5)	28.7	34.9	20.4	20.7	20.9	20.8	19.2	22.1	6.3	6.3	18.0	20.3	19.9
BM25 + SMR (Qwen3)	44.0	42.1	22.9	23.3	18.3	21.1	29.5	14.5	7.8	3.1	6.9	6.4	20.0
EMR: State-based Reasoning with Episodic Memory													
BM25 + EMR (Qwen2.5)	41.2	42.1	24.7	33.1	19.2	23.8	26.1	26.4	8.3	6.8	26.5	29.2	26.4
BM25 + EMR (Qwen3)	53.1	54.6	23.2	37.4	23.4	26.9	31.8	15.8	9.0	7.2	20.9	17.3	26.7

Table 2: Retrieval performance (nDCG@10) on the **BRIGHT** benchmark using **BM25** as the underlying retriever. Best scores per dataset are bolded.

	Bio	Earth	Econ	Psy	Rob	Stack	Sus	Leet	Pony	AoPS	TheoQ	TheoT	Avg
Retriever													
ReasonIR	43.6	42.9	32.7	38.8	20.9	25.8	27.5	31.5	19.6	7.4	33.1	35.7	30.0
CoT-based Reasoning													
ReasonIR + Rank1	49.7	35.8	22.0	37.5	22.5	21.7	35.0	18.8	32.5	10.8	22.9	43.7	29.4
ReasonIR + Rank-R1	50.3	36.1	24.4	38.9	23.7	22.0	34.2	22.6	31.3	10.7	24.5	44.0	30.2
ReasonIR + O1-Pruner	48.6	37.4	24.4	38.8	24.0	21.6	35.3	22.5	31.1	11.2	24.0	44.1	30.3
State-based Reasoning													
ReasonIR + SMR (Qwen2.5)	52.2	51.1	27.5	45.1	22.8	29.4	30.9	27.9	18.1	7.4	36.4	32.6	31.8
ReasonIR + SMR (Qwen3)	53.7	52.3	29.5	46.9	24.5	30.6	32.5	29.0	19.8	9.2	38.0	34.2	33.4
EMR: State-based Reasoning with Episodic Memory													
ReasonIR + EMR (Qwen2.5)	55.9	55.4	33.3	49.4	26.7	33.6	35.1	31.2	22.6	11.1	39.0	35.6	35.7
ReasonIR + EMR (Qwen3)	56.9	56.0	33.5	50.0	28.2	34.1	35.7	32.4	23.1	12.4	40.5	37.2	36.7

Table 3: Retrieval performance (nDCG@10) on **GPT4-Reason queries** of the **BRIGHT** benchmark. Best scores per dataset are bolded.

Figure 4 shows that EMR successfully resolves the cycle in over 90% of these failure cases on average, as indicated by the blue portion of each chart. Moreover, this resolution translates into substantial gain. The number at the center of each chart represents the nDCG@10 gain for these queries, amounting to 8%p increase in average. Notably, the largest performance gains are observed in datasets with natural language queries, such as *EarthScience* (+15.3) and *Psychology* (+16.5).

These results explain the model-specific performance in Table 2. In logical domains like LeetCode and TheoQ, stronger base models often exhibit repetitive refinement behavior. Our analysis reveals that Qwen3 generates 25% more reasoning cycles than Qwen2.5 on LeetCode. This suggests that more capable models are prone to repetitive refinement in specific domains. These findings suggest

that stronger base models can still exhibit repetitive refinement behavior, making explicit trajectory control useful even as model quality improves.

4.2.2 RQ2. Is EMR More Effective and Efficient?

To answer RQ2, we evaluate EMR along two axes: retrieval effectiveness and token efficiency.

Retrieval Effectiveness. Table 2 shows that EMR achieves the best average nDCG@10 under BM25. On average, EMR achieves a 6.7%p gain in nDCG@10 compared to SMR. See Appendix A.3 for a comprehensive evaluation across other metrics such as Recall and MAP.

This performance gain is not just the result of better initial query refinement of LLMs. To verify this, we conducted an experiment using strong initial queries refined by GPT-4, provided by the

BRIGHT benchmark. As detailed in Table 3, EMR remains the top-performing method, still outperforming SMR by a 3.3%p margin in nDCG@10.

