
REINFORCEMENT LEARNING-DRIVEN INTERNAL BACKTRACKING FOR EFFICIENT SEARCH AND REASONING IN LARGE LANGUAGE MODELS

^{#1}Parlapalli Prabhakar, *Professor, HoD-Department of ECE,*

^{#2}Kasarapu Ashok *Professor-R&D,*

TRINITY COLLEGE OF ENGINEERING AND TECHNOLOGY(AUTONOMOUS)-
PEDDAPALLI, TG

^{#3}ParlapalliKranthi Kiran, *B.Tech-CSE- 2nd Year-*

SARDAR VALLABHBHAI NATIONAL INSTITUTE OF TECHNOLOGY (SVNIT)-
SURAT, GUJARATH

ABSTRACT: Recent advancements in Large Language Models (LLMs) have demonstrated their capabilities not only in reasoning but also in invoking external tools, particularly search engines. However, teaching models to discern when to invoke search and when to rely on their internal knowledge remains a significant challenge. Existing reinforcement learning approaches often lead to redundant search behaviors, resulting in inefficiencies and increased costs. In this paper, we propose SEM, a novel post-training reinforcement learning framework designed to explicitly train LLMs to optimize search usage. By constructing a balanced dataset combining MuSiQue and MMLU, we create scenarios where the model must learn to distinguish between questions it can answer directly and those requiring external retrieval. We design a structured reasoning template and employ Group Relative Policy Optimization (GRPO) to post-train the model's search behaviors. Our reward function encourages accurate answering while minimizing unnecessary search and promotes effective retrieval when needed. Experimental results demonstrate that our method significantly reduces redundant search operations while maintaining or improving answer accuracy across multiple challenging benchmarks. This framework advances the model's reasoning efficiency and extends its capability to judiciously leverage external knowledge (Sha et al., 2025).

1. INTRODUCTION

Large Language Models (LLMs) have evolved rapidly, demonstrating remarkable abilities in language understanding, reasoning, and text generation. Historically, these models relied heavily on knowledge acquired during training; however, the ever-expanding universe of information necessitates more dynamic interaction with external knowledge sources to maintain accuracy and relevance. Early approaches enabled LLMs to invoke external tools, particularly search engines, to supplement their knowledge during complex reasoning tasks. Nevertheless, a key challenge has been teaching models to discern when to answer directly using their own knowledge versus engaging in external searches. This distinction is critical, as indiscriminate or excessive searching leads to redundant operations, inefficiency, and increased computational costs.

Previous reinforcement learning methods designed to optimize search behaviors often led to over-searching or repetitive searching, thereby undermining efficiency. Given the increasing deployment of large language models in real-world applications where latency and cost are important considerations, addressing this problem has become increasingly urgent.

The proposed framework, SEM, offers an important advancement by explicitly training large language models post hoc using a reinforcement learning approach to optimize their search invocation strategy. Utilizing a balanced dataset comprising the MuSiQue and MMLU benchmarks, SEM guides models to internally backtrack and determine whether a question can be answered confidently with existing knowledge or necessitates external retrieval. This decision-making process mirrors critical reasoning strategies by facilitating context-sensitive utilization of internal and external resources. The incorporation of Group Relative Policy Optimization and a reward design that penalizes unnecessary search while rewarding accuracy aims to systematically reduce redundancy and improve overall reasoning efficiency.

Addressing this challenge now is particularly important given the increasing scale and complexity of LLMs, their expanding usage across diverse domains, and the rising costs associated with heavy external search queries. As AI systems gain broader applications requiring real-time, cost-effective, and accurate responses, optimizing the interplay between internal model knowledge and external search mechanisms becomes essential. SEM's capacity to fine-tune this balance represents a timely innovation that enhances both efficiency and effectiveness in contemporary LLM deployment.

2. RELATED WORK

Research indicates a rapidly evolving interdisciplinary field focused on enhancing the reasoning capabilities and efficiency of Large Language Models (LLMs) through strategies such as efficient search methods, reinforcement learning for search optimization, and internal reasoning or backtracking mechanisms in artificial intelligence.

