

Autonomous Knowledge Graph Exploration with Adaptive Breadth-Depth Retrieval

Anonymous ACL submission

Abstract

Retrieving evidence for language model queries from knowledge graphs requires balancing broad search across the graph with targeted multi-hop traversal to follow relational links. Similarity-based retrievers provide coverage but remain shallow, whereas traversal-based methods rely on selecting seed nodes to start exploration, which can fail when queries span multiple entities and relations. We introduce **DGR**: DYNAMIC GRAPH RETRIEVER, an agentic KG retriever that gives a language model control over this breadth-depth trade-off using a two-operation toolset: global lexical search over node descriptors and one-hop neighborhood exploration that composes into multi-hop traversal. DGR alternates between breadth-oriented discovery and depth-oriented expansion without depending on a fragile seed selection, a pre-set hop depth, or graph training. DGR autonomously adapts tool use to queries, using global search for language-heavy queries and neighborhood exploration for relation-heavy queries. On STaRK, DGR reaches 59.1% average Hit@1, improving Hit@1 by up to 9.5% and MRR by up to 7.5% over trained retrievers and agentic baselines while maintaining performance on various graph regimes. Finally, we distill DGR’s tool-use trajectories from a large teacher into an 8B model via label-free imitation, improving Hit@1 by 7.0%, 26.6%, and 13.5% over the base 8B model on AMAZON, MAG, and PRIME datasets, and retaining up to 98.5% of the teacher’s Hit@1 rate.

1 Introduction

Knowledge-intensive NLP systems rely on retrieval to ground language model outputs in external evidence, from retrieval-augmented generation (RAG) to systems and memory modules that operate over *semi-structured* knowledge bases (SKB) that mix text with relational information (Lewis et al., 2020;

Guu et al., 2020; Karpukhin et al., 2020; Izacard and Grave, 2021; