

AutoGraph-R1: End-to-End Reinforcement Learning for Knowledge Graph Construction

Hong Ting Tsang¹, Jiaxin Bai^{1,*}, Haoyu Huang¹, Qiao Xiao²,
Tianshi Zheng¹, Baixuan Xu¹, Shujie Liu³, Yangqiu Song¹

¹The Hong Kong University of Science and Technology

²Cornell University

³Microsoft Research

{httsangaj, jbai}@connect.ust.hk, yqsong@cse.ust.hk
shujliu@microsoft.com

Abstract

Building effective knowledge graphs (KGs) for Retrieval-Augmented Generation (RAG) is pivotal for advancing question answering (QA) systems. However, its effectiveness is hindered by a fundamental disconnect: the knowledge graph (KG) construction process is decoupled from its downstream application, yielding suboptimal graph structures. To bridge this gap, we introduce AutoGraph-R1, the first framework to directly optimize KG construction for task performance using Reinforcement Learning (RL). AutoGraph-R1 trains an LLM constructor by framing graph generation as a policy learning problem, where the reward is derived from the graph’s functional utility in a RAG pipeline. We design two novel, task-aware reward functions, one for graphs as knowledge carriers and another as knowledge indices. Across multiple QA benchmarks, **AutoGraph-R1** consistently enables graph RAG methods to achieve significant performance gains over using task-agnostic baseline graphs. Our work shows it is possible to close the loop between construction and application, shifting the paradigm from building intrinsically “good” graphs to building demonstrably “useful” ones.

1 Introduction

Retrieval-Augmented Generation (RAG) has become a cornerstone for enhancing Large Language Models (LLMs), enabling them to ground responses in external knowledge, reduce hallucinations, and incorporate real-time, domain-specific information (Gao et al., 2024a; Peng et al., 2024a). A particularly promising frontier is graph-based RAG, where LLMs leverage the structured nature of Knowledge Graphs (KGs) for complex data sensemaking and reasoning (Edge et al., 2025). The typical pipeline begins by constructing a KG from unstructured text, often using LLM-driven

extractors and heuristics (Han et al., 2024; Lairgi et al., 2024; Bai et al., 2025a), which then supports downstream question answering (Gutiérrez et al., 2025a; Sun et al., 2024).

Despite its potential, the prevailing graph-based RAG paradigm suffers from a fundamental disconnect: the process is split into two isolated phases. First, a construction phase, where a graph is built and evaluated on intrinsic metrics like precision and coverage (Huang et al., 2025a). Second, an application phase, where this static graph is used for a downstream task. The critical flaw in this approach is that a “good” graph by intrinsic standards is not necessarily a “useful” one for the end task (Xue and Zou, 2022). For instance, striving for exhaustive completeness often creates noisy, fragmented graphs where relevance signals are diluted or critical reasoning paths are broken, as illustrated in Figure 1.

This disconnect persists because closing the loop between downstream performance and graph construction is technically challenging. The construction process generating discrete (subject, predicate, object) triples is inherently non-differentiable. Consequently, standard gradient-based optimization cannot backpropagate performance signals from a downstream task, such as question answering accuracy, to guide the graph generation model. The graph, once built, cannot learn from its failures. To bridge this gap, we employ Reinforcement Learning (RL). While prior work has used RL to refine retrieval over search tools or improve query reformulation (Jin et al., 2025; Jiang et al., 2025b; Luo et al., 2025), our work is the first to leverage RL to directly optimize the KG construction process itself. We introduce AutoGraph-R1, a framework that fine-tunes an LLM-based graph generator by optimizing for downstream task performance. As shown in Figure 2, the graph generation model learns a construction policy from raw text. The utility of the generated graph is then evaluated in a

*Corresponding author.

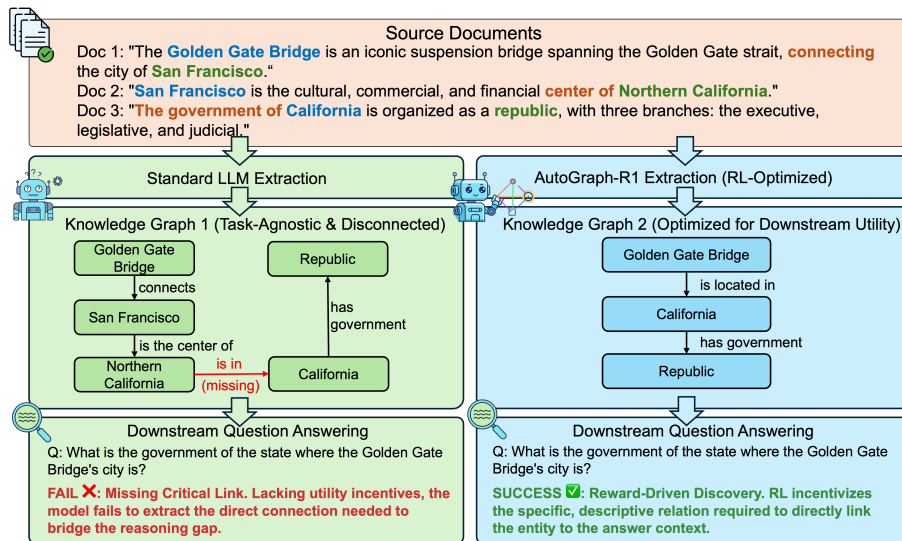


Figure 1: **Bridging the Gap Between KG Construction and Utility.** (A) Traditional extraction is task-agnostic and optimized for intrinsic metrics, often failing to capture implicit connections or diluting relevance signals with noise. (B) **AutoGraph-R1** aligns the graph structure with the specific downstream retriever. Driven by the reward, it learns to construct concrete reasoning pathways, or optimize graph connectivity to ensure robust signal propagation, delivering the specific structure required for the target task.

downstream RAG pipeline, yielding a reward signal. This task-aware reward is used to update the generator’s parameters via policy gradient methods, guiding it to produce graphs that are not just factually accurate but functionally optimal, for instance, by creating valid paths that facilitate complex reasoning.

We specifically focus on question answering over benchmark datasets requiring multi-document reasoning (Yang et al., 2018; Trivedi et al., 2022; Mallen et al., 2023). Our framework’s reward function is based on the utility of the resulting graph, evaluating whether it serves as an effective knowledge index for retrieving useful text chunks or provides subgraphs that directly support the reasoning process. By designing these task-aware rewards, we directly align the objectives of KG construction with end-task performance.

In summary, our contributions are:

- We introduce AutoGraph-R1, a novel RL framework that directly optimizes knowledge graph construction for downstream utility, bridging the critical gap between graph quality and task performance.
- We design and implement task-aware reward functions that successfully align KG structure with the demands of complex reasoning tasks, compelling the model to build functionally superior graphs.
- Through extensive experiments, we demonstrate that integrating AutoGraph-R1’s graphs into a state-of-the-art RAG pipeline yields

significant performance gains on QA benchmarks, validating that RL-driven graph construction improves downstream task utility.

2 Related Work

2.1 Graph-based Retrieval-Augmented Generation

Large Language Models (LLMs), despite demonstrating strong reasoning capabilities (DeepSeek-AI et al., 2025a), remain susceptible to factual hallucinations (Ji et al., 2023; Huang et al., 2025b) and knowledge incompleteness (Peng et al., 2023). Retrieval-Augmented Generation (RAG) (Lewis et al., 2021; Gao et al., 2024b) mitigates these issues by grounding LLMs in external knowledge sources, thereby improving factual accuracy and reasoning. A burgeoning area of research extends RAG with graph-structured knowledge (Peng et al., 2024b; Xiang et al., 2025; Zhang et al., 2025a; Han et al., 2025). In these pipelines, graphs serve two primary functions. First, as knowledge indices, where the graph organizes and connects raw text chunks, and its structural properties are leveraged for more sophisticated retrieval strategies (Liang et al., 2024; Liu et al., 2024b; Zhang et al., 2024b; Wang et al., 2023; Li et al., 2024a). Methods like HippoRAG (Gutiérrez et al., 2025a,b) exemplify this by exploiting structural connections to access relevant information more effectively (Xiao et al., 2025). Second, as knowledge carriers, where the graph itself is the primary information source, and the model reasons directly over recovered subgraphs (Shen et al., 2025b; Liu et al., 2024a). This

paradigm is adopted by approaches such as Think-on-Graph (Ma et al., 2025; Sun et al., 2024), SubgraphRAG (Li et al., 2025a), StructRAG (Li et al., 2024b), and KnowGPT (Zhang et al., 2024a).

The construction of KGs has evolved from traditional rule-based systems like OpenIE (Angeli et al., 2015) to more flexible LLM-based pipelines such as PiVE (Han et al., 2024), iText2KG (Lairgi et al., 2024), KGEN (Mo et al., 2025), GraphRAG (Edge et al., 2025), LightRAG (Guo et al., 2025) and AutoSchemaKG (Bai et al., 2025a). While powerful, these LLM-driven methods typically generate a static graph based on fixed prompts or heuristics, often evaluated using intrinsic metrics. However, the optimal structure of a KG is highly dependent on variety of downstream applications (Gubanov et al., 2024; Wu et al., 2024; Zhao et al., 2024; He et al., 2024; Liu et al., 2024c; Bai et al., 2025b; Huang et al., 2026). For instance, a graph acting as a text index may prioritize fine-grained partitioning, while one used for reasoning chains requires long-range connectivity (Jin et al., 2024; Huang et al., 2024). This creates the disconnect we identified earlier: a graph built to be “good” in isolation may be functionally poor for a specific task. Our work addresses this gap by optimizing graph construction directly for downstream performance, a problem that, to our knowledge, has not been systematically investigated, despite progress in KG refinement and completion techniques (Chen et al., 2024b,a; Zhang et al., 2022; Dong et al., 2023b).

2.2 Reinforcement Learning for Language Model Optimization

Reinforcement learning (RL) (Kaelbling et al., 1996) offers a powerful framework for optimizing the sequential decision-making capabilities of LLMs by enabling them to learn through environmental interaction and reward feedback (Kaufmann et al., 2024; Xi et al., 2025). As LLMs have become more powerful through fine-tuning, methodological advances—from Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022) to more scalable and cost-effective algorithms like Proximal Policy Optimization (PPO) (Schulman et al., 2017), Dynamic Sampling Policy Optimization (DAPO) (Yu et al., 2025), and Group Relative Policy Optimization (GRPO) (DeepSeek-AI et al., 2025a)—have enabled successful applications in diverse domains, including open-domain retrieval and scientific discovery. (Zheng et al., 2025; Yu et al., 2024; Zhu et al., 2025; Shen et al., 2025a).

Prior to its widespread adoption for LLM alignment, RL had been explored for knowledge base tasks, such as bidirectional text-to-graph conversion (Dognin et al., 2021) and abductive reasoning over knowledge graphs (Bai et al., 2024). More recently, RL has proven effective in training LLMs to interact with external tools, such as search engines (Jin et al., 2025; Jiang et al., 2025b; Li et al., 2025b; Zhang et al., 2025b; Jiang et al., 2025a). Notably, frameworks like Graph-R1 (Luo et al., 2025) have shown that RL can teach an LLM to effectively navigate graph-structured tools to improve retrieval. However, these works use RL to learn a policy for navigating or querying an existing knowledge source. Our approach is fundamentally different: we use RL to learn a policy for constructing the knowledge source from raw text. While work like (Zhang and Zhao, 2025) utilizes RL to complete or refine existing KGs during reasoning, to our knowledge, this is the first application of RL to directly optimize end-to-end graph generation from corpora based on its measured utility in a downstream task. This distinction forms the motivation for AutoGraph-R1.

3 Preliminaries

In this section, we formalize the key concepts underlying AutoGraph-R1, including knowledge graph construction, graph-based retrieval, and answer generation within a RAG pipeline.

3.1 Knowledge Graph Construction

We define a knowledge graph (KG) as a directed, labeled graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathcal{R})$, constructed from a set of documents \mathcal{D} . Here, \mathcal{V} is the set of nodes, \mathcal{R} the set of relation types, and $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{R} \times \mathcal{V}$ the set of edges represented as triples (s, r, o) . Nodes $s, o \in \mathcal{V}$ may correspond to entities, events, or concepts, and $r \in \mathcal{R}$ denotes a relation type. Following prior work (Zhang et al., 2025a), we consider two principal configurations for the graph’s role.

Graphs as Knowledge Carriers In this configuration, the graph serves as a self-contained knowledge base. It consists of factual triples $(s, r, o) \in \mathcal{E}$, where nodes $s, o \in \mathcal{V}$ are entities and $r \in \mathcal{R}$ is the relation. These triples act as discrete, structured knowledge units that are retrieved and processed directly by the downstream model.

Graphs for Knowledge Indexing Alternatively, the graph functions as a structured index over the raw document corpus \mathcal{D} . Nodes are augmented

with pointers to text spans, denoted $\tau(v)$. Formally, the node set can be partitioned into entity nodes \mathcal{V}_e and document or chunk nodes \mathcal{V}_d , such that $\mathcal{V} = \mathcal{V}_e \cup \mathcal{V}_d$. This hybrid structure allows the graph to guide retrieval not only of structured facts but also of the original, unstructured text passages.

3.2 Retrieval Module

Given a query q , the goal of the retrieval module is to produce a set of evidence units $\mathcal{C}(q)$ that will be passed to the LLM for answer generation. We consider two complementary retrieval strategies.

Graph Knowledge Retriever This retriever, denoted $\mathcal{R}_{\text{graph}}$, operates directly on the graph structure. Given a query q , it identifies relevant components such as individual triples, multi-hop paths, or entire subgraphs. Formally, we define its output as a set of structured evidence $\mathcal{P}(q) = \mathcal{R}_{\text{graph}}(q, \mathcal{G})$. The elements of $\mathcal{P}(q)$ are then linearized into text to serve as context for the LLM.