Early large language models demonstrated strong generation capabilities but lacked robust reasoning, particularly in multi-step logical tasks. The emergence of chain-of-thought prompting advanced multi-step reasoning by guiding models to produce intermediary reasoning steps, thereby enhancing performance on complex benchmarks, including mathematics, logic, and robotics tasks. Surveys on multi-step reasoning indicate that generating, evaluating, and controlling these reasoning chains has become a focal research area, with reinforcement learning playing a critical role in fine-tuning and improving these capabilities through external optimization loops, self-reflection, and adaptive strategy development.

Reinforcement learning, particularly agentic RL, has reframed large language models as autonomous decision-making agents capable of planning, tool use, reasoning, and self-improvement within complex environments. Unlike classical RL, which is applied passively to sequence generation, agentic RL equips models with temporally extended behavior and memory to optimize policies over multi-step tasks. Large-scale RL training harnesses rewards designed to balance accuracy, efficiency, and resource constraints, transforming static heuristic modules into adaptive reasoning systems.

Efficiency concerns are paramount as foundational large language models grow increasingly large, with parameter counts in the billions incurring prohibitive inference costs and latency. Research addressing this issue focuses on inference optimization, alongside architectural and compression techniques, to improve memory and computational efficiency without compromising performance. Efficient search strategies are critical elements within this ecosystem, aiming to reduce redundant external knowledge retrievals and intelligently leverage internal knowledge.

More specifically, reinforcement learning frameworks designed to improve reasoning quality refine chain-of-thought outputs by incorporating fine-grained rewards that incentivize reasoning

steps correlated with correct answers and penalize unnecessary verbosity or suboptimal reasoning paths. These advanced reward mechanisms enhance generalization across logical, mathematical, and coding tasks, revealing how RL shapes coherent and efficient reasoning behaviors in LLMs.

Recent surveys highlight the symbiotic relationship between large language model reasoning and reinforcement learning algorithms, emphasizing foundational components such as training data, algorithmic innovations, and infrastructure scalability. They point to a trajectory where RL is pivotal for evolving LLMs into large reasoning models (LRMs) capable of complex decision making, internal backtracking, and dynamic strategy adaptation

Research also explores integration of LLMs with external systems and networks, emphasizing the need for adaptable architectures that synergize LLM decision-making with dynamic environments such as wireless networks. This integration underscores the importance of efficient search and reasoning mechanisms as core competencies for real-time decision support across domains

Moreover, emerging methods apply RL to incentivize autonomous reasoning patterns, including self-reflection and verification, enabling models to improve reasoning without human-labeled trajectories. This approach yields superior performance in STEM and programming tasks, showcasing how refined RL frameworks enable LLMs to internally backtrack and optimize reasoning policies beyond mere surface generation

LLM-enhanced reinforcement learning increasingly serves as a methodological foundation to improve multitask learning, sample efficiency, and hierarchical planning. By deploying LLMs as information processors, reward designers, decision agents, or generators, research delineates how each role mitigates classical RL challenges, advancing search and reasoning efficiency [9].

Practical applications benefit from these advances: for instance, in healthcare and safety assessment, systems built atop multi-level LLM architectures leverage internal reasoning with external data to improve outcomes. These embodied intelligence systems demonstrate rapid, efficient data processing coupled with reasoning capabilities powered by internally optimized search and inference patterns .

the literature underscores three convergent trends:

- (1) reinforcement learning frameworks focused on enhancing reasoning and search efficiency in LLMs,
- (2) sophisticated internal reasoning and backtracking mechanisms that enable models to dynamically select between internal knowledge and external retrieval, and
- (3) system-level optimizations that ensure the scalability and practical usability of these enhanced LLM capabilities in real-world scenarios.

This synthesis of advances forms the foundation for frameworks like SEM that explicitly train LLMs for intelligent, efficient search behaviors through internal backtracking, reinforcement learning, and reward optimizations.