Regarding the limited performance gain of CoT over ReasonIR in Table 4, we find that ReasonIR has effectively internalized the reasoning process of CoT. Since ReasonIR is trained specifically for reasoning-intensive retrieval, adding explicit CoT steps during inference becomes redundant. In contrast, EMR still achieves significant gains on top of ReasonIR. This demonstrates that EMR provides orthogonal benefits that even a sophisticated reasoning retriever does not possess.

Furthermore, the gains are not simply from using a more powerful LLM. Across both Table 2 and Table 3, while the Qwen3-based models generally outperform Qwen2.5-based models, EMR maintains the highest average score regardless of the underlying LLM. This indicates that episodic memory provides a robust advantage that is not dependent on initial query quality or biased towards a specific LLM.

Robustness Across Base Models. The gains of EMR are not limited to the 32B Qwen models used in our main experiments. On DeepSeek-R1, EMR improves over SMR by +5.3 nDCG@10 (17.6 to 22.9). The same trend also holds in a smaller model, where EMR improves over SMR with Qwen2.5-7B from 29.2 to 31.9. These results suggest that the effectiveness of EMR stems from its trajectory-aware memory, rather than from model scale alone. Appendix A.4 gives more details about the experiments.

Token Efficiency. Beyond performance, EMR demonstrates superior token efficiency. Figure 5 shows retrieval performance against the number of tokens consumed for the representative four datasets of the BRIGHT benchmark. In this figure, each color represents a distinct dataset, while solid lines denote EMR and dashed lines indicate SMR. See Appendix A.5 for more details.

The figure shows that EMR reaches higher nDCG@10 than SMR at comparable or lower token budgets. The added memory tokens are offset by fewer redundant reasoning steps and by memory compression.

We also perform an ablation on the compression rate by varying the top-k filtering parameter. EMR remains stable across all tested values, with less than 1.0 nDCG@10 variance overall. This suggests

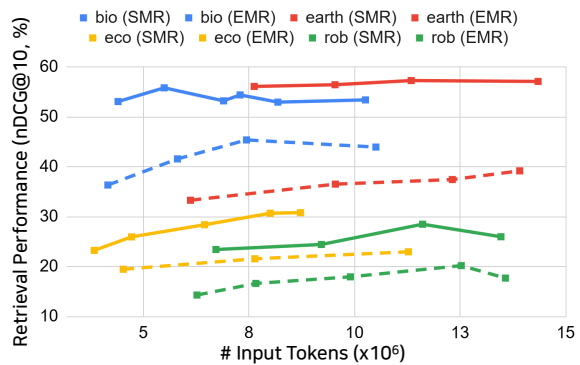


Figure 5: Retrieval performance (nDCG@10) versus cumulative input tokens on four representative BRIGHT datasets. Solid lines denote EMR and dashed lines denote SMR. Each color represents a different dataset. Across datasets, EMR achieves higher retrieval performance at comparable or lower token budgets.

that the benefits of EMR do not depend on delicate tuning of the memory budget (Appendix A.6).

To broaden the comparison, we measure cumulative token usage for CoT-based baselines. As shown in Table 5, EMR remains the most efficient method while also achieving the best retrieval performance. CoT-based approaches such as Rank1, Rank-R1, and O1-Pruner require roughly 8 times more input tokens than EMR, yet still underperform it in nDCG@10. This result confirms that EMR outperforms both SMR and CoT-style methods with improved computational efficiency.

End-to-End Latency. We measured the end-to-end wall-clock latency per query on an NVIDIA A6000 GPU using vLLM. Despite the additional memory management, EMR achieves a $\times 2.7$ speedup over SMR. While prompt construction adds a small overhead, EMR prevents the model from falling into redundant reasoning loops. These loops are the primary bottleneck in SMR. Consequently, the time saved by reducing reasoning steps far outweighs the cost of memory compression. This suggests that the efficiency gains of EMR extend beyond token counts to end-to-end runtime.