Graph-based Text Retriever This retriever, denoted $\mathcal{R}_{\text{text}}$, uses the graph as an index to find relevant text passages from the source corpus. It leverages graph connectivity to identify promising document nodes. Formally, $\mathcal{R}_{\text{text}}(q, \mathcal{G}) \mapsto \mathcal{T}(q)$, where $\mathcal{T}(q) \subseteq \{\tau(v) \mid v \in \mathcal{V}_d\}$ is a set of raw text passages linked from document nodes.

3.3 Answer Generation

The final answer generation step uses a large language model π^{ans} to synthesize an answer \hat{y} from the query q and the retrieved evidence $\mathcal{C}(q)$. The evidence context $\mathcal{C}(q)$ is composed of either the linearized graph structures $\mathcal{P}(q)$ from the graph knowledge retriever or the text passages $\mathcal{T}(q)$ from the graph-based text retriever. The final answer \hat{y} is generated by conditioning the LLM on the query and evidence: $\hat{y} = \pi^{\text{ans}}(q, \mathcal{C}(q))$. This unified framework allows our optimization process to apply to both types of graph construction, directly linking the structure of \mathcal{G} to its utility in the final QA-task.

4 AutoGraph-R1

4.1 RL for Graph Construction

AutoGraph-R1, an end-to-end reinforcement learning (RL) framework that directly optimizes knowledge graph (KG) construction with downstream task performance as the reward signal. The framework unifies two common graph-augmented

retrieval paradigms: *Graph RAG* (retrieval over entity triples) and *Graph Text RAG* (retrieval over text nodes through graph index).

As shown in Figure 2, AutoGraph-R1 consists of three components: (1) a KG construction policy model π_{θ}^{KG} , instantiated as a large language model (LLM), which maps a list of documents \mathcal{D} into a graph \mathcal{G} ; (2) a frozen RAG server with a fixed answer generator π^{ans} , which retrieves from \mathcal{G} and produces an answer \hat{y} to the input query q ; (3) a task-specific reward function $R(q, \hat{y}, y, \mathcal{G})$ that evaluates how well the constructed graph supports QA, where y is the gold answer.

Task-Aware Training Loop A central design choice in AutoGraph-R1, inspired by s3 (Jiang et al., 2025b), is to freeze the retrieval module while the KG construction policy π_{θ}^{KG} adapts. During training, each sample (q, y, \mathcal{D}_q) comprising a query, a gold answer, and relevant documents, triggers a full end-to-end loop of KG construction, retrieval, and answer generation. Crucially, the definition of a "useful" graph is contingent on the retrieval paradigm. We therefore tailor the training process for two distinct scenarios, aligning the KG’s structure with its intended function.

Training with a Graph Knowledge Retriever. When the KG acts as a knowledge carrier, retrieval quality is measured by its ability to provide a self-contained, structured context for reasoning. To isolate the impact of graph structure, we employ a simple subgraph retriever. Given a query q , we extract its named entities to serve as anchors and retrieve the n -hop neighborhood surrounding them to form the context $\mathcal{P}(q)$. This design intentionally bypasses dense vector similarity, forcing the reward signal to reflect the graph’s relational completeness and structural integrity. The policy is rewarded for creating graphs where the correct answer is directly *deducible* from the retrieved subgraph.

Training with a Graph-based Text Retriever. When the KG serves as a knowledge index, its utility is determined by how well it guides the retriever to relevant text passages. We adapt the HippoRAG-2 retriever (Gutiérrez et al., 2025b) for this purpose. First, candidate triples (s, r, o) are selected based on embedding similarity to the query q (using Qwen-3-0.6B). However, unlike the original method, we then use *only* these triple-level similarities to initialize a Personalized PageRank algorithm over the constructed graph \mathcal{G} . This process propagates relevance scores through the graph structure, ultimately identifying text nodes that are struc-

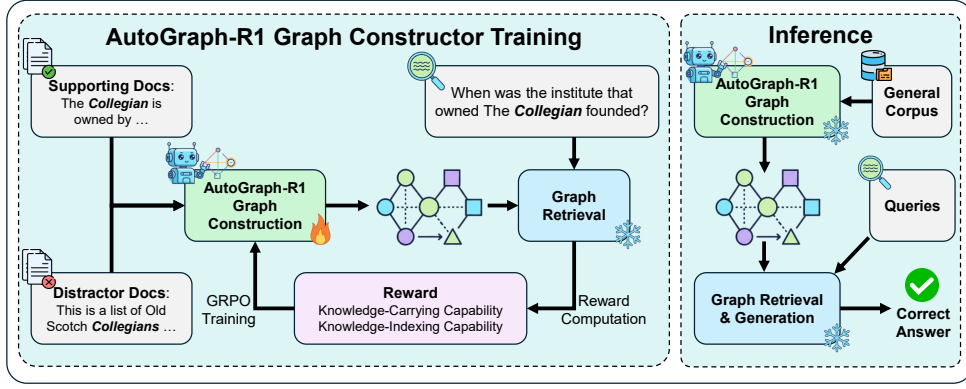


Figure 2: Overview of the AutoGraph-R1 Framework. AutoGraph-R1 optimizes knowledge graph construction for downstream utility using reinforcement learning. During the training phase (left), a graph constructor is fine-tuned with GRPO. The reward signal is derived from the performance of a graph retriever on the generated KG, directly measuring the graph’s functional quality. During the inference phase (right), the trained constructor is used to build a large-scale KG from a general corpus, which then serves a downstream graph-based RAG system.

turally connected to the most pertinent facts, from which the top- N passages are returned. The policy is therefore incentivized to build graphs where structural connectivity, not just semantic similarity, is a reliable signal for identifying crucial evidence passages, including both direct and complementary information that might otherwise be overlooked.

In both scenarios, by freezing the retriever and aligning the reward with its specific mechanism, AutoGraph-R1 ensures the graph construction policy learns to produce graphs that are optimized for a given downstream retrieval strategy.

4.2 Reward Design for Functional Graph Construction

A primary challenge in optimizing this end-to-end pipeline is the sparse and indirect nature of the learning signal, where a single reward is given after a long sequence of construction actions. This creates a severe credit assignment problem and places a heavy burden on the quality of the reward signal itself. While the final answer’s F1 score has been explored as a reward in prior work (Jin et al., 2025), we find its properties make it a challenging choice for our goal of guiding graph construction. The F1 score is brittle; minor phrasing variations in the LLM output can cause swings in the metric, an issue that persists even with deterministic decoding. This instability results in a noisy reward that can impede or destabilize policy optimization. Our approach is therefore motivated by the need to design more direct and stable, task-specific rewards better suited to our problem.

To overcome these challenges, we design two distinct reward functions that provide a more direct and stable learning signal by measuring the functional utility of the graph for a specific retrieval task.

Graph Knowledge Retriever Retrieval operates on subgraphs or relation paths. The key requirement is that the gold answer y should be *deducible* from \mathcal{G} . An answer y is considered deducible from \mathcal{G} if the retrieved triples or subgraphs contain sufficient relational information to logically infer the gold answer y for a given query q , either directly through explicit facts or indirectly through reasoning over connected triples. We therefore define a binary reward, the **Knowledge-Carrying Reward**, R_C , which measures the *deducibility* of gold answer in the constructed KG. An external LLM judge is prompted with (q, y, \mathcal{G}) and determines whether y can be deduced from the retrieved triples:

$$R_C(q, y, \mathcal{G}) = \mathbb{I}[\text{deducible}(q, y \mid \mathcal{G})] \quad (1)$$

Graph-based Text Retriever Retrieval operates over the graph structure of the KG to locate relevant text passages. To enhance the *knowledge-indexing capability*, we use **Knowledge-Indexing Reward**, R_I , as the reward function. This aligns with the fundamental objective of these retrievers by measuring the effectiveness of the retrieved passages in capturing the relevant information. The reward is defined as follows:

$$R_I(q, \mathcal{D}_{\text{gold}}, \mathcal{G}) = \frac{|\text{Top-}k(\mathcal{G}, q) \cap \mathcal{D}_{\text{gold}}|}{|\mathcal{D}_{\text{gold}}|} \quad (2)$$

where $\mathcal{D}_{\text{gold}}$ denotes the gold passages for q , and $\text{Top-}k(\mathcal{G}, q)$ are the retrieved passages.

4.3 GRPO for Graph Construction

To optimize the knowledge graph (KG) constructor policy π_{θ}^{KG} , we employ **Group-Relative Policy Optimization (GRPO)** (Shao et al., 2024), a memory-efficient alternative to Proximal Policy

Optimization (PPO). GRPO is well-suited for our LLM-based framework, as it eliminates the need for a separate value model by using a relative reward baseline derived from a group of sampled graph outputs. This approach reduces computational overhead and memory usage, enabling scalable training for large-scale graph construction tasks.

The GRPO objective updates π_{θ}^{KG} to favor graphs that maximize downstream QA performance, incorporating a clipping mechanism to ensure stable updates. We simplify the training procedure by removing the KL divergence term to lower the computational overhead and save memory usage without damaging the training (Liu et al., 2025; Hu et al., 2025). The objective is defined as:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{\substack{\mathbf{s} \sim \mathcal{D} \\ \{\mathbf{a}_i\}_i \sim \pi_{\theta, \text{old}}}} \left[\frac{1}{G} \sum_{i=1}^G \sum_{t=1}^{|\mathbf{a}_i|} \min \left(r_{i,t}(\theta) \hat{A}_i, \text{clip}(r_{i,t}(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_i \right) \right] \quad (3)$$

Here, the policy π_{θ}^{KG} takes the document list \mathbf{s} as input. In this equation, \mathbf{s} corresponds to the input document list drawn from the training distribution \mathcal{D} .

Construction of Training Inputs (\mathcal{D}): The composition of \mathbf{s} is tailored to the specific retrieval paradigm to shape the learning signal. To create a challenging training environment for the **text retrieval scenario**, we implement a hard negative mining strategy. For each query in the training set (HotpotQA and Musique), we use a Qwen3-8B (Zhang et al., 2025c) embedding model to identify the most semantically similar non-gold passage from the corpus, which is then added as a distractor. In contrast, for the **graph knowledge retriever scenario**, no distractors are used, as the primary objective is to optimize the informational completeness of the graph constructed from source documents, improving the knowledge-carrying capability of graph rather than its ability to filter irrelevant content. $\hat{A}_i = \frac{R_i - \mu_R}{\sigma_R}$ represents the Group-Relative Advantage for the entire graph \mathbf{a}_i , which is derived by normalizing its reward relative to the group’s mean μ_R and standard deviation σ_R . The downstream reward signal, R_i , for the i -th graph sample is determined by the specific training setup: $R = R_C$ (Eq. 1) when employing a graph knowl-

edge retriever, or $R = R_I$ (Eq. 2) when a graph-based text retriever is used. ϵ is a small clipping hyperparameter that ensures stable updates by preventing the new policy from straying too far from the old policy.

5 Experiments

Our experiments are designed to answer three primary research questions: **RQ1.** Does optimizing KG construction with a downstream task reward lead to better end-to-end RAG performance compared to standard, task-agnostic KG construction? **RQ2.** Is this performance improvement consistent across different graph-based RAG paradigms (i.e., graph as a *knowledge carrier* vs. a *knowledge index*) and across different model scales? **RQ3.** Does optimizing for downstream utility also improve the *intrinsic quality* (e.g., factual precision and recall) of the graph, and how do different reward functions bias the final graph structure?

5.1 Datasets and Corpora

Training Datasets For the reinforcement learning phase, we utilize two multi-hop QA datasets: HotpotQA (Yang et al., 2018) and Musique (Trivedi et al., 2022). These datasets provide the complex reasoning challenges necessary to train our policy model.

Evaluation Datasets For final RAG evaluation, we use a diverse set of five QA datasets, each comprising 1,000 samples. These include two general QA benchmarks, Natural Questions (NQ) (Kwiatkowski et al., 2019) and PopQA (Mallen et al., 2023), and three multi-hop QA benchmarks, HotpotQA (Yang et al., 2018), 2WikiMultihopQA (Ho et al., 2020), and Musique (Trivedi et al., 2022). These datasets enable a comprehensive evaluation of AutoGraph-R1 on downstream task-utility.

Corpora. For the general QA datasets (NQ, PopQA), the knowledge corpus is built from the introductory sections of the December 2021 Wikipedia dump Izacard et al. (2022). For the multi-hop QA datasets, the corpus for each is constructed from the documents associated with its 1,000 evaluation samples, following the methodology in Gutiérrez et al. (2025a).

5.2 Experiments Configs

Models We experiment with fine-tuning both Qwen2.5-3B and Qwen2.5-7B (Qwen et al., 2025) as the KG construction model (π_{θ}^{KG}). For all RAG

evaluations, we use a frozen Qwen2.5-7B as the answer generation LLM. The Qwen3-0.6B model is used consistently for all embedding tasks in both training and evaluation. **RL Training Configuration** We fine-tune the KG construction policy using the GRPO algorithm (Shao et al., 2024) on two H100 GPUs. For each training sample, the policy iteratively generates triples for each document. The training setup is tailored to the retrieval paradigm. For the graph-based text retriever, the model operates on a fixed pool of 15 documents per query, retrieving the top- N passages, where N equals the number of gold supporting passages. For the graph knowledge retriever, we retrieve a 3-hop subgraph, to ensure that a sufficient context is always retrieved for reasoning. **Evaluation Protocol** For evaluation, a KG is first constructed over the entire document corpus for each dataset. Then, depending on the type of retriever, the corresponding RAG is performed using this static graph. We report the final answer F1 score as the primary metric, consistent with prior work.

5.3 Baselines

To rigorously isolate the impact of downstream-aware optimization, we benchmark AutoGraph-R1 against a controlled, task-agnostic baseline using SOTA retrieval methods, while holding all other pipeline components constant.