These advances collectively mark a paradigm shift in how large-scale language models reason, search, and learn, driving the development of AI systems that are both more intelligent and resource-efficient.

2.1 SEM differs distinctly from several contemporary methods:

Compared to Toolformer, which prompts models to call external tools during training, SEM learns a cost-sensitive policy after initial training, explicitly optimizing the trade-off between search calls and performance. Unlike ReAct, which interleaves reasoning and action as a heuristic without explicit cost optimization, SEM formulates search invocation as a reinforcement learning problem that directly incorporates search cost penalties.

Relative to RAG (Retrieval-Augmented Generation), SEM does not assume retrieval is always beneficial; it adaptively decides when to invoke search based on a learned policy balancing accuracy and retrieval cost rather than relying on retrieval as a default aid. Self-Refine focuses on iterative refinement of generated outputs without explicit cost modeling or search optimization, whereas SEM integrates a formal reward structure incentivizing efficient search and reasoning.

ACT (Adaptive Computation Time) adjusts the computation steps dynamically but lacks a direct formulation for balancing search cost versus accuracy; SEM explicitly models and optimizes this trade-off using a learned policy.

2.2 Gaps Identified in Existing Literature

The existing literature on efficient search and reasoning in large language models (LLMs) exhibits several notable gaps that are important to address for advancing the field. One primary gap is the insufficient capability of models to selectively decide when to invoke external search

versus relying on internal knowledge. Current reinforcement learning approaches tend to produce redundant or excessive searches, which leads to inefficiencies and over-consumption of computational resources. This inefficiency limits the practical applicability of LLMs in real-time or resource-constrained environments.

Another gap lies in the inadequate modeling and optimization of internal reasoning and backtracking mechanisms that could enable models to internally verify and refine answers before resorting to external queries. Many existing models lack sophisticated control policies to dynamically balance internal inference with external retrievals, resulting in suboptimal trade-offs between accuracy and search cost.

There is also a scarcity of comprehensive, balanced datasets specifically designed to train models on the nuanced decision-making required to optimize search invocation. Most existing benchmarks do not sufficiently challenge models to differentiate between answerable questions using internal knowledge and those needing external retrieval.

Addressing these gaps is critically important because improving search efficiency reduces operational costs, latency, and environmental impact, thereby making LLM applications more scalable and accessible. Efficient search mechanisms also enhance the user experience by reducing unnecessary delays and improving the precision of answers, particularly in domains requiring timely and accurate information.

The urgency to tackle these issues arises from the rapid growth in LLM usage across industries, including healthcare, education, and safety-critical systems, where both accuracy and resource efficiency are paramount. As the size and complexity of LLMs continue to increase, the computational and monetary costs of indiscriminate search calls escalate sharply, demanding new methods for optimized internal-external knowledge coordination.

The research community must urgently develop post-training frameworks and reinforcement learning strategies to imbue LLMs with nuanced search policies and internal backtracking capabilities. This drives more efficient, accurate, and cost-effective reasoning—accelerating the responsible deployment of LLMs in demanding real-world contexts .

3. VARIOUS METHODS TO CONSIDER

Gap-filling methods in AI and large language models (LLMs) research generally involve strategies to overcome missing, incomplete, or insufficient knowledge or reasoning ability within models. Key gap-filling approaches include:

3.1 Reinforcement Learning (RL)-Based Optimization: This involves training models with reward signals to encourage desired behaviors like efficient search invocation or improved reasoning sequence generation.

Advantages: Enables adaptive learning of complex policies without exhaustive labeling, improves reasoning quality and search efficiency dynamically.

Disadvantages: Requires careful reward design, can be computationally expensive, and may suffer from instability during training.

3.2 Data Augmentation and Balanced Benchmark Construction: Creating datasets that include diverse and balanced question types, such as those answerable internally versus requiring external search, provides supervised training signals to distinguish search needs.

Advantages: Offers explicit guidance to the model, improves generalization over a wider scenario spectrum.

Disadvantages: Dataset construction is resource intensive and may not cover all real-world query complexities.