4.2.3 RQ3. Is EMR Complementary to Stronger Retrievers?

For RQ3, we test whether EMR remains effective when paired with a stronger retriever.

We test if the performance gains from EMR persist when using a more powerful retriever. We replaced the sparse retriever BM25 with ReasonIR, a strong dense retriever trained for reasoning tasks.

	Bio	Earth	Econ	Psy	Rob	Stack	Sus	Leet	Pony	AoPS	TheoQ	TheoT	Avg
Retriever													
ReasonIR	26.3	31.5	23.3	30.3	17.8	24.0	20.6	35.0	10.3	14.3	31.6	27.2	24.4
CoT-based Reasoning													
ReasonIR + Rank1	32.0	30.0	22.3	31.1	16.7	25.8	22.4	31.4	15.5	12.1	27.8	26.3	24.5
ReasonIR + Rank-R1	33.0	34.1	24.2	33.5	20.2	25.4	22.8	33.5	14.1	10.7	30.3	27.8	25.8
ReasonIR + O1-Pruner	31.6	34.9	24.5	33.2	21.0	24.9	24.7	33.3	11.8	12.9	29.9	27.1	25.8
State-based Reasoning													
ReasonIR + SMR (Qwen2.5)	34.7	35.1	26.2	32.8	20.9	25.2	24.2	30.8	10.4	13.5	30.1	28.6	26.0
ReasonIR + SMR (Qwen3)	41.7	37.8	26.4	31.2	20.3	26.0	32.5	28.8	10.5	14.8	33.0	28.5	27.6
EMR: State-based Reasoning with Episodic Memory													
ReasonIR + EMR (Qwen2.5)	41.7	41.5	31.4	36.6	22.8	25.8	34.9	32.0	11.2	15.1	38.4	34.3	30.5
ReasonIR + EMR (Qwen3)	52.4	46.8	32.2	39.9	24.7	30.2	35.3	36.1	12.2	15.3	40.7	34.6	33.4

Table 4: Retrieval performance (nDCG@10) on the **BRIGHT** benchmark using **ReasonIR** as the underlying retriever. Best scores per dataset are bolded.

Method	Tokens	nDCG@10
Rank1	×8.3	29.4
Rank-R1	×8.1	30.2
O1-Pruner	×8.2	30.3
SMR	×3.6	33.4
EMR	×1.0	36.7

Table 5: Comparison of retrieval performance and relative input token usage across CoT-based, state-based, and memory-augmented reasoning methods. ×1.0 denotes the input token cost of EMR.

Table 4 shows that the advantage of EMR remains with a significant margin. Specifically, EMR with Qwen3 achieves an average nDCG@10 of 33.4, surpassing the strongest baseline (SMR with Qwen3) by 5.8%p. This suggests that the gains from EMR and the gains from a stronger retriever are complementary.

We also compare EMR against GRITHopper (Erker et al., 2026), a strong retriever+LLM baseline in the broader multi-hop retrieval paradigm. EMR still achieves superior performance, improving nDCG@10 from 30.2 to 36.7. This comparison is notable because EMR is training-free and applied at test time, whereas GRITHopper relies on additional training. The result suggests that better trajectory control can remain competitive with trained retrieval-and-reasoning pipelines.

5 Conclusion

In this work, we address reasoning cycles in state-based approaches for reasoning-intensive IR by

introducing EMR, a framework that integrates an episodic memory module to guide the reasoning process. By conditioning the policy on its entire history of visited states, EMR reduces the cycle rate by about 94%. Experiments on the BRIGHT benchmark demonstrate that this approach consistently outperforms existing CoT-based and state-based baselines in both retrieval effectiveness and token efficiency.

Acknowledgements

This work was supported by the InnoCORE program of the Ministry of Science and ICT (N10250156), and by Institute of Information & communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. 2022-0-00077/RS-2022-II220077, AI Technology Development for Commonsense Extraction, Reasoning, and Inference from Heterogeneous Data).