KG Construction Baseline To benchmark the performance of our RL-optimized constructor, we establish a baseline using a zero-shot approach. Specifically, we construct the baseline knowledge graphs using Qwen models guided by the identical prompt used during fine-tuning. This standard, task-agnostic method allows us to isolate the impact to directly measure the gains attributable to our downstream-aware optimization, isolating the impact of the construction policy by excluding auxiliary refinement steps used in complex pipelines.

RAG Method Baselines We evaluate the KGs constructed by both AutoGraph-R1 and the zero-shot baseline using a suite of state-of-the-art RAG methods to measure their functional utility. For **graph knowledge retrieval**, where the graph itself is the source of information, we test three distinct approaches. First, we use **ToG** (Sun et al., 2024), setting both the width and depth to 3. It performs beam search on the graph and uses the Qwen3-0.6B model for relation pruning to discover meaningful paths. Second, we employ a **Subgraph Retriever**, which first performs Named Entity Recog-

inition (NER) on the query and then expands 1-hop from the identified entities to form the retrieval context. Third, we include a **Dense Triple Retriever**, which uses the Qwen3-0.6B model to retrieve triples based on the semantic similarity between their embeddings and the query embedding. For these methods, the top-10 retrieved paths or triples are used as context. For **graph-based text retrieval**, where the graph serves as an index over a text corpus, we use **HippoRAG** (Gutiérrez et al., 2025a) and **HippoRAG-2** (Gutiérrez et al., 2025b). Both methods perform query-to-edge retrieval to identify relevant triples, which then seed a Personalized PageRank (PPR) algorithm over the KG to score and rank text passages. The top-5 ranked passages are then returned as evidence.

5.4 Results and Analysis

AutoGraph-R1 consistently improves RAG performance across different paradigms, models, and scales. Our primary finding is that optimizing KG construction for downstream utility leads to significant end-to-end F1 score improvements over a standard zero-shot constructor, and this holds true for both Qwen and Llama (see Appendix A.5) model families. As shown in Table 1, when the KG acts as a *knowledge carrier*, our method yields substantial average F1 gains, up to +4.39 (Qwen-3B) and +2.97 (Qwen-7B). Similarly, when the KG is a *knowledge index*, performance increases by up to +2.09 (Qwen-3B) and +1.48 (Qwen-7B) average F1 points. This robust improvement across different models and RAG paradigms confirms that task-aware optimization is broadly effective.

AutoGraph-R1 demonstrably improves the graph’s core function as a knowledge index. The magnitude of F1 gains is more modest in the text retrieval setting. This is an expected outcome, stemming from the dual nature of using full text passages as evidence. On one hand, their rich context can enable the generator to succeed even with imperfect retrieval, masking some F1 gains. On the other hand, this verbosity can introduce noise, unlike the concise, structured triples provided by the graph knowledge retriever. To isolate the direct impact on retrieval quality, we evaluate passage recall@5. Table 1 reveals significant improvements across all models. We observe consistent gains in average recall (e.g., + over 4.0 points for Qwen-3B), confirming that our RL framework successfully optimizes the graph structure to act as a more effective index, reliably guiding text retrievers to

Table 1: Comprehensive performance evaluation of AutoGraph-R1 across Qwen models. The table is divided into two sections: (Top) Graph Knowledge Retrievers evaluated on F1 score, and (Bottom) Graph-based Text Retrievers evaluated on both F1 score and Passage Recall@5 (R@5). AutoGraph-R1 (“Ours”) consistently outperforms the zero-shot baselines across all metrics and datasets.

Methods	NQ* F1	PopQA* F1	HotpotQA F1	2WikiMultihopQA F1	Musique F1	Avg. F1						
Graph Knowledge Retrievers (F1 Only)												
<i>Qwen2.5-3B</i>												
Subgraph (Base)	26.43	54.48	39.02	34.15	13.82	33.58						
Subgraph (Ours)	28.03 ^{↑1.60}	59.46 ^{↑4.98}	40.77 ^{↑1.75}	34.71 ^{↑0.56}	15.13 ^{↑1.31}	35.62 ^{↑2.04}						
Triples (Base)	30.53	51.67	40.76	32.18	17.81	34.58						
Triples (Ours)	33.67 ^{↑3.14}	56.76 ^{↑5.09}	46.94 ^{↑6.18}	36.09 ^{↑3.91}	21.41 ^{↑3.60}	38.97 ^{↑4.39}						
ToG (Base)	26.32	54.92	41.77	43.54	18.21	36.95						
ToG (Ours)	29.27 ^{↑2.95}	61.40 ^{↑6.48}	44.56 ^{↑2.79}	49.33 ^{↑5.79}	18.42 ^{↑0.21}	40.60 ^{↑3.65}						
<i>Qwen2.5-7B</i>												
Subgraph (Base)	28.07	55.43	41.66	33.97	15.24	34.87						
Subgraph (Ours)	28.54 ^{↑0.47}	60.94 ^{↑5.51}	43.59 ^{↑1.93}	37.43 ^{↑3.46}	15.65 ^{↑0.41}	37.23 ^{↑2.36}						
Triples (Base)	33.26	55.56	44.99	35.57	20.43	37.96						
Triples (Ours)	33.98 ^{↑0.72}	58.02 ^{↑2.46}	48.28 ^{↑3.29}	36.04 ^{↑0.47}	20.56 ^{↑0.13}	39.38 ^{↑1.42}						
ToG (Base)	25.59	57.53	43.93	46.03	18.46	38.31						
ToG (Ours)	29.36 ^{↑3.77}	62.85 ^{↑5.32}	44.68 ^{↑0.75}	50.20 ^{↑4.17}	19.31 ^{↑0.85}	41.28 ^{↑2.97}						
Graph-based Text Retrievers (F1 & Recall@5)												
	F1	R@5	F1	R@5	F1	R@5	F1	R@5	F1	R@5		
<i>Qwen2.5-3B</i>												
HippoRAG (Base)	36.28	79.50	65.55	92.10	53.22	68.41	48.97	70.84	27.44	46.50	46.29	71.47
HippoRAG (Ours)	38.28 ^{↑2.00}	93.00 ^{↑13.5}	65.93 ^{↑0.38}	95.60 ^{↑3.50}	55.39 ^{↑2.17}	68.82 ^{↑0.41}	51.69 ^{↑2.72}	74.02 ^{↑3.18}	28.11 ^{↑0.67}	47.65 ^{↑1.15}	47.88 ^{↑1.59}	75.82 ^{↑4.35}
HippoRAG2 (Base)	35.88	82.20	65.02	92.20	53.70	70.06	50.98	73.49	25.70	46.93	46.25	72.98
HippoRAG2 (Ours)	38.45 ^{↑2.57}	94.00 ^{↑11.8}	66.23 ^{↑1.21}	95.40 ^{↑3.20}	56.28 ^{↑2.58}	71.21 ^{↑1.15}	52.80 ^{↑1.82}	76.42 ^{↑2.93}	27.93 ^{↑2.23}	49.13 ^{↑2.20}	48.34 ^{↑2.09}	77.23 ^{↑4.25}
<i>Qwen2.5-7B</i>												
HippoRAG (Base)	37.16	93.30	65.95	92.90	55.50	69.97	53.01	72.87	26.03	47.19	47.53	75.24
HippoRAG (Ours)	38.80 ^{↑1.64}	94.30 ^{↑1.00}	67.85 ^{↑1.90}	95.80 ^{↑2.90}	57.19 ^{↑1.69}	71.40 ^{↑1.43}	53.60 ^{↑0.59}	76.16 ^{↑3.29}	26.97 ^{↑0.94}	48.44 ^{↑1.25}	48.88 ^{↑1.35}	77.22 ^{↑1.98}
HippoRAG2 (Base)	37.02	94.10	65.74	92.80	57.08	72.03	54.99	75.98	26.77	48.55	48.32	76.69
HippoRAG2 (Ours)	38.68 ^{↑1.66}	95.00 ^{↑0.90}	67.72 ^{↑1.97}	96.30 ^{↑3.50}	58.98 ^{↑1.90}	73.61 ^{↑1.58}	56.46 ^{↑1.47}	78.66 ^{↑2.68}	27.18 ^{↑0.41}	49.23 ^{↑0.68}	49.80 ^{↑1.48}	78.56 ^{↑1.87}

Table 2: Further analysis on whether GRPO training increases the triple extraction intrinsic quality, measured by Precision, Recall and, F1 defined in previous work (Huang et al., 2025a; Bai et al., 2025a).

KG Construction Model	HotpotQA			2WikiMultihopQA			Musique			2021Wiki		
	Acc	Recall	F1	Acc	Recall	F1	Acc	Recall	F1	Acc	Recall	F1
Qwen2.5-7B-Instruct	98.50	93.68	95.65	94.80	91.19	92.68	96.77	95.27	95.73	95.03	91.39	92.92
+ GRPO with Knowledge-Carrying Reward	97.53	96.66	96.71	95.51	96.55	95.25	97.14	96.75	96.45	96.17	96.66	96.15
+ GRPO with Knowledge-Indexing Reward	98.96	94.81	96.59	98.35	94.54	96.16	99.53	93.14	95.81	97.44	95.01	95.99
Qwen2.5-3B-Instruct	94.41	91.00	91.92	83.53	79.34	81.01	92.07	89.52	90.31	87.79	86.01	86.63
+ GRPO with Knowledge-Carrying Reward	96.52	94.24	94.80	95.91	96.28	95.80	97.01	94.74	95.22	96.70	95.58	95.76
+ GRPO with Knowledge-Indexing Reward	97.11	93.15	94.64	96.19	93.66	94.48	96.20	93.87	94.55	98.22	96.04	96.85

the correct evidence.

Optimizing for downstream utility also enhances the intrinsic factual quality of the graph.

We investigated whether extrinsic optimization comes at the cost of intrinsic quality by measuring the precision, recall, and F1 score of the extracted triples against the source text (Huang et al., 2025a) using Deepseek-V3 model as a judge (DeepSeek-AI et al., 2025b). The results in Table 2 show a clear positive correlation. Across all datasets, KGs fine-tuned with AutoGraph-R1 exhibit higher intrinsic F1 scores than their zero-shot counterparts. This indicates our RL framework does not sacrifice factual accuracy for functional utility; rather, it improves both simultaneously.

The choice of reward function induces specific and beneficial structural biases in the KG. A deeper analysis of Table 2 reveals that the two reward functions specialize the graph’s structure. The **Knowledge-Carrying Reward** (R_C), optimized for graph knowledge retrieval, consistently pro-

duces graphs with higher recall, aligning with its goal of ensuring all necessary facts for reasoning are present. In contrast, the **Knowledge-Indexing Reward** (R_I), optimized for text retrieval, yields graphs with higher precision, reflecting its need for a clean, high-fidelity index. This finding highlights that AutoGraph-R1 not only improves graph quality but also tailors the graph’s structure to its specific downstream function. For analysis on comparisons with 70B+ baselines, On-Policy RL versus offline fine-tuning, and the structural analysis of finetuned LLM constructed KG, please refer to the ablation studies in Appendix A.

5.5 Ablation: Impact of Reward Function Design

To validate our task-specific rewards (R_C and R_I), we compare them against the final answer’s F1 score as a direct reward signal. We trained two additional KG constructor models using this F1 reward, keeping all other hyperparameters identi-

Table 3: Ablation Study on Reward Functions for Qwen2.5-7B. **(Top)** Graph Knowledge Retrievers evaluated on F1. **(Bottom)** Graph-based Text Retrievers evaluated on F1 and Recall@5. Our task-specific rewards (R_C , R_I) consistently outperform the naive F1 Reward and zero-shot baseline across settings.

Methods	NQ* F1		PopQA* F1		HotpotQA F1		2WikiMultihopQA F1		Musique F1		Avg. F1	
Graph Knowledge Retrievers (F1 Only)												
Subgraph (Base)	28.07		55.43		41.66		33.97		15.24		34.87	
Subgraph (F1 Reward)	27.36 _{↓0.71}		55.09 _{↓0.34}		41.15 _{↓0.51}		35.13		15.70		34.89	
Subgraph (R_C)	28.54		60.94		43.59		37.43		15.65		37.23	
Triples (Base)	33.26		55.56		44.99		35.57		20.43		37.96	
Triples (F1 Reward)	31.52 _{↓1.74}		53.85 _{↓1.71}		44.52 _{↓0.47}		30.68 _{↓4.89}		18.00 _{↓2.43}		35.71 _{↓2.25}	
Triples (R_C)	33.98		58.02		48.28		36.04		20.56		39.38	
ToG (Base)	25.59		57.53		43.93		46.03		18.46		38.31	
ToG (F1 Reward)	27.64		56.95 _{↓0.58}		45.19		51.10		18.37 _{↓0.09}		39.85	
ToG (R_C)	29.36		62.85		44.68		50.20		19.31		41.28	
Graph-based Text Retrievers (F1 & Recall@5)												
	F1 R@5		F1 R@5		F1 R@5		F1 R@5		F1 R@5		F1 R@5	
HippoRAG (Base)	37.16	93.30	65.95	92.90	55.50	69.97	53.01	72.87	26.03	47.19	47.53	75.25
HippoRAG (F1 Reward)	39.66	93.90	63.74 _{↓2.21}	92.10 _{↓0.80}	53.74 _{↓1.76}	69.18 _{↓0.79}	49.58 _{↓3.43}	72.71 _{↓0.16}	28.68	47.71	47.08 _{↓0.45}	75.12 _{↓0.13}
HippoRAG (R_I)	38.80	94.30	67.85	95.80	57.19	71.40	53.60	76.16	26.97	48.44	48.88	77.22
HippoRAG2 (Base)	37.02	94.10	65.74	92.80	57.08	72.03	54.99	75.98	26.77	48.55	48.32	76.69
HippoRAG2 (F1 Reward)	38.33	94.80	62.92 _{↓2.82}	92.00 _{↓0.80}	55.19 _{↓1.89}	71.07 _{↓0.96}	50.23 _{↓4.76}	72.81 _{↓3.17}	27.51	48.75	46.83 _{↓1.49}	75.88 _{↓0.81}
HippoRAG2 (R_I)	38.68	95.00	67.72	96.30	58.98	73.61	56.46	78.66	27.18	49.23	49.80	78.56

cal to our main setup. As shown in Table 3, the F1-rewarded model underperforms the zero-shot baseline by over -2.2 average F1 on the Triples Retriever and actively degrades both HippoRAG and HippoRAG2 performance. In contrast, our task-specific rewards deliver consistent gains across all settings, confirming that measuring functional graph utility directly is essential for stable RL-based KG construction. Detailed training dynamics are provided in Appendix B.