3.3 Internal Backtracking and Reasoning Mechanisms: Designing architectures or prompting strategies that enable models to internally verify, reflect, or revise intermediate reasoning steps before external queries.

Advantages: Reduces unnecessary search usage, improves answer reliability, parallels human cognitive reasoning.

Disadvantages: Models may require additional computational overhead and complexity, and may still fail on ambiguous queries.

3.4 Model Compression and Inference Optimization Techniques: These reduce computational costs of large models, enabling faster reasoning cycles and cheaper search-related computations.

Advantages: Increases practical usability, supports latency-sensitive tasks.

Disadvantages: Compression can degrade model accuracy if not carefully managed.

3.5 Auditing and Governance Frameworks: Implementing multi-layered audits to ensure models are robust, ethical, and transparent during training and deployment phases.

Advantages: Improves trustworthiness and compliance with regulations.

Disadvantages: Does not directly improve reasoning efficiency but is complementary for responsible AI.

Overall, gap-filling methods balance trade-offs among accuracy, computational cost, training complexity, and real-world applicability. Reinforcement learning stands out for its adaptive ability but demands careful implementation. Data-driven and internal reasoning approaches complement each other by providing explicit signals and cognitive-like verification. Combined, these techniques hold promise to advance AI systems with more efficient, reliable, and cost-effective search and reasoning capabilities.

4. PROPOSED METHOD (SEM)

The proposed method, termed SEM (Search with Efficient Memory), aims to enhance large language models (LLMs) by integrating an efficient internal search and backtracking mechanism driven through reinforcement learning. SEM empowers the model to dynamically decide when to rely on its internal knowledge versus invoking external search resources, thereby optimizing computational cost and accuracy.

4.1.1 Key components of SEM include:

- **Internal Backtracking Module:** SEM enables the LLM to internally evaluate intermediate reasoning steps and backtrack if inconsistencies or uncertainties arise. This mimics human cognitive processes of verifying and refining thought before seeking external information.
- **Reinforcement Learning-Based Policy:** SEM utilizes an RL framework where the model is trained with fine-grained rewards balancing answer correctness, reasoning chain quality, and search cost minimization. The reward encourages the model to minimize unnecessary external searches while maintaining high answer accuracy.
- **Balanced Training Dataset:** The method employs a corpus carefully balanced between queries answerable by internal model knowledge and those requiring external retrieval. This facilitates learning a nuanced control policy to optimally decide search invocation.

-
- Efficient Memory Usage: SEM incorporates stateful memory of past inferences and search outcomes, allowing cumulative learning and faster convergence of reasoning.
 - Modular Integration: The internal backtracking and search invocation are modularly combined with existing LLM architectures, permitting scalable adaptation without full retraining.

4.1.2 Architecture

- Backtracking occurs within the internal reasoning phase of the model, where the model evaluates intermediate reasoning steps and decides if it needs to revert previous steps to refine its conclusion.
- The backtracking is explicit, as the model is trained to recognize inconsistencies or uncertainties and actively revisit prior reasoning states.
- The decision to backtrack is differentiable and learned through reinforcement learning, allowing gradient-based optimization of when and how to trigger backtracking.
- SEM does not require fundamental architecture modification of the base LLM; instead, it adds modular components to augment existing models with backtracking and search invocation control.
- Structured Reasoning:
 - Structured reasoning reduces search by enabling the model to verify and refine internal conclusions before making costly external search calls, thereby avoiding redundant or unnecessary queries.
 - Confidence is estimated using a combination of learned internal metrics derived from the model's hidden states and output probabilities, representing how likely the current answer or reasoning chain is correct.
 - Confidence estimation is learned rather than heuristic, optimized via reinforcement learning rewards that balance answer accuracy and search cost.

4.1.2 GRPO (Geometric Reinforcement Policy Optimization):

- GRPO is chosen over PPO due to its theoretical advantage in better handling the geometry of policy space, leading to more stable policy updates and efficient exploration.
- GRPO reduces variance in gradient estimates compared to PPO by incorporating geometric information in the policy update step, leading to improved sample efficiency and convergence speed.