Limitations

While EMR demonstrates robust performance, it has a modular design that allows for several promising extensions. For instance, the current action space is intentionally concise (Refine, Rerank, Stop). This set could easily be expanded to include more specialized tools like external API calls and user feedback, transforming EMR into a more comprehensive information searcher. Such interactions would enhance our episodic memory with a more comprehensive history.

References

- Debrup Das, Sam O’ Nuallain, and Razieh Rahimi. 2025. Rader: Reasoning-aware dense retrieval models. *arXiv preprint arXiv:2505.18405*.
- Chad DeChant. 2025. Episodic memory in ai agents poses risks that should be studied and mitigated. In *2025 IEEE Conference on Secure and Trustworthy Machine Learning (SaTML)*, pages 321–332. IEEE.
- Justus-Jonas Erker, Nils Reimers, and Iryna Gurevych. 2026. Grithopper: Decomposition-free multi-hop dense retrieval. In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 77–94.
- Yair Feldman and Ran El-Yaniv. 2019. Multi-hop paragraph retrieval for open-domain question answering. *arXiv preprint arXiv:1906.06606*.
- Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*.
- Dohyeon Lee, Yeonseok Jeong, and Seung-won Hwang. 2025. From token to action: State machine reasoning to mitigate overthinking in information retrieval. *arXiv preprint arXiv:2505.23059*.
- Haotian Luo, Li Shen, Haiying He, Yibo Wang, Shiwei Liu, Wei Li, Naiqiang Tan, Xiaochun Cao, and Dacheng Tao. 2025. O1-pruner: Length-harmonizing fine-tuning for o1-like reasoning pruning. *arXiv preprint arXiv:2501.12570*.
- Xueguang Ma, Liang Wang, Nan Yang, Furu Wei, and Jimmy Lin. 2024. Fine-tuning llama for multi-stage text retrieval. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2421–2425.
- Andrew M Nuxoll and John E Laird. 2012. Enhancing intelligent agents with episodic memory. *Cognitive Systems Research*, 17:34–48.
- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Rulin Shao, Rui Qiao, Varsha Kishore, Niklas Muennighoff, Xi Victoria Lin, Daniela Rus, Bryan Kian Hsiang Low, Sewon Min, Wen-tau Yih, Pang Wei Koh, et al. 2025. Reasonir: Training retrievers for reasoning tasks. *arXiv preprint arXiv:2504.20595*.
- Hongjin Su, Howard Yen, Mengzhou Xia, Weijia Shi, Niklas Muennighoff, Han-yu Wang, Haisu Liu, Quan Shi, Zachary S Siegel, Michael Tang, et al. 2024. Bright: A realistic and challenging benchmark for reasoning-intensive retrieval. *arXiv preprint arXiv:2407.12883*.
- Orion Weller, Kathryn Ricci, Eugene Yang, Andrew Yates, Dawn Lawrie, and Benjamin Van Durme. 2025. Rank1: Test-time compute for reranking in information retrieval. *arXiv preprint arXiv:2502.18418*.
- Chenghao Xiao, G Thomas Hudson, and Noura Al Moubayed. 2024. Rar-b: Reasoning as retrieval benchmark. *arXiv preprint arXiv:2404.06347*.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. 2018. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*.
- Shengyao Zhuang, Xueguang Ma, Bevan Koopman, Jimmy Lin, and Guido Zuccon. 2025. Rank-r1: Enhancing reasoning in llm-based document rerankers via reinforcement learning. *arXiv preprint arXiv:2503.06034*.

A Appendix

A.1 Baseline Methods

To provide a comprehensive evaluation of EMR, we compare it against various baseline methods.

CoT-based Reasoning. We include CoT-based baselines as they are standard approach in the literature and represent a widely adopted methodology, ensuring our evaluation is grounded in established practices.

- **Rank1** (Weller et al., 2025) and **Rank-R1** (Zhuang et al., 2025) are strong CoT-based reasoning models. Rank1 is fine-tuned on reasoning traces, while Rank-R1 further enhances this capability through reinforcement learning. We select these models over earlier methods like ReAct (Yao et al., 2023) because they are specifically adapted for IR and represent the leading CoT-based methods in this domain.
- **O1-Pruner** (Luo et al., 2025) focuses on the efficiency of CoT. It employs reinforcement learning to compress the lengthy reasoning paths generated by CoT models, aiming to reduce computational overhead without sacrificing the final answer quality.