5.6 Computational Cost

The per-sample training cost is bounded: graph construction operates over a small, fixed document pool (15 documents for text retrieval; 2–5 for graph retrieval), and the R_C judge prompt (Figure 11) is a minimal binary Yes/No query that can be batched efficiently. Table 4 reports wall-clock times on $2 \times \text{H100}$ GPUs, which are well within academic lab budgets.

Table 4: Wall-clock training times for AutoGraph-R1 on $2 \times \text{H100}$ GPUs.

Experiment Setting	Time
7B Graph-based Text Retriever	~8h 38m
3B Graph-based Text Retriever	~5h 11m
7B Graph Knowledge Retriever	~3h 15m
3B Graph Knowledge Retriever	~2h 18m

6 Conclusion

In this work, we introduced **AutoGraph-R1**, the first reinforcement learning framework for knowledge graph construction that directly optimizes downstream RAG performance. By incorporating task-aware rewards, our approach bridges the gap

between traditional graph quality metrics and end-task utility. Experiments across five QA benchmarks demonstrate consistent improvements over strong baselines in both graph knowledge and graph-based text retrieval. Overall, our work shows that reinforcement learning can effectively connect the graph construction process with downstream QA performance, ensuring that knowledge graphs are optimized for their intended applications.

7 Limitations

While AutoGraph-R1 advances the state of task-aware knowledge graph construction, several limitations warrant discussion:

Dependency on Frozen Retrievers. Our framework optimizes the graph construction policy to maximize utility for a *specific, frozen* retrieval module. While this ensures high performance for the target pipeline, it implicitly couples the graph structure to the retriever’s inductive biases. A graph optimized for a graph knowledge retriever effectively learns to construct reasoning pathways, which may not be the optimal topology for a distinct retrieval paradigm (e.g., a BM25 retriever).

Reward Computation Bottleneck. Our use of task-aware rewards (R_C and R_I) necessitates executing a full retrieval and evaluation loop during training. For the Knowledge-Carrying reward specifically, this involves querying a LLM judge for every generated sample to verify deducibility. This introduces computational overhead compared to offline or heuristic-based training objectives, potentially limiting the scalability of our method to extremely large training datasets without substantial compute resources.

Model Scale and Generalization. Due to computational constraints, our experiments were primarily conducted on LLMs in the 1B to 7B parameter range (Qwen2.5 and Llama-3.2). While we demonstrate that fine-tuned smaller models can outperform larger zero-shot baselines, we have not yet empirically verified the scaling laws of our RL framework on massive foundation models (e.g., 70B+ parameters). Additionally, our evaluation is limited to English-language benchmarks; the effectiveness of our reward-driven policy in multilingual or low-resource settings remains an open question.

Dependency on Associated Passages. Our training pipeline assumes access to associated passages alongside query-answer pairs, i.e., training data of the form (*query*, *gold_answer*, *associated_passages*). While such data is available in standard multi-hop QA benchmarks (e.g., HotpotQA, Musique), this assumption may not hold in more general settings where only (*query*, *gold_answer*) pairs are available. Collecting associated passages constitutes an additional annotation cost, which may limit the applicability of our framework in low-resource or domain-specific scenarios. We leave the exploration of passage-free training strategies as an important direction for future work.

8 Ethics Statement

This work complies with the ACL Code of Ethics. Our research focuses on optimizing the structural utility of knowledge graphs and does not involve human subjects or the collection of non-public data. All experiments utilize established, publicly available datasets (HotpotQA, Musique, NQ, PopQA). We acknowledge that generative models and retrieval-augmented systems can inadvertently propagate biases present in their training corpora (e.g., Wikipedia) or pre-trained base models. Furthermore, while our method aims to improve retrieval accuracy, the automated construction of knowledge graphs carries a risk of generating or amplifying hallucinations if the source text is factually incorrect. We encourage practitioners to audit both the source corpora and the constructed graphs, particularly when applying this framework in sensitive domains such as healthcare or finance. This paper uses LLMs to assist with grammar and writing quality.

Acknowledgments

The authors of this paper were supported by the ITSP Platform Research Project (ITS/189/23FP) from ITC of Hong Kong, SAR, China, and the AoE (AoE/E-601/24-N), the RIF (R6021-20), the GRF (16205322) and the JRFS (JRFS2526-6S10) from RGC of Hong Kong, SAR, China. We also thank the support from the Microsoft Accelerate Foundation Models Research (AFMR) grant program.

References

- Gabor Angeli, Melvin Jose Johnson Premkumar, and Christopher D. Manning. 2015. [Leveraging Linguistic Structure For Open Domain Information Extraction](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 344–354, Beijing, China. Association for Computational Linguistics.
- Jiaxin Bai, Wei Fan, Qi Hu, Qing Zong, Chunyang Li, Hong Ting Tsang, Hongyu Luo, Yauwai Yim, Haoyu Huang, Xiao Zhou, Feng Qin, Tianshi Zheng, Xi Peng, Xin Yao, Huiwen Yang, Leijie Wu, Yi Ji, Gong Zhang, Renhai Chen, and Yangqiu Song. 2025a. [AutoSchemaKG: Autonomous Knowledge Graph Construction through Dynamic Schema Induction from Web-Scale Corpora](#). *arXiv preprint*. ArXiv:2505.23628 [cs].
- Jiaxin Bai, Yicheng Wang, Tianshi Zheng, Yue Guo, Xin Liu, and Yangqiu Song. 2024. [Advancing Abductive Reasoning in Knowledge Graphs through Complex Logical Hypothesis Generation](#). *arXiv preprint*. ArXiv:2312.15643 [cs].
- Jiaxin Bai, Zihao Wang, Yukun Zhou, Hang Yin, Weizhi Fei, Qi Hu, Zheyang Deng, Jiayang Cheng, Tianshi Zheng, Hong Ting Tsang, Yisen Gao, Zhongwei Xie, Yufei Li, Lixin Fan, Binhang Yuan, Wei Wang, Lei Chen, Xiaofang Zhou, and Yangqiu Song. 2025b. [Top Ten Challenges Towards Agentic Neural Graph Databases](#). *arXiv preprint*. ArXiv:2501.14224 [cs].
- Shengyuan Chen, Yunfeng Cai, Huang Fang, Xiao Huang, and Mingming Sun. Differentiable Neuro-Symbolic Reasoning on Large-Scale Knowledge Graphs.
- Shengyuan Chen, Qinggang Zhang, Junnan Dong, Wen Hua, Jiannong Cao, and Xiao Huang. 2024a. [Neuro-Symbolic Entity Alignment via Variational Inference](#). *arXiv preprint*. ArXiv:2410.04153 [cs].
- Shengyuan Chen, Qinggang Zhang, Junnan Dong, Wen Hua, Qing Li, and Xiao Huang. 2024b. [Entity Alignment with Noisy Annotations from Large Language Models](#). *arXiv preprint*. ArXiv:2405.16806 [cs].

- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025a. **DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning**. *arXiv preprint*. ArXiv:2501.12948 [cs].
- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 181 others. 2025b. **DeepSeek-V3 Technical Report**. *arXiv preprint*. ArXiv:2412.19437 [cs].
- Pierre Dognin, Inkit Padhi, Igor Melnyk, and Payel Das. 2021. **ReGen: Reinforcement Learning for Text and Knowledge Base Generation using Pretrained Language Models**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1084–1099, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Hanze Dong, Wei Xiong, Deepanshu Goyal, Yihan Zhang, Winnie Chow, Rui Pan, Shizhe Diao, Jipeng Zhang, Kashun Shum, and Tong Zhang. 2023a. **RAFT: Reward rAnked FineTuning for Generative Foundation Model Alignment**. *arXiv preprint*. ArXiv:2304.06767 [cs].
- Junnan Dong, Qinggang Zhang, Xiao Huang, Qiaoyu Tan, Daochen Zha, and Zhao Zihao. 2023b. **Active Ensemble Learning for Knowledge Graph Error Detection**. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, pages 877–885, Singapore Singapore. ACM.
- Darren Edge, Ha Trinh, Newman Cheng, Joshua Bradley, Alex Chao, Apurva Mody, Steven Truitt, Dasha Metropolitan, Robert Osazuwa Ness, and Jonathan Larson. 2025. **From Local to Global: A Graph RAG Approach to Query-Focused Summarization**. *arXiv preprint*. ArXiv:2404.16130 [cs].
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024a. **Retrieval-augmented generation for large language models: A survey**. *Preprint*, arXiv:2312.10997.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. 2024b. **Retrieval-Augmented Generation for Large Language Models: A Survey**. *arXiv preprint*. ArXiv:2312.10997 [cs].
- Michael Gubanov, Anna Pyayt, and Aleksandra Karolak. 2024. **CancerKG.ORG A Web-scale, Interactive, Verifiable Knowledge Graph-LLM Hybrid for Assisting with Optimal Cancer Treatment and Care**. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 4497–4505. ArXiv:2501.00223 [cs].
- Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. 2025. **LightRAG: Simple and Fast Retrieval-Augmented Generation**. *arXiv preprint*. ArXiv:2410.05779 [cs].
- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2025a. **HippoRAG: Neurobiologically Inspired Long-Term Memory for Large Language Models**. *arXiv preprint*. ArXiv:2405.14831 [cs].
- Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. 2025b. **From RAG to Memory: Non-Parametric Continual Learning for Large Language Models**. *arXiv preprint*. ArXiv:2502.14802 [cs].
- Haoyu Han, Harry Shomer, Yu Wang, Yongjia Lei, Kai Guo, Zhigang Hua, Bo Long, Hui Liu, and Jiliang Tang. 2025. **RAG vs. GraphRAG: A Systematic Evaluation and Key Insights**. *arXiv preprint*. ArXiv:2502.11371 [cs].
- Jiuzhou Han, Nigel Collier, Wray Buntine, and Ehsan Shareghi. 2024. **PiVe: Prompting with Iterative Verification Improving Graph-based Generative Capability of LLMs**. *arXiv preprint*. ArXiv:2305.12392 [cs].
- Mengliang He, Aimin Zhou, and Xiaoming Shi. **Enhancing Textbook Question Answering with Knowledge Graph-Augmented Large Language Models**.
- Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. **G-Retriever: Retrieval-Augmented Generation for Textual Graph Understanding and Question Answering**. *arXiv preprint*. ArXiv:2402.07630 [cs].
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. **Constructing A Multi-hop QA Dataset for Comprehensive Evaluation of Reasoning Steps**. *arXiv preprint*. ArXiv:2011.01060 [cs].
- Jingcheng Hu, Yinmin Zhang, Qi Han, Daxin Jiang, Xianguyu Zhang, and Heung-Yeung Shum. 2025. **OpenReasoner-Zero: An Open Source Approach to Scaling Up Reinforcement Learning on the Base Model**. *arXiv preprint*. ArXiv:2503.24290 [cs].
- Haoyu Huang, Chong Chen, Zeang Sheng, Yang Li, and Wentao Zhang. 2025a. **Can LLMs be Good Graph Judge for Knowledge Graph Construction?** *arXiv preprint*. ArXiv:2411.17388 [cs].
- Haoyu Huang, Hong Ting Tsang, Jiabin Bai, Xi Peng, Gong Zhang, and Yangqiu Song. 2026. **AtlasKV: Augmenting LLMs with Billion-Scale Knowledge Graphs in 20GB VRAM**. *arXiv preprint*. ArXiv:2510.17934 [cs].
- Jiatan Huang, Mingchen Li, Zonghai Yao, Zhichao Yang, Yongkang Xiao, Feiyun Ouyang, Xiaohan Li, Shuo Han, and Hong Yu. 2024. **RiTeK: A Dataset**