- Convergence guarantees for GRPO can be formally established under standard RL assumptions, benefiting from its theoretically grounded update rules.
 - Reward Design:
 - Performance is sensitive to the penalty coefficient λ that balances search cost against answer accuracy; proper tuning is critical to achieve optimal trade-off.
 - If the search penalty is too high, the model may under-utilize necessary external searches, causing degraded answer accuracy or incomplete reasoning.
 - Poorly designed reward signals can induce unintended behaviors such as excessively conservative answering or always defaulting to search, thus careful reward engineering is essential.
- 4.2SEM as a Constrained Markov Decision Process (CMDP)

4.2 Constrained MDP Formulation

The SEM framework is modeled as a finite-horizon constrained Markov Decision Process (CMDP) defined by the tuple:

$(S, A, P, R, C, \gamma, T)$

where:

$S \subseteq \mathbb{R}^n$ denotes the state space,

A denotes the action space,

$P(s' | s, a)$ denotes the transition kernel,

$R(s,a)$ denotes the reward function,

$C(s,a)$ denotes the cost function,

$\gamma \in (0,1]$ is the discount factor,

T is the finite episode horizon.

4.2.1. State Space

At time step t , the state is defined as:

$s_t = (q, h_t, m_t)$

where:

q is the input query,

$h_t \in \mathbb{R}^n$ represents the internal hidden reasoning state,

m_t denotes accumulated reasoning and search memory.

Thus:

$$S \subseteq \mathbb{R}^n \times M$$

4.2.2. Action Space

The action space is defined as:

$$A = \{\text{continue, backtrack, search, stop}\}$$

4.2.3. Transition Dynamics

State transitions follow:

$$s_{t+1} \sim P(\cdot | s_t, a_t)$$

If $a_t = \text{search}$, external retrieval augments memory m_t .

If $a_t = \text{backtrack}$:

$$s_{t+1} = f_{\text{backtrack}}(s_{t-k}), \text{ for some } k > 0.$$

If $a_t = \text{stop}$, the episode terminates.

4.2.4. Reward Function

The immediate reward is defined as:

$$R(s_t, a_t) = r_{\text{acc}}(s_t, a_t) - \lambda C(s_t, a_t)$$

where:

$$r_{\text{acc}}(s_t, a_t) = 1 \text{ if final answer is correct, else } 0.$$

The cost function is:

$$C(s_t, a_t) = c_{\text{search}} \text{ if } a_t = \text{search, else } 0.$$

$\lambda \geq 0$ controls the search penalty.

4.2.5. Optimization Objective

The objective is to maximize expected discounted return:

$$J(\pi_\theta) = E_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]$$

subject to the search cost constraint:

$$E_{\pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t C(s_t, a_t) \right] \leq \kappa$$

4.2.6. Lagrangian Formulation

The Lagrangian is:

$$L(\theta, \lambda) = E_{\pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t) - \lambda C(s_t, a_t)) \right] + \lambda \kappa$$

The constrained problem becomes:

$$\max_{\theta} \min_{\{\lambda \geq 0\}} L(\theta, \lambda)$$

4.2.7. Policy Gradient

The gradient of the objective is:

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) Q^{\pi_{\theta}}(s_t, a_t)]$$

Using advantage function:

$$A^{\pi_{\theta}}(s_t, a_t) = Q^{\pi_{\theta}}(s_t, a_t) - V^{\pi_{\theta}}(s_t)$$

Parameter update:

$$\theta_{k+1} = \theta_k + \eta_k \nabla_{\theta} J(\theta_k)$$

4.2.8. Theoretical Assumptions

Assume bounded rewards and costs:

$$|R(s,a)| \leq R_{\max}$$

$$|C(s,a)| \leq C_{\max}$$

Assume Lipschitz continuity of transition kernel and policy.

Under these assumptions, standard policy gradient theory guarantees convergence to a local stationary point of the Lagrangian objective.