State-based Reasoning. We include state-based baselines to compare EMR, aiming to demonstrate superior overall performance and how our proposed episodic memory component enhances the efficiency of state transitions.

- **SMR** (Lee et al., 2025) frames the multi-step retrieval process as a Markov Decision Process. At each step, a policy model observes the current state (e.g., the query and retrieved documents) and selects a discrete action to transition to the next state, finally building a set of relevant documents for the answer.

Retriever. We include both a traditional sparse retriever and a strong dense retriever to investigate the synergy between the retrieval backbone and EMR. This comparison is crucial to demonstrate that the benefits of EMR are orthogonal to the retriever’s capabilities. Our method provides consistent improvements even when paired with a more powerful retriever.

- **BM25** (Robertson et al., 2009) serves as our traditional sparse retrieval baseline. It is a keyword-based algorithm that has been a long-standing and robust baseline in IR for decades, relying on term frequency and inverse document frequency.
- **ReasonIR** (Shao et al., 2025) represents a strong baseline in dense retrieval for reasoning-intensive IR. It is highly effective for tasks that require understanding context and intent beyond simple keyword matching. Its strong performance over other notable dense retrievers like Contriever (Izacard et al., 2021) and RankLLaMA (Ma et al., 2024) makes it a challenging and relevant baseline.

A.2 Implementation Details

Models and Adaptations. The baseline models used in our experiments are derived from various LLMs. Rank1-32B utilizes Qwen2.5-32B¹ as its backbone, while Rank-R1-14B is built upon Qwen2.5-14B². Due to the absence of a 32B version for Rank-R1, our experiments employ the 14B model. For O1-Pruner, which is not inherently a retrieval model, we implemented a custom prompt to enable query rewriting and document reranking functionalities; its 32B version originates from QwQ-32B-preview³. To ensure a fair and comprehensive evaluation, we also benchmark against the base LLMs, Qwen2.5-32B¹ and Qwen3-32B⁴.

Experimental Settings. We conducted all experiments on a single NVIDIA A6000 GPU using the vLLM engine. Unless otherwise noted, hyperparameters remained consistent across all datasets: a batch size of 8, LLM temperature of 0.1, and a maximum of 16 reasoning steps. For retrieval, we set the top- k value to 10. We define computational cost as the total number of input tokens consumed throughout the entire reasoning process.

To evaluate temporal efficiency, we measured the end-to-end wall-clock latency per query. This metric covers the entire pipeline, including document retrieval, prompt construction, and the compression process. We also used the same system prompt (Table 6) as SMR for our policy model to ensure a fair comparison and isolate the improvements of our method. This setup prevents performance gains from being attributed to prompt variations rather than our core mechanism.

A.3 Evaluation on Other Metrics

To further demonstrate the robustness of our proposed method, this section supplements the primary nDCG@10 evaluation with results on additional metrics. We evaluated performance on the BRIGHT benchmark using BM25 as the retriever, measuring results with two additional metrics: MAP@10 (shown in Table 7) and Recall@10 (shown in Table 8). Consistent with our main findings, EMR maintains a clear performance advantage over both CoT and compressed CoT baselines across these metrics.

A.4 Additional Backbone and Baseline Comparisons

To verify that the gains of EMR are not specific to the 32B Qwen backbones used in our main experiments, we evaluate additional backbone settings in Table 9. First, on DeepSeek-R1, EMR improves nDCG@10 from 17.6 to 22.9 when replacing SMR. Second, in a smaller-model setting with Qwen2.5-7B, EMR improves over SMR from 29.2 to 31.9. These results suggest that the benefits of episodic memory are not a simple consequence of model scale, but arise from a more general improvement to the state-based reasoning process.

A.5 Token Efficiency

We assess the token efficiency of EMR by comparing its token consumption to that of SMR under the same experimental setup (top- $k=10$). Table 10 shows that EMR is more efficient, requiring about 72% fewer tokens on average. This demonstrates that the episodic memory in EMR enhances retrieval performance while achieving superior token efficiency.