- for Large Language Models Complex Reasoning over Textual Knowledge Graphs. *arXiv preprint*. ArXiv:2410.13987 [cs].
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025b. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Transactions on Information Systems*, 43(2):1–55. ArXiv:2311.05232 [cs].
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2022. Atlas: Few-shot Learning with Retrieval Augmented Language Models. *arXiv preprint*. ArXiv:2208.03299 [cs].
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Delong Chen, Wenliang Dai, Ho Shu Chan, Andrea Madotto, and Pascale Fung. 2023. Survey of Hallucination in Natural Language Generation. *ACM Computing Surveys*, 55(12):1–38. ArXiv:2202.03629 [cs].
- Pengcheng Jiang, Jiacheng Lin, Lang Cao, Runchu Tian, SeongKu Kang, Zifeng Wang, Jimeng Sun, and Jiawei Han. 2025a. DeepRetrieval: Hacking Real Search Engines and Retrievers with Large Language Models via Reinforcement Learning. *arXiv preprint*. ArXiv:2503.00223 [cs].
- Pengcheng Jiang, Xueqiang Xu, Jiacheng Lin, Jinfeng Xiao, Zifeng Wang, Jimeng Sun, and Jiawei Han. 2025b. s3: You Don't Need That Much Data to Train a Search Agent via RL. *arXiv preprint*. ArXiv:2505.14146 [cs].
- Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng Tang, Suhang Wang, Yu Meng, and Jiawei Han. 2024. Graph Chain-of-Thought: Augmenting Large Language Models by Reasoning on Graphs. *arXiv preprint*. ArXiv:2404.07103 [cs].
- Bowen Jin, Hansi Zeng, Zhenrui Yue, Jinsung Yoon, Sercan Arik, Dong Wang, Hamed Zamani, and Jiawei Han. 2025. Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning. *arXiv preprint*. ArXiv:2503.09516 [cs].
- L. P. Kaelbling, M. L. Littman, and A. W. Moore. 1996. Reinforcement Learning: A Survey. *arXiv preprint*. ArXiv:cs/9605103.
- Timo Kaufmann, Paul Weng, Viktor Bengs, and Eyke Hüllermeier. 2024. A Survey of Reinforcement Learning from Human Feedback. *arXiv preprint*. ArXiv:2312.14925 [cs].
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics*, 7:452–466. Place: Cambridge, MA Publisher: MIT Press.
- Yassir Lairgi, Ludovic Moncla, Rémy Cazabet, Khalid Benabdeslem, and Pierre Cléau. 2024. iText2KG: Incremental Knowledge Graphs Construction Using Large Language Models. *arXiv preprint*. ArXiv:2409.03284 [cs].
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *arXiv preprint*. ArXiv:2005.11401 [cs].
- Mufei Li, Siqi Miao, and Pan Li. 2025a. Simple Is Effective: The Roles of Graphs and Large Language Models in Knowledge-Graph-Based Retrieval-Augmented Generation. *arXiv preprint*. ArXiv:2410.20724 [cs].
- Shilong Li, Yancheng He, Hangyu Guo, Xingyuan Bu, Ge Bai, Jie Liu, Jiaheng Liu, Xingwei Qu, Yangguang Li, Wanli Ouyang, Wenbo Su, and Bo Zheng. 2024a. GraphReader: Building Graph-based Agent to Enhance Long-Context Abilities of Large Language Models. *arXiv preprint*. ArXiv:2406.14550 [cs].
- Yangning Li, Weizhi Zhang, Yuyao Yang, Wei-Chieh Huang, Yaozu Wu, Junyu Luo, Yuanchen Bei, Henry Peng Zou, Xiao Luo, Yusheng Zhao, Chunkit Chan, Yankai Chen, Zhongfen Deng, Yinghui Li, Haitao Zheng, Dongyuan Li, Renhe Jiang, Ming Zhang, Yangqiu Song, and Philip S. Yu. 2025b. Towards Agentic RAG with Deep Reasoning: A Survey of RAG-Reasoning Systems in LLMs. *arXiv preprint*. ArXiv:2507.09477 [cs].
- Zhuoqun Li, Xuanang Chen, Haiyang Yu, Hongyu Lin, Yaojie Lu, Qiaoyu Tang, Fei Huang, Xianpei Han, Le Sun, and Yongbin Li. 2024b. StructRAG: Boosting Knowledge Intensive Reasoning of LLMs via Inference-time Hybrid Information Structurization. *arXiv preprint*. ArXiv:2410.08815 [cs].
- Xun Liang, Simin Niu, Zhiyu li, Sensen Zhang, Shichao Song, Hanyu Wang, Jiawei Yang, Feiyu Xiong, Bo Tang, and Chenyang Xi. 2024. Empowering Large Language Models to Set up a Knowledge Retrieval Indexer via Self-Learning. *arXiv preprint*. ArXiv:2405.16933 [cs].
- Haochen Liu, Song Wang, Yaochen Zhu, Yushun Dong, and Jundong Li. 2024a. Knowledge Graph-Enhanced Large Language Models via Path Selection. *arXiv preprint*. ArXiv:2406.13862 [cs].
- Wei Liu, Ailun Yu, Daoguang Zan, Bo Shen, Wei Zhang, Haiyan Zhao, Zhi Jin, and Qianxiang Wang. 2024b. GraphCoder: Enhancing Repository-Level

- Code Completion via Code Context Graph-based Retrieval and Language Model. *arXiv preprint*. ArXiv:2406.07003 [cs].
- Xiangyan Liu, Bo Lan, Zhiyuan Hu, Yang Liu, Zhicheng Zhang, Fei Wang, Michael Shieh, and Wenmeng Zhou. 2024c. **CodexGraph: Bridging Large Language Models and Code Repositories via Code Graph Databases**. *arXiv preprint*. ArXiv:2408.03910 [cs].
- Zichen Liu, Changyu Chen, Wenjun Li, Penghui Qi, Tianyu Pang, Chao Du, Wee Sun Lee, and Min Lin. 2025. **Understanding R1-Zero-Like Training: A Critical Perspective**. *arXiv preprint*. ArXiv:2503.20783 [cs].
- Haoran Luo, Haihong E, Guanting Chen, Qika Lin, Yikai Guo, Fangzhi Xu, Zemin Kuang, Meina Song, Xiaobao Wu, Yifan Zhu, and Luu Anh Tuan. 2025. **Graph-R1: Towards Agentic GraphRAG Framework via End-to-end Reinforcement Learning**. *arXiv preprint*. ArXiv:2507.21892 [cs].
- Shengjie Ma, Chengjin Xu, Xuhui Jiang, Muzhi Li, Huaren Qu, Cehao Yang, Jiaxin Mao, and Jian Guo. 2025. **Think-on-Graph 2.0: Deep and Faithful Large Language Model Reasoning with Knowledge-guided Retrieval Augmented Generation**. *arXiv preprint*. ArXiv:2407.10805 [cs].
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. **When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories**. *arXiv preprint*. ArXiv:2212.10511 [cs].
- Belinda Mo, Kyssen Yu, Joshua Kazdan, Proud Mpala, Lisa Yu, Chris Cundy, Charilaos Kanatsoulis, and Sanmi Koyejo. 2025. **KGGen: Extracting Knowledge Graphs from Plain Text with Language Models**. *arXiv preprint*. ArXiv:2502.09956 [cs].
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. **Training language models to follow instructions with human feedback**. *arXiv preprint*. ArXiv:2203.02155 [cs].
- Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. 2024a. **Graph retrieval-augmented generation: A survey**. *Preprint*, arXiv:2408.08921.
- Boci Peng, Yun Zhu, Yongchao Liu, Xiaohe Bo, Haizhou Shi, Chuntao Hong, Yan Zhang, and Siliang Tang. 2024b. **Graph Retrieval-Augmented Generation: A Survey**. *arXiv preprint*. ArXiv:2408.08921 [cs].
- Cheng Peng, Xi Yang, Aokun Chen, Kaleb E. Smith, Nima PourNejatian, Anthony B. Costa, Cheryl Martin, Mona G. Flores, Ying Zhang, Tanja Magoc, Gloria Lipori, Duane A. Mitchell, Naykky S. Ospina, Mustafa M. Ahmed, William R. Hogan, Elizabeth A. Shenkman, Yi Guo, Jiang Bian, and Yonghui Wu. 2023. **A Study of Generative Large Language Model for Medical Research and Healthcare**. *npj Digital Medicine*, 6(1):210. ArXiv:2305.13523 [cs].
- Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jingren Zhou, Junyang Lin, and 24 others. 2025. **Qwen2.5 Technical Report**. *arXiv preprint*. ArXiv:2412.15115 [cs].
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. **Proximal Policy Optimization Algorithms**. *arXiv preprint*. ArXiv:1707.06347 [cs].
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. **DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models**. *arXiv preprint*. ArXiv:2402.03300 [cs].
- Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, Ruochen Xu, and Tiancheng Zhao. 2025a. **VLM-R1: A Stable and Generalizable R1-style Large Vision-Language Model**. *arXiv preprint*. ArXiv:2504.07615 [cs].
- Xiangqing Shen, Fanfan Wang, and Rui Xia. 2025b. **Reason-Align-Respond: Aligning LLM Reasoning with Knowledge Graphs for KGQA**. *arXiv preprint*. ArXiv:2505.20971 [cs].
- Jiashuo Sun, Chengjin Xu, Luminyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel M. Ni, Heung-Yeung Shum, and Jian Guo. 2024. **Think-on-Graph: Deep and Responsible Reasoning of Large Language Model on Knowledge Graph**. *arXiv preprint*. ArXiv:2307.07697 [cs].
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2022. **MuSiQue: Multi-hop Questions via Single-hop Question Composition**. *arXiv preprint*. ArXiv:2108.00573 [cs].
- Yu Wang, Nedim Lipka, Ryan A. Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. 2023. **Knowledge Graph Prompting for Multi-Document Question Answering**. *arXiv preprint*. ArXiv:2308.11730 [cs] version: 3.
- Junde Wu, Jiayuan Zhu, Yunli Qi, Jingkun Chen, Min Xu, Filippo Menolascina, and Vicente Grau. 2024. **Medical Graph RAG: Towards Safe Medical Large Language Model via Graph Retrieval-Augmented Generation**. *arXiv preprint*. ArXiv:2408.04187 [cs].

- Yunjia Xi, Jianghao Lin, Yongzhao Xiao, Zheli Zhou, Rong Shan, Te Gao, Jiachen Zhu, Weiwen Liu, Yong Yu, and Weinan Zhang. 2025. [A Survey of LLM-based Deep Search Agents: Paradigm, Optimization, Evaluation, and Challenges](#). *arXiv preprint*. ArXiv:2508.05668 [cs].
- Zhishang Xiang, Chuanjie Wu, Qinggang Zhang, Shengyuan Chen, Zijin Hong, Xiao Huang, and Jinsong Su. 2025. [When to use Graphs in RAG: A Comprehensive Analysis for Graph Retrieval-Augmented Generation](#). *arXiv preprint*. ArXiv:2506.05690 [cs].
- Qiao Xiao, Hong Ting Tsang, and Jiaxin Bai. 2025. [TERAG: Token-Efficient Graph-Based Retrieval-Augmented Generation](#). *arXiv preprint*. ArXiv:2509.18667 [cs].
- Bingcong Xue and Lei Zou. 2022. [Knowledge Graph Quality Management: a Comprehensive Survey](#). *IEEE Transactions on Knowledge and Data Engineering*, pages 1–1.
- 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 [cs].
- Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, and 16 others. 2025. [DAPO: An Open-Source LLM Reinforcement Learning System at Scale](#). *arXiv preprint*. ArXiv:2503.14476 [cs].
- Zishun Yu, Yunzhe Tao, Liyu Chen, Tao Sun, and Hongxia Yang. 2024. [\$\mathcal{B}\$ -Coder: Value-Based Deep Reinforcement Learning for Program Synthesis](#). *arXiv preprint*. ArXiv:2310.03173 [cs].
- Qinggang Zhang, Shengyuan Chen, Yuanchen Bei, Zheng Yuan, Huachi Zhou, Zijin Hong, Junnan Dong, Hao Chen, Yi Chang, and Xiao Huang. 2025a. [A Survey of Graph Retrieval-Augmented Generation for Customized Large Language Models](#). *arXiv preprint*. ArXiv:2501.13958 [cs] version: 1.
- Qinggang Zhang, Junnan Dong, Hao Chen, Daochen Zha, Zailiang Yu, and Xiao Huang. 2024a. [KnowGPT: Knowledge Graph based Prompting for Large Language Models](#). *arXiv preprint*. ArXiv:2312.06185 [cs].
- Qinggang Zhang, Junnan Dong, Keyu Duan, Xiao Huang, Yezi Liu, and Linchuan Xu. 2022. [Contrastive Knowledge Graph Error Detection](#). In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 2590–2599. ArXiv:2211.10030 [cs].
- Weizhi Zhang, Yangning Li, Yuanchen Bei, Junyu Luo, Guancheng Wan, Liangwei Yang, Chenxuan Xie, Yuyao Yang, Wei-Chieh Huang, Chunyu Miao, Henry Peng Zou, Xiao Luo, Yusheng Zhao, Yankai Chen, Chunkit Chan, Peilin Zhou, Xinyang Zhang, Chenwei Zhang, Jingbo Shang, and 4 others. 2025b. [From Web Search towards Agentic Deep Research: Incentivizing Search with Reasoning Agents](#). *arXiv preprint*. ArXiv:2506.18959 [cs].
- Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. 2025c. [Qwen3 Embedding: Advancing Text Embedding and Reranking Through Foundation Models](#). *arXiv preprint*. ArXiv:2506.05176 [cs].
- Yichi Zhang, Binbin Hu, Zhuo Chen, Lingbing Guo, Ziqi Liu, Zhiqiang Zhang, Lei Liang, Huajun Chen, and Wen Zhang. 2024b. [Multi-domain Knowledge Graph Collaborative Pre-training and Prompt Tuning for Diverse Downstream Tasks](#). *arXiv preprint*. ArXiv:2405.13085 [cs] version: 1.
- Zhiqiang Zhang and Wen Zhao. 2025. [A Collaborative Reasoning Framework Powered by Reinforcement Learning and Large Language Models for Complex Questions Answering over Knowledge Graph](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 10672–10684, Abu Dhabi, UAE. Association for Computational Linguistics.
- Xinjie Zhao, Moritz Blum, Rui Yang, Boming Yang, Luis Márquez Carpio, Mónica Pina-Navarro, Tony Wang, Xin Li, Huitao Li, Yanran Fu, Rongrong Wang, Juntao Zhang, and Irene Li. 2024. [AGENTi-Graph: An Interactive Knowledge Graph Platform for LLM-based Chatbots Utilizing Private Data](#). *arXiv preprint*. ArXiv:2410.11531 [cs].
- Tianshi Zheng, Zheyang Deng, Hong Ting Tsang, Weiqi Wang, Jiabin Bai, Zihao Wang, and Yangqiu Song. 2025. [From Automation to Autonomy: A Survey on Large Language Models in Scientific Discovery](#). *arXiv preprint*. ArXiv:2505.13259 [cs].
- Changtai Zhu, Siyin Wang, Ruijun Feng, Kai Song, and Xipeng Qiu. 2025. [ConvSearch-R1: Enhancing Query Reformulation for Conversational Search with Reasoning via Reinforcement Learning](#). *arXiv preprint*. ArXiv:2505.15776 [cs].

A Ablation Study

A.1 Comparison with State-of-the-Art Large Models

To assess the efficacy of our reward-driven optimization, we compared our fine-tuned 3B and 7B models against a suite of state-of-the-art (SOTA) large-scale models, including Llama-3.3-70B, Qwen2.5-72B, and GPT-5-mini and etc.