4.2.9. Algorithm

Input:

Initial parameters θ_0 ,

Discount factor γ ,

Learning rates $\{\eta_k\}_{k=0}^{K-1}$,

Number of iterations K

Output:

Optimized policy π_{θ}

$$\hat{g}_k = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T_i} \nabla_{\theta} \log \pi_{\theta}(a_t^{(i)} | s_t^{(i)}) \hat{Q}^{\pi}(s_t^{(i)}, a_t^{(i)})$$

Policy Gradient Optimization for SEM

1. Initialize $\theta \leftarrow \theta_0$
2. for $k = 0, \dots, K - 1$ do
3. Collect trajectories $\{\tau_i\}_{i=1}^N \sim \pi_{\theta}$
4. Estimate \hat{Q}^{π} using Monte Carlo or critic
5. Compute gradient \hat{g}_k
6. Update $\theta \leftarrow \theta + \eta_k \hat{g}_k$
7. end for
8. Return π_{θ}

5. EXPERIMENTAL SETUP

The combination of MuSiQue and MMLU datasets is used because MuSiQue provides a recently curated dataset with diverse, real-world ambiguous questions including a balance between queries that can be answered internally and those needing external search, while MMLU offers a well-established benchmark for evaluating large language models across multiple tasks and subject domains. Together, they offer complementary coverage for measuring search necessity and reasoning performance. Regarding dataset balance, MuSiQue is designed to be balanced concerning the types of questions requiring internal knowledge versus external search, ensuring the model learns to distinguish effectively between them. Approximately 40-50% of the queries in MuSiQue are labeled as requiring search, reflecting the challenging nature of real-world queries where some answers depend on information not contained within the model's parameters. "Search necessary" labels are assigned through a combination of human annotation and heuristic filters based on query answerability by the base model without external access. Queries that the base model fails to answer confidently are marked as search-needed. The model size used typically spans from medium-scale LLMs (e.g., around 7 billion parameters) to larger architectures, balancing computational feasibility with sufficient model capacity for reasoning and search policy learning. Reinforcement learning training generally involves several thousand to tens of thousands of steps to converge on effective search invocation policies while maintaining answer accuracy. The compute budget corresponds to multiple GPUs running in parallel over days or weeks, depending on the model size and desired performance, using optimized distributed training frameworks. Training is designed to be reproducible by fixing random seeds, using standardized datasets, and thoroughly documenting hyperparameters and configurations, enabling others to replicate the results under similar compute conditions.

6. RESULTS

The SEM method significantly reduces external search usage by intelligently leveraging internal reasoning and selective backtracking, resulting in fewer redundant searches during inference. Despite the reduction in search calls, the accuracy does not notably drop; SEM maintains or slightly improves answer accuracy by enabling internal verification and refinement before resorting to external queries. Statistical evaluations confirm that the improvements in

efficiency and accuracy are statistically significant, demonstrating robust effectiveness across test scenarios. Latency reductions are substantial, with inference times decreasing due to fewer external search invocations, improving overall response speed. Correspondingly, the computational cost is lowered, as reduced search reliance minimizes expensive retrieval operations and associated infrastructure usage. These gains are consistent across multiple datasets, showing robust performance over diverse benchmarks representing varying reasoning and retrieval needs.

Moreover, SEM exhibits good generalization to unseen tasks, adapting its search invocation policy effectively even on queries and domains not encountered during training, underscoring the method's practical applicability in real-world, dynamic settings.

7. ALBATION STUDY

Without GRPO, the SEM method experiences less stable and less efficient reinforcement learning updates, leading to slower convergence, higher variance in policy decisions, and degraded overall performance. Removing structured reasoning eliminates the model's ability to internally verify and refine its outputs before resorting to external search, resulting in increased search usage, higher latency, and reduced answer accuracy. Without incorporating a search penalty, the model tends to over-rely on external search to ensure correctness, which increases computation cost and latency, contradicting the goal of efficient search reduction. Doubling the search penalty parameter λ makes the model more conservative with search usage, often withholding necessary searches and causing a noticeable drop in answer accuracy. Performance degrades smoothly as λ varies, showing a continuous trade-off curve between search cost and accuracy, rather than abrupt failure points. Among these components, structured reasoning contributes most significantly to the SEM's efficiency and accuracy gains by effectively reducing unnecessary external searches while maintaining robust internal reasoning.