A.6 Ablation on Memory Compression Rate

To evaluate the sensitivity of EMR to the memory compression budget, we perform an ablation on the top- k filtering parameter used in Section 3.2.3. Table 11 reports the retrieval performance for different values of k . Across all tested settings, EMR remains stable with less than 1.0 nDCG@10 variance. Even at the lowest-performing setting, EMR still outperforms the strongest baseline SMR. This result indicates that EMR is robust to the compression rate and does not require sensitive hyperparameter tuning.

A.7 Use of AI Assistants

We used AI assistants only for limited support, including code debugging, grammar checking, and minor phrasing suggestions.

¹Qwen/Qwen2.5-32B-Instruct

²Qwen/Qwen2.5-14B-Instruct

³Qwen/QwQ-32B-Preview

⁴Qwen/Qwen3-32B-FP8

You are a highly intelligent artificial agent responsible for managing a search system. Your role is to either refine the given query or re-rank retrieved search results, thereby enhancing both recall and precision of the search. You can output exactly one of the following operations, after which another agent will execute it and return the results to you.

Input Format

The input provided to you will have the following structure:

```
...
{
  "query": "<current version of a query>",
  "ranks": [
    (<docid>, "document contents"),
    (<docid>, "document contents"),
    ...
  ]
}
...
```

Decision policy (check in order):

1. Query Refinement

Choose "action = refine" if any of the following are met:

- The query is ambiguous or generic
- The retrieved search results are unsatisfactory
- The query is short
- Key domain terms are missing in the query

2. Re-ranking

Choose "action = rerank" only if the query already looks good and at least one retrieved document seems on-topic.

3. Stop

Choose "action = stop" only when you are *certain* that no further improvement is possible.

Possible Outputs (select exactly one)

Query Refinement

You may refine the query by rewriting it into a clear, specific, and formal version that is better suited for retrieving relevant information from a list of passages. Only return the document IDs (``docid``) in the ``reranked`` list. Do not include document contents. Output format:

```
...
{
  "action": "refine",
  "query": "<refined version of a query>",
  "reason": "<reason for this action>"
}
...
```

Re-ranking

You may reorder the retrieved documents (do not remove non-relevant ones). The results should be sorted in descending order of relevance. Output format:

```
...
{
  "action": "rerank",
  "ranks": [<docid>, <docid>, ...],
  "reason": "<reason for this action>"
}
...
```

Stop

You may stop this iteration when the results are satisfactory. Output format:

```
...
{
  "action": "stop",
  "reason": "<reason for this action>"
}
...
```

Table 6: System prompt used in EMR.

	Bio	Earth	Econ	Psy	Rob	Stack	Sus	Leet	Pony	AoPS	TheoQ	TheoT	Avg
Retriever													
BM25	13.1	18.1	8.9	7.0	8.5	13.4	10.0	19.7	1.6	6.2	10.4	4.9	10.2
CoT-based Reasoning													
BM25 + Rank1	16.1	23.0	7.9	11.3	10.5	12.3	16.2	17.6	2.2	3.1	10.4	8.6	11.6
BM25 + Rank-R1	17.7	25.0	10.3	10.7	14.1	13.1	16.1	20.2	2.1	4.9	10.8	7.7	12.7
BM25 + O1-Pruner	17.2	22.6	10.7	12.9	12.7	14.9	15.5	19.5	1.8	3.7	8.9	5.6	12.2
State-based Reasoning													
BM25 + SMR (Qwen2.5)	22.0	25.7	13.6	15.4	15.0	15.5	16.6	21.0	1.6	3.1	14.1	15.2	14.9
BM25 + SMR (Qwen3)	23.1	24.4	14.0	15.7	15.3	16.2	17.0	21.8	1.9	3.4	16.2	15.4	15.4
EMR: State-based Reasoning with Episodic Memory													
BM25 + EMR (Qwen2.5)	33.9	42.1	16.1	26.3	15.2	17.3	19.7	22.0	1.9	4.0	24.2	25.8	20.7
BM25 + EMR (Qwen3)	44.3	44.5	15.4	28.2	17.3	20.5	24.6	10.4	2.2	4.2	19.1	15.2	20.5

Table 7: Retrieval performance (MAP@10) on the **BRIGHT** benchmark using **BM25** as the underlying retriever. Best scores per dataset are bolded.