Table 5 present the results. Our findings demonstrate that AutoGraph-R1 allows smaller models to achieve parity with, and often surpass, models that are significantly larger. Specifically:

- On Graph Retrievers, our Qwen-7B model demonstrates remarkable efficiency, surpassing most zero-shot models on both Subgraph and Triples retriever, while maintaining competitive parity on the ToG retriever.
- On Text Retrievers, our Qwen-7B model outperforms the baselines for HippoRAG and HippoRAG2, validating that our method builds more effective knowledge indices than general-purpose large models.

A.2 Ablation Study on Training Methodology (RL vs. RAFT)

To determine if our performance gains stem specifically from On-Policy RL (GRPO) or simply from reward-guided data selection, we conducted an ablation study comparing GRPO against a RAFT (Reward rAnked FineTuning) approach (Dong et al., 2023a). In the RAFT setup, we utilized our reward functions to filter generated KGs offline and fine-tuned the model on high-reward samples.

Table 6 presents the results on the Qwen2.5-7B model. Both strategies significantly outperform the zero-shot baseline, confirming that our reward functions capture genuine signal for downstream utility. Comparing the two, GRPO consistently yields superior results, achieving the highest average F1 score in 4 out of 5 retrieval settings. While RAFT remains a highly competitive and computationally efficient alternative, even slightly outperforming GRPO on the ToG retriever the results suggest that the active exploration inherent in on-policy RL enables the model to discover more robust structural strategies than offline filtering alone.

A.3 Impact of Judge Model Capacity

We investigated the sensitivity of our training framework to the capacity of the Judge model used

to compute rewards. Table 7 compares the performance when training with a 3B-parameter Judge versus a 7B-parameter Judge. The results indicate that a larger Judge generally provides a higher-fidelity signal, leading to better final RAG performance across the board.

A.4 Integration with Existing KG-RAG Pipelines

To verify the transferability of our method, we integrated our fine-tuned Qwen-7B model into the official HippoRAG (Gutiérrez et al., 2025a) codebase, replacing their default zero-shot construction model. As shown in Table 8, our model functions as a plug-and-play enhancement, improving both F1 and Recall metrics across datasets.

A.5 Additional Results on Llama Models

In this section, we present the performance improvements of AutoGraph-R1 when applied to the Llama-3.2 model family. These results (Tables 9, 10, and 11) demonstrate that our method generalizes effectively across different LLM architectures.

A.6 Structural Analysis of Constructed KGs

To understand how our reward functions alter the topology of the generated graphs, we analyzed the structural metrics of KGs generated by the Qwen model family across four datasets: 2021Wiki, 2WikiMultihopQA, HotpotQA, and Musique. Table 12 presents the granular statistics for both 3B and 7B models across all evaluated datasets. The analysis reveals two distinct, reward-driven strategies:

- Optimizing for Completeness (Knowledge-Carrying Reward): The model trained for graph retrieval (Knowledge-Carrying) maintains a high number of triples but generates a dramatically larger vocabulary of unique relation types (e.g., +34% on average for Qwen-7B). This is direct evidence that the model learns to be more descriptive, discovering nuanced, long-tail relations instead of generic ones. This creates a richer graph structure necessary for multi-hop reasoning.
- Optimizing for Precision (Knowledge-Indexing Reward): Conversely, the model trained for the text retriever (Knowledge-Indexing) learns to be more concise. It produces fewer total triples per document (e.g., -12% on average for Qwen-7B) and a

Table 5: Comprehensive performance evaluation of various large-scale models as KG constructors. The table compares zero-shot KG construction performance across Graph Knowledge Retrievers (F1 Only) and Graph-based Text Retrievers (F1 & Recall). Results show that even powerful models like GPT-5-mini and Llama-3.3-70B exhibit performance variances, highlighting the importance of task-specific optimization.

Methods	NQ* F1	PopQA* F1	HotpotQA F1	2WikiMultihopQA F1	Musique F1	Avg. F1						
Graph Knowledge Retrievers (F1 Only)												
<i>DeepSeek-V3</i>												
Subgraph Retriever	27.20	61.42	42.36	34.96	13.29	35.85						
Triples Retriever	33.89	54.61	44.84	32.67	19.13	37.03						
ToG Retriever	28.29	64.02	44.72	50.98	17.99	41.20						
<i>GPT-5-mini</i>												
Subgraph Retriever	27.09	58.62	41.25	32.78	14.34	34.82						
Triples Retriever	34.62	54.99	43.07	29.66	18.29	36.13						
ToG Retriever	28.20	62.65	44.14	49.11	18.78	40.58						
<i>Llama-3.3-70B</i>												
Subgraph Retriever	26.90	62.05	41.81	34.56	16.11	36.29						
Triples Retriever	32.95	54.65	44.00	31.84	19.55	36.60						
ToG Retriever	26.31	62.48	45.45	50.42	18.59	40.65						
<i>gpt-oss-120b</i>												
Subgraph Retriever	26.92	58.83	40.60	33.35	13.39	34.62						
Triples Retriever	32.95	54.43	42.91	29.66	19.10	35.81						
ToG Retriever	26.00	60.82	43.52	48.43	18.30	39.41						
<i>Qwen2.5-72B-Instruct</i>												
Subgraph Retriever	28.12	62.23	42.57	36.17	14.85	36.79						
Triples Retriever	34.23	54.68	44.10	32.47	20.07	37.11						
ToG Retriever	28.73	64.21	45.54	50.19	18.37	41.41						
<i>Qwen3-235B</i>												
Subgraph Retriever	27.62	60.90	41.49	35.59	15.83	36.29						
Triples Retriever	34.12	55.12	41.67	32.81	20.71	36.89						
ToG Retriever	27.68	63.54	45.35	49.68	21.77	41.60						
<i>gemini-2.5-flash-lite</i>												
Subgraph Retriever	27.81	60.42	40.95	35.21	15.18	35.91						
Triples Retriever	33.52	54.73	44.02	33.02	19.91	37.04						
ToG Retriever	25.90	62.55	42.69	49.55	19.60	40.06						
AutoGraph-R1 (Qwen-3B)												
Subgraph Retriever	28.03	59.46	40.77	34.71	15.13	35.62						
Triples Retriever	33.67	56.76	46.94	36.09	21.41	38.97						
ToG Retriever	29.27	61.40	44.56	49.33	18.42	40.60						
AutoGraph-R1 (Qwen-7B)												
Subgraph Retriever	28.54	60.94	43.59	37.43	15.65	37.23						
Triples Retriever	33.98	58.02	48.28	36.04	20.56	39.38						
ToG Retriever	29.36	62.85	44.68	50.20	19.31	41.28						
Graph-based Text Retrievers (F1 & Recall@5)												
	F1	R@5	F1	R@5	F1	R@5	F1	R@5	F1	R@5	F1	R@5
<i>DeepSeek-V3</i>												
HippoRAG	39.82	93.90	66.29	94.90	54.33	68.30	49.59	72.69	27.79	47.23	47.56	75.40
HippoRAG2	38.45	93.70	65.73	96.10	54.89	69.68	52.23	75.17	28.14	47.76	47.89	76.48
<i>GPT-5-mini</i>												
HippoRAG	39.28	93.80	66.12	94.70	54.22	68.60	47.59	69.92	28.70	47.80	47.18	74.96
HippoRAG2	39.20	94.20	65.47	94.20	55.85	69.28	49.44	69.77	28.80	49.00	47.75	75.29
<i>Llama-3.3-70B</i>												
HippoRAG	39.29	94.30	63.57	94.00	56.14	60.00	52.70	73.12	28.32	47.41	48.00	73.77
HippoRAG2	39.29	94.50	65.18	95.50	56.29	70.71	53.93	75.89	26.97	48.27	48.33	76.97
<i>gpt-oss-120b</i>												
HippoRAG	39.52	93.70	64.82	92.90	53.11	65.61	47.79	68.57	25.37	45.43	46.12	73.24
HippoRAG2	39.45	93.60	64.78	93.10	54.53	66.51	49.49	69.49	27.03	46.97	47.06	73.93
<i>Qwen2.5-72B-Instruct</i>												
HippoRAG	39.73	94.60	65.18	94.00	53.98	67.92	50.68	73.06	26.42	46.93	47.20	75.30
HippoRAG2	39.29	95.10	65.71	94.60	57.34	70.42	53.21	76.00	26.71	48.63	48.45	76.95
<i>Qwen3-235B</i>												
HippoRAG	39.56	93.90	65.32	95.40	55.07	68.02	51.94	72.67	27.72	47.58	47.92	75.51
HippoRAG2	39.68	94.30	67.25	95.60	54.81	68.52	51.76	74.37	28.91	48.13	48.48	76.18
<i>gemini-2.5-flash-lite</i>												
HippoRAG	40.06	94.40	64.54	94.40	55.39	66.88	51.55	72.19	27.32	46.36	47.77	74.85
HippoRAG2	39.73	94.10	65.48	95.30	54.55	68.93	51.73	74.37	28.34	47.81	47.97	76.10
AutoGraph-R1 (Qwen-3B)												
HippoRAG	38.28	93.00	65.93	95.60	55.39	68.82	51.69	74.02	28.11	47.65	47.88	75.82
HippoRAG2	38.45	94.00	66.23	95.40	56.28	71.21	52.80	76.42	27.93	49.13	48.34	77.23
AutoGraph-R1 (Qwen-7B)												
HippoRAG	38.80	94.30	67.85	95.80	57.19	71.40	53.60	76.16	26.97	48.44	48.88	77.22
HippoRAG2	38.68	95.00	67.72	96.30	58.98	73.61	56.46	78.66	27.18	49.23	49.80	78.56

Table 6: Ablation study comparing On-Policy RL (GRPO) against Offline Filtering (RAFT) for Qwen2.5-7B. Both methods leverage our reward functions. Results show that reward-guided training consistently improves over the baseline, with GRPO and RAFT offering competitive trade-offs.

Retriever	Method	NQ	PopQA	HotpotQA	2Wiki	Musique	Avg.
<i>Graph Retrievers (Knowledge-Carrying Reward)</i>							
Subgraph	Base	28.07	55.43	41.66	33.97	15.24	34.87
	GRPO (On-Policy)	28.54	60.94	43.59	37.43	15.65	37.23
	RAFT (Offline)	28.72	61.27	42.14	35.78	15.41	36.66
Triples	Base	33.26	55.56	44.99	35.57	20.43	37.96
	GRPO (On-Policy)	33.98	58.02	48.28	36.04	20.56	39.38
	RAFT (Offline)	34.72	56.42	44.91	31.78	20.21	37.61
ToG	Base	25.59	57.53	43.93	46.03	18.46	38.31
	GRPO (On-Policy)	29.36	62.85	44.68	50.20	19.31	41.28
	RAFT (Offline)	29.65	63.30	45.47	50.10	20.32	41.77
<i>Text Retrievers (Knowledge-Indexing Reward)</i>							
HippoRAG	Base	37.16	65.95	55.50	53.01	26.03	47.53
	GRPO (On-Policy)	38.80	67.85	57.19	53.60	26.97	48.88
	RAFT (Offline)	39.60	66.22	55.77	51.01	29.29	48.38
HippoRAG2	Base	37.02	65.74	57.08	54.99	26.77	48.32
	GRPO (On-Policy)	38.68	67.72	58.98	56.46	27.18	49.80
	RAFT (Offline)	38.83	65.21	56.49	52.47	29.44	48.49

Table 7: Ablation on Judge Model Size. We report the downstream RAG F1 scores using Qwen-7B as the KG constructor. Training with a stronger Judge (7B) generally yields better downstream performance than a smaller Judge (3B).

Retriever	Judge	NQ	PopQA	HotpotQA	2wiki	Musique	Avg
Subgraph	3B	27.10	61.76	39.71	37.31	14.87	36.15
	7B	28.54	60.94	43.59	37.43	15.65	37.23
Triples	3B	34.20	59.20	48.40	31.12	19.04	38.39
	7B	33.98	58.02	48.28	36.04	20.56	39.38
ToG	3B	28.95	62.86	43.68	47.43	20.05	40.59
	7B	29.36	62.85	44.68	50.20	19.31	41.28

Table 8: Performance of the HippoRAG KG Construction pipeline when replacing the default KG constructor with our fine-tuned Qwen-7B model. The consistent improvement confirms that our optimized KG construction transfers benefits to established third-party RAG systems.

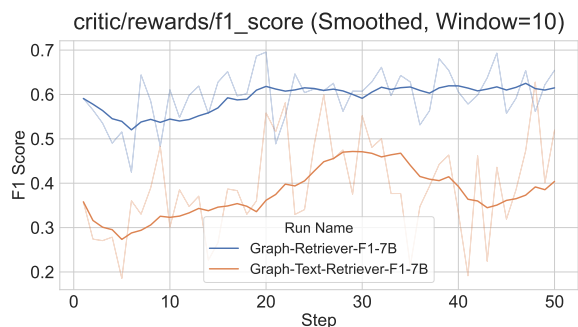
Method	NQ	PopQA	HotpotQA	2wiki	Musique	Avg
<i>HippoRAG F1 Score</i>						
Original (Base 7B)	28.78	37.83	41.34	33.66	18.85	32.09
+ Our Fine-tuned 7B	28.47	38.04	42.84	32.24	20.74	32.47
<i>HippoRAG Recall Score</i>						
Original (Base 7B)	63.89	51.65	83.60	76.76	53.88	65.96
+ Our Fine-tuned 7B	63.28	52.10	84.00	78.35	54.01	66.35

more focused set of relations. This demonstrates an optimization for efficiency, where the model actively filters out noise to create a cleaner, high-signal "index" for text retrieval.