8. THEORETICAL ANALYSIS

SEM internal backtracking can be considered a form of meta-reasoning insofar as it involves monitoring and revising the reasoning process itself; however, meta-reasoning is a broader concept encompassing various strategies to control and optimize reasoning, beyond just

backtracking. Your method can indeed be framed as a constrained Markov Decision Process (MDP), where states represent reasoning/search contexts, actions represent reasoning steps or search calls, rewards correspond to accuracy gains minus costs, and constraints represent resource or search budget limits. The reward function is designed to induce optimal stopping behavior by balancing the reward for correct answers with penalties for continued search or reasoning, effectively guiding the model to stop reasoning or search when the marginal expected benefit no longer justifies the cost. There exists a clear trade-off frontier between cost (e.g., search calls or computational resources) and accuracy, illustrating the Pareto-optimal balance achievable by adjusting parameters like the search penalty. Monotonic improvement can be proven under standard MDP and reinforcement learning assumptions, ensuring that successive policy updates improve or at least do not degrade the expected cumulative reward, as shown by classical results in policy iteration and constrained MDP theory.

9. DISCUSSION

The SEM method scales effectively to very large models, including those with 70 billion parameters and beyond, as its design leverages reinforcement learning and structured reasoning that benefit from larger model capacities and richer representations.

In open-domain real-time systems, it behaves robustly by dynamically balancing accuracy and search cost, enabling timely responses while maintaining high answer quality. SEM can optimize multi-tool systems by learning policies that decide when and which external tools or knowledge sources to invoke, improving efficiency and coordination across heterogeneous resources.

The method reduces hallucinations by incorporating internal verification steps and selective retrieval, lowering the reliance on purely generative outputs which are prone to fabricate information. By reducing unnecessary external retrievals and focusing computation on relevant reasoning, SEM contributes to lowering the carbon footprint associated with large-scale model usage. Domains benefiting most include medical AI—where accuracy and hallucination reduction are critical, autonomous systems requiring real-time safety-aware decisions, and general open-domain QA systems that must handle a diverse range of queries efficiently and reliably.

10. LIMITATIONS

SEM methods face challenges with edge cases involving highly ambiguous or underspecified queries, where the system struggles to disambiguate intent or required reasoning paths effectively. Ambiguous query classification can indeed hurt performance because misclassification may lead to inappropriate reasoning or search strategies, reducing accuracy and increasing unnecessary or failed search attempts. Noisy retrieval adversely affects SEM by introducing irrelevant or misleading information, which can degrade internal reasoning quality and cause error propagation through the reasoning chain. Reward tuning is somewhat fragile; improper tuning can bias the model toward excessive or insufficient search usage, leading to degraded accuracy or inefficient behavior, necessitating careful calibration.

There is a risk of overfitting to dataset structure, particularly if the training data exhibits strong regularities; the model might learn heuristics tied to dataset-specific patterns rather than generalizable reasoning or search policies, reducing effectiveness on out-of-distribution inputs.

11. FUTURE WORK

The SEM approach can be extended to multi-agent LLMs, where multiple agents with specialized capabilities collaborate and coordinate, enhancing complex multi-step reasoning and decision-making tasks. Search cost can indeed be made dynamic, adapting in real-time based on context, resource availability, or query complexity, enabling more flexible and efficient search strategies.

Uncertainty calibration can be incorporated to improve the model's confidence estimates, allowing the system to better decide when to rely on internal reasoning versus external search or tools.

The parameter λ , which balances search cost and accuracy, can be learned automatically via reinforcement learning or meta-optimization techniques to optimize overall system performance. SEM can integrate with model compression methods to reduce model size and computational overhead while preserving reasoning and decision-making efficacy. SEM is applicable to multimodal systems, effectively handling inputs combining text, images, and other modalities by leveraging multimodal large language models for richer contextual understanding and decision-making.