	Bio	Earth	Econ	Psy	Rob	Stack	Sus	Leet	Pony	AoPS	TheoQ	TheoT	Avg
Retriever													
BM25	21.7	31.9	16.8	15.6	19.6	21.3	21.3	29.5	4.1	6.0	8.6	9.2	17.1
CoT-based Reasoning													
BM25 + Rank1	21.7	31.9	16.8	15.6	19.6	21.3	21.3	29.5	4.1	6.0	8.6	9.2	17.1
BM25 + Rank-R1	21.7	31.9	16.8	15.6	19.6	21.3	21.3	29.5	4.1	6.0	8.6	9.2	17.1
BM25 + O1-Pruner	23.8	34.4	20.2	19.8	19.6	21.3	21.5	29.5	2.3	6.0	12.9	10.3	18.5
State-based Reasoning													
BM25 + SMR (Qwen2.5)	25.3	33.0	19.5	21.7	21.4	24.6	21.6	26.6	3.6	6.1	20.9	24.8	20.8
BM25 + SMR (Qwen3)	25.9	34.1	22.9	22.4	20.8	25.0	22.1	27.2	3.7	6.6	21.3	25.1	21.1
EMR: State-based Reasoning with Episodic Memory													
BM25 + EMR (Qwen2.5)	38.3	47.0	23.3	34.5	17.4	29.1	26.8	30.7	4.3	6.4	26.0	30.7	26.2
BM25 + EMR (Qwen3)	52.8	53.7	23.2	38.3	24.3	32.5	31.2	25.6	4.3	6.9	21.0	17.9	27.6

Table 8: Retrieval performance (Recall@10) on the **BRIGHT** benchmark using **BM25** as the underlying retriever. Best scores per dataset are bolded.

	Bio	Earth	Econ	Psy	Rob	Stack	Sus	Leet	Pony	AoPS	TheoQ	TheoT	Avg
On DeepSeek-R1													
DeepSeek-R1 + SMR	28.5	39.2	18.8	17.1	11.9	17.6	20.7	26.2	6.5	10.0	7.7	6.7	17.6
DeepSeek-R1 + EMR	49.7	50.7	18.8	34.1	19.8	24.4	27.5	11.2	4.8	2.5	17.6	13.7	22.9
On Smaller Backbone (7B)													
Qwen2.5 (7B) + SMR	47.9	49.8	25.8	44.4	19.6	27.3	26.0	26.3	16.0	4.9	35.2	27.6	29.2
Qwen2.5 (7B) + EMR	51.6	52.4	30.1	43.7	24.2	27.2	30.6	27.6	18.2	7.3	35.3	34.3	31.9

Table 9: Retrieval performance (nDCG@10) on additional backbone settings. EMR consistently improves over SMR on both DeepSeek-R1 and Qwen2.5-7B.

Dataset	SMR	EMR	Reduction
bio	10.5	4.4	-58.1%
earth	40.7	5.0	-87.7%
econ	14.3	3.8	-73.2%
psy	22.7	3.4	-85.1%
rob	18.5	6.7	-63.8%
stack	30.2	5.1	-83.0%
sus	13.8	5.1	-62.8%
leet	6.5	5.0	-22.6%
pony	5.0	2.0	-60.1%
aops	3.3	2.9	-11.0%
theoq	6.2	4.3	-31.0%
theot	3.9	1.3	-66.9%
Average	14.6	4.1	-72.1%

Table 10: Comparison of token usage (in millions) between EMR and SMR across various datasets in the BRIGHT benchmark. The values represent the total number of input tokens required per task. The final column shows the average of all datasets.

k	nDCG@10
3	26.2
5	26.7
7	26.8
9	26.3
11	25.9

Table 11: Ablation on the top-k memory compression rate. EMR remains stable across different values of k, with less than 1.0 nDCG@10 variance.