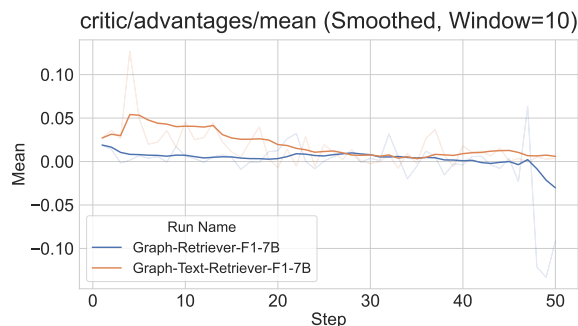
These results confirm that **AutoGraph-R1** is not a monolithic filter but a targeted optimization framework that reshapes the KG topology to maximize downstream utility.

B Impact of Using F1 Reward on Training Dynamics of AutoGraph-R1

To validate our choice of using task-specific rewards (R_C and R_I), we conducted an ablation study comparing them against the final answer’s F1 score as a direct reward signal. While the downstream performance results are reported in Table 3, here we analyze the training dynamics to further explain why the F1 reward is an inferior choice for guiding graph construction.



(a) Volatile reward signal.



(b) Stagnant advantage gain.

Figure 3: Unstable Training Dynamics Using a Naive F1 Reward. Training curves for the ablation study where the final RAG F1 score is used as the reward. (a) The F1 reward signal is highly volatile and shows no clear upward trend, providing a noisy and ineffective learning signal. (b) Consequently, the advantage gain remains flat and centered around zero, confirming that the policy is failing to find a consistent direction for improvement. This leads to stalled optimization, as reflected in the poor downstream results.

F1-based RL leads to unstable training Figure 3a illustrates the instability inherent in using the final F1 score as a reward. The reward curve (Figure 3a) exhibits high variance and lacks a clear, monotonic upward trend compared with using task specific reward, indicating a noisy learning signal.

C Training Dynamics of AutoGraph-R1

To verify the stability and effectiveness of our RL framework, we analyze the learning curves for both retrieval scenarios. Figure 4 illustrates the training process for the Graph Knowledge Retriever using the Deducible Reward (R_C). We observe a steady increase in the reward signal, indicating that the model progressively learns to construct graphs with higher deductive validity. Similarly, Figure 5 depicts the dynamics for the Graph-based Text Retriever using the Knowledge-Indexing Reward (R_I). The clear upward trend in recall-based reward, coupled with stable advantage estimates, confirms that AutoGraph-R1 successfully optimizes the graph structure for indexing utility without suffering from the volatility often associated with sparse RL rewards.

D Training Dynamics of Llama Models

Our initial experiments with the Llama model series showed a tendency to generate repetitive triples, leading to unstable training dynamics. As seen in Figure 6, the model’s generation length became erratic, with the mean response length fluctuating and the maximum often saturating the context window, indicating severe repetition loops.

To mitigate this issue, we augmented the reward for Llama experiments with a repetition penalty. The modified reward is defined as $\text{Reward} = R_{C/I} - \lambda_{\text{rep}} \cdot P_{\text{rep}}$, where λ_{rep} is a scaling hyperparameter controlling the strength of the penalty, and the repetition penalty is calculated as $P_{\text{rep}} = (|T_{\text{gen}}| - |T_{\text{unique}}|) / |T_{\text{gen}}|$. Furthermore, we applied a hard constraint that sets the reward to 0 if $P_{\text{rep}} > 0.3$. This approach effectively discourages repetition, encouraging the model to produce diverse and high-quality facts.

Table 9: Impact of our Knowledge-Carrying reward finetuned KG constructor on Graph Knowledge Retriever performance across the Llama model family.

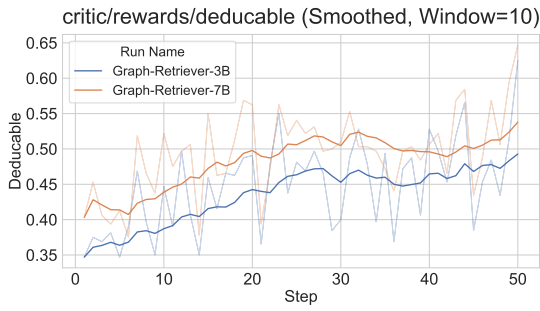
Methods	Simple QA		Multihop QA			Overall
	NQ*	PopQA*	HotpotQA	2WikiMultihopQA	Musique	Avg.
<i>Llama-3.2-1B</i>						
Subgraph (Base)	23.17	26.14	27.11	19.36	9.13	20.98
Subgraph (Ours)	25.72 ^{↑2.55}	47.32 ^{↑21.18}	35.84 ^{↑8.73}	25.05 ^{↑5.69}	11.90 ^{↑2.77}	29.17 ^{↑8.19}
Triples Retriever (Base)	21.51	24.55	29.86	18.54	10.79	21.05
Triples Retriever (Ours)	24.96 ^{↑3.45}	41.23 ^{↑16.68}	36.25 ^{↑6.39}	23.88 ^{↑5.34}	13.89 ^{↑3.10}	28.04 ^{↑6.99}
ToG (Base)	20.77	18.17	25.63	15.29	7.79	17.53
ToG (Ours)	19.74	34.20 ^{↑16.03}	30.70 ^{↑5.07}	21.35 ^{↑6.06}	9.60 ^{↑1.81}	23.12 ^{↑5.59}
<i>Llama-3.2-3B</i>						
Subgraph (Base)	25.19	44.07	34.86	26.23	12.26	28.52
Subgraph (Ours)	27.34 ^{↑2.15}	53.50 ^{↑9.43}	40.89 ^{↑6.03}	33.91 ^{↑7.68}	15.50 ^{↑3.24}	34.23 ^{↑5.71}
Triples Retriever (Base)	26.18	43.80	39.11	27.03	15.21	30.27
Triples Retriever (Ours)	30.93 ^{↑4.75}	51.50 ^{↑7.70}	43.67 ^{↑4.56}	31.49 ^{↑4.46}	18.21 ^{↑3.00}	35.16 ^{↑4.89}
ToG (Base)	20.00	31.33	31.93	27.20	11.14	24.32
ToG (Ours)	23.78 ^{↑3.78}	48.55 ^{↑17.22}	40.60 ^{↑8.67}	39.79 ^{↑12.59}	17.31 ^{↑6.17}	34.01 ^{↑9.69}

Table 10: Impact of our Knowledge-Indexing reward finetuned KG constructor on Graph-based Text Retriever performance across the Llama model family.

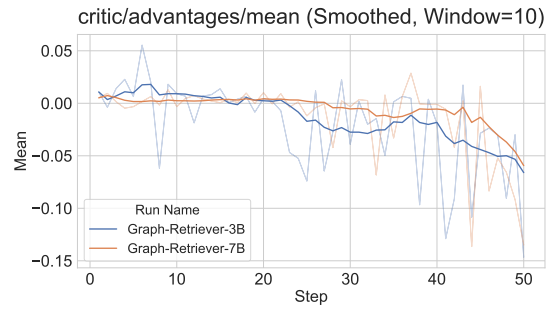
Methods	Simple QA		Multihop QA			Overall
	NQ*	PopQA*	HotpotQA	2WikiMultihopQA	Musique	Avg.
<i>Llama-3.2-1B</i>						
HippoRAG (Base)	28.94	39.16	39.37	27.26	18.89	30.72
HippoRAG (Ours)	35.49 ^{↑6.55}	52.63 ^{↑13.47}	47.31 ^{↑7.94}	38.13 ^{↑10.87}	20.17 ^{↑1.28}	38.75 ^{↑8.03}
HippoRAG2 (Base)	29.79	40.64	37.90	29.82	20.21	31.67
HippoRAG2 (Ours)	35.61 ^{↑5.82}	51.66 ^{↑11.02}	50.97 ^{↑13.07}	41.24 ^{↑11.42}	21.58 ^{↑1.37}	40.21 ^{↑8.54}
<i>Llama-3.2-3B</i>						
HippoRAG (Base)	38.35	58.81	50.68	44.26	23.63	43.15
HippoRAG (Ours)	38.77 ^{↑0.42}	59.62 ^{↑0.81}	54.82 ^{↑4.14}	47.03 ^{↑2.77}	25.26 ^{↑1.63}	45.10 ^{↑1.95}
HippoRAG2 (Base)	38.20	60.18	52.86	46.44	23.78	44.29
HippoRAG2 (Ours)	38.54 ^{↑0.34}	61.05 ^{↑0.87}	54.43 ^{↑1.57}	48.79 ^{↑2.35}	23.44	45.25 ^{↑0.96}

Table 11: Evaluating Knowledge Indexing Quality via Passage Recall across the Llama model family.

Methods	Simple QA		Multihop QA			Overall
	NQ*	PopQA*	HotpotQA	2WikiMultihopQA	Musique	Avg.
<i>Llama-3.2-1B</i>						
HippoRAG (Base)	57.50	51.30	37.86	36.78	24.35	41.56
HippoRAG (Ours)	66.30 ^{↑8.80}	71.60 ^{↑20.30}	56.56 ^{↑18.70}	56.29 ^{↑19.51}	33.99 ^{↑9.64}	56.95 ^{↑15.39}
HippoRAG2 (Base)	65.00	54.80	42.02	38.90	26.48	45.44
HippoRAG2 (Ours)	83.00 ^{↑18.00}	73.70 ^{↑18.90}	60.51 ^{↑18.49}	58.19 ^{↑19.29}	36.44 ^{↑9.96}	62.37 ^{↑16.93}
<i>Llama-3.2-3B</i>						
HippoRAG (Base)	87.70	85.40	60.80	63.97	39.02	67.38
HippoRAG (Ours)	90.90 ^{↑3.20}	88.50 ^{↑3.10}	66.23 ^{↑5.43}	68.77 ^{↑4.80}	41.63 ^{↑2.61}	71.21 ^{↑3.83}
HippoRAG2 (Base)	89.70	86.70	63.85	66.42	40.69	69.47
HippoRAG2 (Ours)	91.20 ^{↑1.50}	89.00 ^{↑2.30}	68.00 ^{↑4.15}	71.57 ^{↑5.15}	43.23 ^{↑2.54}	72.60 ^{↑3.13}

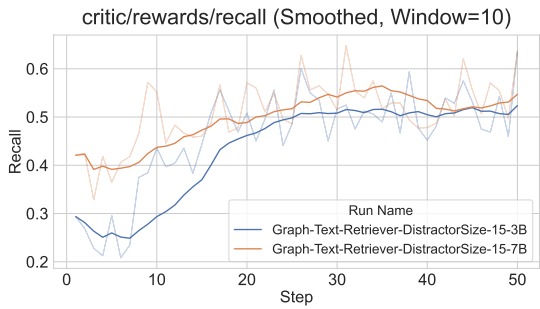


(a) Reward convergence.

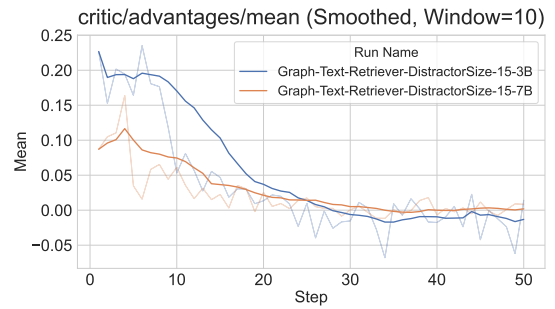


(b) Positive advantage gain.

Figure 4: **Effective Training Dynamics with the Deducible Reward (R_C)**. Training curves for the Graph Knowledge Retriever setting. (a) The reward, measuring answer deducibility, steadily increases and converges, demonstrating the policy is successfully learning its objective. (b) The advantage gain trends towards a small negative value, indicating that the value function’s estimate of expected reward is rising quickly while the policy makes stable, incremental improvements. This dynamic, coupled with the rising absolute reward, points to effective and controlled optimization.

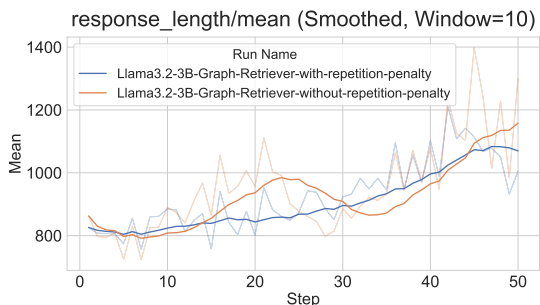


(a) Reward convergence.

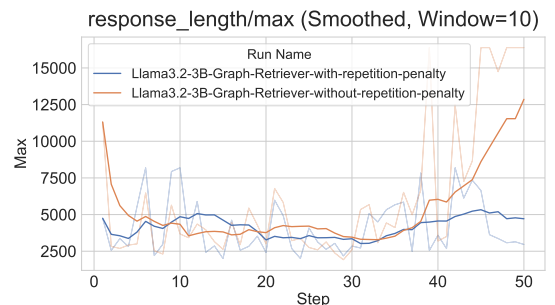


(b) Positive advantage gain.

Figure 5: **Effective Training Dynamics with the Knowledge-Indexing Reward (R_I)**. Training curves for the Graph-based Text Retriever setting. (a) The reward, measuring passage recall, shows a clear upward trend of improvement. (b) The advantage gain dynamic, paired with the rising reward curve, confirms that the policy is effectively learning from this stable, task-specific signal.



(a) Mean response length during training.



(b) Maximum response length during training.

Figure 6: **Response Length Dynamics of Llama Models**. The Llama models showed unstable response lengths due to repetition issues. (a) Mean response length fluctuated significantly. (b) Maximum response length often reached the context window limit, indicating severe repetition loops.

E Case Studies: The Functional Advantage of AutoGraph-R1

To qualitatively illustrate the benefits of our task-aware optimization, we present two case studies from the 2WikiMultiHopQA dataset. These examples demonstrate the **adaptive nature** of the utility objective: whereas Figure 1 showed the model learning to *establish concrete pathways*, these case studies highlight its ability to *bridge disconnected reasoning chains*. Together, they confirm that AutoGraph-R1 does not simply maximize or minimize graph size, but actively constructs the specific topology required for the downstream reasoning task.