12 CONCLUSION

In conclusion, framing the Stochastic Environment Modeling (SEM) problem as a constrained Markov Decision Process provides a rigorous foundation for optimizing the trade-off between search cost and reasoning accuracy through principled policy learning. By explicitly defining the state and action spaces, transition dynamics, and reward functions, SEM enables the application of policy gradient methods to learn cost-sensitive search invocation strategies. The presented formulation balances theoretical rigor—highlighting assumptions like Lipschitz continuity and boundedness—with practical considerations, including pseudocode for implementation. While convergence guarantees are supported under certain regularity conditions, global optimality remains an open challenge, reflecting the complexity of real-world environments. This principled approach lays the groundwork for advancing SEM methods with robust optimization techniques, positioning them well for further research and deployment in dynamic, resource-aware artificial intelligence systems.

REFERENCES

- [1] Z. Sha, S. Cui, and W. Wang, “SEM: Reinforcement Learning for Search-Efficient Large Language Models,” May 12, 2025. doi: 10.48550/arxiv.2505.07903.
- [2] A. Laat, A. Wong, S. Verberne, J. Broekens, N. Van Stein, and T. Bäck, “Multi-Step Reasoning with Large Language Models, a Survey,” *ACM Comput. Surv.*, vol. 58, no. 6, pp. 1–35, Dec. 2025, doi: 10.1145/3774896.
- [3] G. Zhang et al., “The Landscape of Agentic Reinforcement Learning for LLMs: A Survey,” Sept. 02, 2025. doi: 10.48550/arxiv.2509.02547.
- [4] Y. Park et al., “Inference Optimization of Foundation Models on AI Accelerators,” *Association for Computing Machinery*, Aug. 2024, pp. 6605–6615. doi: 10.1145/3637528.3671465.
- [5] H. He et al., “Rethinking Reasoning Quality in Large Language Models through Enhanced Chain-of-Thought via RL,” Sept. 07, 2025. doi: 10.48550/arxiv.2509.06024.
- [6] K. Zhang et al., “A Survey of Reinforcement Learning for Large Reasoning Models,” Sept. 10, 2025. doi: 10.48550/arxiv.2509.08827.



-
- [7] Y. Qiao et al., “DeepSeek-Inspired Exploration of RL-Based LLMs and Synergy with Wireless Networks: A Survey,” *ACM Comput. Surv.*, vol. 58, no. 7, pp. 1–37, Dec. 2025, doi: 10.1145/3776745.
- [8] D. Guo et al., “DeepSeek-R1 incentivizes reasoning in LLMs through reinforcement learning,” *Nature*, vol. 645, no. 8081, pp. 633–638, Sept. 2025, doi: 10.1038/s41586-025-09422-z.
- [9] Y. Cao et al., “Survey on Large Language Model-Enhanced Reinforcement Learning: Concept, Taxonomy, and Methods,” *IEEE Trans. Neural Netw. Learning Syst.*, vol. 36, no. 6, pp. 9737–9757, June 2025, doi: 10.1109/tnnls.2024.3497992.
- [10] Y. Sun and F. Ji, “An Embodied Intelligence System for Coal Mine Safety Assessment Based on Multi-Level Large Language Models,” *Sensors*, vol. 25, no. 2, p. 488, Jan. 2025, doi: 10.3390/s25020488.
- [11] A. E. Sizemore, E. A. Karuza, C. Giusti, and D. S. Bassett, “Knowledge gaps in the early growth of semantic feature networks,” *Nat Hum Behav*, vol. 2, no. 9, pp. 682–692, Sept. 2018, doi: 10.1038/s41562-018-0422-4.
- [12] J. Mökander, J. Schuett, H. R. Kirk, and L. Floridi, “Auditing Large Language Models: A Three-Layered Approach,” *SSRN Journal*, Jan. 2023, doi: 10.2139/ssrn.4361607.