E.1 Case Study 1: Comparative Reasoning

The first case study examines a question requiring a comparison between the death dates of two film directors. This task requires the KG to contain specific, comparable facts (i.e., dates) for multiple entities. As shown in Figure 7, the zero-shot KG fails because it does not extract the specific death dates needed for comparison. In contrast, the KG constructed by AutoGraph-R1 contains the necessary date information, as the RL training has taught the constructor that dates are critical for such questions. This complete evidence enables the LLM to easily answer the question correctly.

E.2 Case Study 2: Path-Based Reasoning

The second case study involves a 2-hop question that requires finding a path from a film to its director, and then from the director to their child. This task depends on the structural connectivity of the graph.

As shown in Figure 8, the zero-shot KG (top) fails critically. While it successfully extracts the first link in the path—‘(Los Pagares de Mendieta, directed by, Leopoldo Torres Ríos)’—it fails to extract the second, crucial link about the director’s child. The reasoning path is broken after the first hop, causing the QA system to fail. In contrast, the AutoGraph-R1 KG (bottom) explicitly contains the complete 2-hop reasoning path. It successfully extracts both ‘(Los Pagares de Mendieta, directed by, Leopoldo Torres Ríos)’ and ‘(Leopoldo Torres Ríos, father of, Leopoldo Torre Nilsson)’. The RL process has rewarded the constructor for building these essential connective trails, recognizing that entity linkage across different relationships is crucial for multi-hop QA.

Case Study 1: Zero-Shot KG (ToG Retriever) - Failed Answer

Question: Which film has the director who died first, The Goose Woman or You Can No Longer Remain Silent?

Retrieved Triples:

"(You Can No Longer Remain Silent, directed by, Robert A. Stemmle)",
"(Robert A. Stemmle, died in, Baden-Baden, Germany)",
"(The Goose Woman, directed by, Clarence Brown)",
"(Clarence Brown, was a, American film director)"
... (and other irrelevant triples)

Case Study 1: AutoGraph-R1 KG (ToG Retriever) - Correct Answer

Question: Which film has the director who died first, The Goose Woman or You Can No Longer Remain Silent?

Retrieved Triples:

"(You Can No Longer Remain Silent, directed by, Robert A. Stemmle)",
"**(Robert A. Stemmle, died on, 24 February 1974)**",
"(The Goose Woman, directed by, Clarence Brown)",
"**(Clarence Brown, died on, August 17, 1987)**",
...

Figure 7: Qualitative comparison for a **comparative reasoning** question. The zero-shot KG lacks specific death dates, leading to failure. The AutoGraph-R1 KG, optimized for task utility, successfully extracts the critical dates needed for comparison.

Case Study 2: Zero-Shot KG (ToG Retriever) - Failed Answer

Question: Who is the child of the director of film Los Pagares De Mendieta?

Retrieved Triples:

"(Los Pagares de Mendieta, directed by, Leopoldo Torres Ríos)",
"(Leopoldo Torres Ríos, age at death, 60)",
... (and other facts about the director, but not their child)

Case Study 2: AutoGraph-R1 KG (ToG Retriever) - Correct Answer

Question: Who is the child of the director of film Los Pagares De Mendieta?

Retrieved Triples:

"(Los Pagares de Mendieta, directed by, Leopoldo Torres Ríos)",
"(Leopoldo Torres Ríos, father of, Leopoldo Torre Nilsson)",
...

Figure 8: Qualitative comparison for a **2-hop path-based** question. The zero-shot KG extracts the first link (director of the film) but misses the second (child of the director), breaking the reasoning path. The AutoGraph-R1 KG successfully constructs the full path.

F AutoGraph-R1 Training Algorithm

The end-to-end training process for the AutoGraph-R1 KG constructor is formalized in Algorithm 1. The core idea is to iteratively construct a knowledge graph for a given query and its context documents, evaluate the graph's utility using a task-specific reward function, and then update the constructor's policy using the collected rewards.

G Prompts

This section details the specific prompts used in our experimental pipeline. The process begins with the graph construction prompt (Figure 10), which guides the LLM to extract triples from raw text. During RL training, the **Knowledge-Carrying Reward** (R_C) is determined using the deducibility judge prompt shown in Figure 11. For the final RAG answer generation step, we use distinct prompts tailored to the retrieved context: one for linearized graph triples (Figure 12) and another for raw text passages (Figure 13). Finally, Figure 9 shows the prompts used for our intrinsic graph quality analysis, where an LLM judge generates and answers multiple-choice questions to evaluate factual coverage.

Multiple-Choice Question Generation and Answering

MCQ Generation Prompt:
You are an expert in generating multiple-choice questions (MCQs) from scientific texts. Your task is to generate 5 multiple-choice questions based on the following passage.
Each question should:
- Focus on factual claims, numerical data, definitions, or relational knowledge from the passage.
- Have 4 options (one correct answer and three plausible distractors).
- Clearly indicate the correct answer.
The output should be in JSON format, with each question as a dictionary containing:
- "question": The MCQ question.
- "options": A list of 4 options (e.g., ["A: ..", "B: ..", "C: ..", "D: .."]).
- "answer": The correct answer (e.g., "A").
Passage: {passage}

MCQ Answering Prompt:
Given the contexts or evidences: {contexts}
Here is a multiple-choice question:
Question: {question}
Options: A. {options_0} B. {options_1} C. {options_2} D. {options_3}
Please select the correct answer by choosing A, B, C, or D. Respond with only the letter of your choice.

Figure 9: The prompt provided to the LLM judge (DeepSeek-V3) to evaluate triples extraction quality

Graph Construction

Graph Generation System Prompt:
You are an expert knowledge graph constructor. Your task is to extract factual information from the provided text and represent it strictly as a JSON array of knowledge graph triples.
Output Format
- The output must be a ****JSON array****.
- Each element in the array must be a ****JSON object**** with exactly three non-empty keys:
- "subject": the main entity, concept, event, or attribute.
- "relation": a concise, descriptive phrase or verb that describes the relationship (e.g., "founded by", "started on", "is a", "has circulation of").
- "object": the entity, concept, value, event, or attribute that the subject has a relationship with.
Constraints
- ****Do not include any text other than the JSON output.****
- Do not add explanations, comments, or formatting outside of the JSON array.
- Extract ****all possible and relevant triples****.
- All keys must exist and all values must be non-empty strings.
- The "subject" and "object" can be specific entities (e.g., "Radio City", "Football in Albania", "Echsmith") or specific values (e.g., "3 July 2001", "1,310,696").
- If no triples can be extracted, return exactly: '[]'.
Extracts for: {passage}

Figure 10: The prompt used for both zero-shot KG construction and fine-tuning KG constructor model during RL.

Deducible Judge

Deducible Judge Prompt:
As an advanced reading comprehension assistant, your task is to evaluate whether the provided knowledge graph (KG) context contains sufficient information to deduce the given true answer to the question. Analyze the KG context carefully and determine if it fully supports the true answer without requiring external knowledge. Respond with only 'Yes' or 'No', indicating whether the true answer can be deduced from the KG context.
Knowledge graph (KG) context: {triples string}
Question: {query}
True Answer: {answer}
Can the true answer be deduced from the KG context?
Answer 'Yes' or 'No' only.

Figure 11: The prompts for freeze LLM to determine the **Knowledge-Carrying Reward** (R_C). The 'Yes' or 'No' response serves as the binary reward signal.

Graph Retriever Answer Generation

Answer Generation Prompt For Graph Retriever:
As an advanced reading comprehension assistant, your task is to analyze extracted information and corresponding questions meticulously. If the knowledge graph information is not enough, you can use your own knowledge to answer the question. Your response start after "Thought: ", where you will methodically break down the reasoning process, illustrating how you arrive at conclusions. Conclude with "Answer: " to present a concise, definitive response as a noun phrase, no elaborations.
{triples string}
{question}
Thought:

Figure 12: The prompt used by the final answer generator when the retrieved evidence consists of linearized knowledge graph triples.

Graph Text Retriever Answer Generation

Answer Generation Prompt:
As an advanced reading comprehension assistant, your task is to analyze text passages and corresponding questions meticulously. If the information is not enough, you can use your own knowledge to answer the question. Your response start after "Thought: ", where you will methodically break down the reasoning process, illustrating how you arrive at conclusions. Conclude with "Answer: " to present a concise, definitive response as a noun phrase, no elaborations.
{Retrieved Texts}
{question}
Thought:

Figure 13: The prompt used by the final answer generator when the retrieved evidence consists of raw text passages.

Algorithm 1 AutoGraph-R1 Training Loop

- 1: **Input:** Training dataset $\mathcal{S} = \{(q_i, y_i, \mathcal{D}_{q_i})\}_{i=1}^N$, where \mathcal{D}_{q_i} are the context documents for query q_i .
 - 2: **Input:** KG constructor policy π_{θ}^{KG} (an LLM).
 - 3: **Input:** Frozen retriever $\mathcal{R}_{\text{frozen}}$ (either a graph knowledge retriever or a graph-based text retriever).
 - 4: **Input:** Chosen reward function R_{task} (either R_C or R_I).
 - 5: **Initialize:** Policy parameters θ .
 - 6: **for** each training step **do**
 - 7: Sample a minibatch of data $\{(q, y, \mathcal{D}_q)\}$ from \mathcal{S} .
 - 8: Initialize an empty list of trajectories ‘trajectories’.
 - 9: **for** each sample (q, y, \mathcal{D}_q) in the minibatch **do**
 - 10: \triangleright **Step 1: Construct the Knowledge Graph**
 - 11: Generate the graph by sampling from the policy: $\mathcal{G} \sim \pi_{\theta}^{KG}(\cdot | \mathcal{D}_q)$.
 - 12: \triangleright **Step 2: Determine Task-Specific Reward**
 - 13: **if** R_{task} is Knowledge-Carrying Reward (R_C) **then**
 - 14: Use the frozen retriever $\mathcal{R}_{\text{graph}}$ to get evidence $\mathcal{P}(q)$ from \mathcal{G} .
 - 15: Calculate reward $r = R_C(q, y, \mathcal{P}(q))$ using Eq. (1).
 - 16: **else if** R_{task} is Knowledge-Indexing Reward (R_I) **then**
 - 17: Use the frozen retriever $\mathcal{R}_{\text{text}}$ to get passages $\mathcal{T}(q)$ from \mathcal{G} .
 - 18: Calculate reward $r = R_I(q, y, \mathcal{T}(q))$ using Eq. (2).
 - 19: **end if**
 - 20: Store the generation trajectory (actions taken to build \mathcal{G}) and the final reward r in ‘trajectories’.
 - 21: **end for**
 - 22: \triangleright **Step 3: Update Policy Parameters**
 - 23: Compute the policy gradient $\nabla_{\theta} J(\theta)$ using the stored ‘trajectories’ and a policy optimization algorithm (e.g., GRPO).
 - 24: Update the policy parameters: $\theta \leftarrow \theta - \eta \cdot \nabla_{\theta} J(\theta)$.
 - 25: **end for**
 - 26: **Return:** Optimized KG constructor parameters θ .
-

Table 12: Detailed statistics of Knowledge Graphs constructed by different models across four benchmarks. The table tracks the total count of extracted entities, relations, triples, the density of extraction (Triples/Doc), and the diversity of the relation vocabulary (Unique Relation Types).

Dataset	Model Variant	Total Ent.	Total Rel.	Total Trip.	Trip./Doc	Unique Rel.
2021Wiki	Qwen-3B (Base)	57,369	56,908	56,908	8.13	20,879
	Qwen-3B (Knwl-Carrying)	67,660	65,413	65,413	9.34	28,933
	Qwen-3B (Knwl-Indexing)	68,631	68,670	68,670	9.81	19,246
	Qwen-7B (Base)	75,521	70,264	70,264	10.04	23,887
	Qwen-7B (Knwl-Carrying)	65,324	66,370	66,370	9.48	32,173
	Qwen-7B (Knwl-Indexing)	72,529	61,324	61,324	8.76	22,201
2WikiMultihop	Qwen-3B (Base)	37,333	38,140	38,140	6.23	9,881
	Qwen-3B (Knwl-Carrying)	44,520	43,623	43,623	7.13	13,177
	Qwen-3B (Knwl-Indexing)	42,681	43,632	43,632	7.13	8,940
	Qwen-7B (Base)	45,740	46,074	46,074	7.53	11,362
	Qwen-7B (Knwl-Carrying)	42,995	44,997	44,997	7.35	16,920
	Qwen-7B (Knwl-Indexing)	44,675	40,192	40,192	6.57	10,609
HotpotQA	Qwen-3B (Base)	70,592	73,674	73,674	7.99	23,994
	Qwen-3B (Knwl-Carrying)	77,593	76,943	76,943	8.34	30,816
	Qwen-3B (Knwl-Indexing)	80,355	83,920	83,920	9.10	21,411
	Qwen-7B (Base)	82,253	81,237	81,237	8.81	24,996
	Qwen-7B (Knwl-Carrying)	74,647	76,700	76,700	8.32	33,335
	Qwen-7B (Knwl-Indexing)	79,861	71,598	71,598	7.76	23,677
Musique	Qwen-3B (Base)	77,462	78,988	78,988	6.78	29,255
	Qwen-3B (Knwl-Carrying)	89,247	88,331	88,331	7.58	39,115
	Qwen-3B (Knwl-Indexing)	91,223	93,930	93,930	8.06	27,132
	Qwen-7B (Base)	97,294	91,395	91,395	7.84	32,310
	Qwen-7B (Knwl-Carrying)	89,547	92,605	92,605	7.94	40,078
	Qwen-7B (Knwl-Indexing)	96,242	82,145	82,145	7.05	30,521

Average Structural Metrics (Qwen-7B Family)

Metric	Base	Carrying	Indexing
Avg. Triples per Doc	8.56	8.27	7.54
Total Unique Rel. Types	23,139	30,627	21,752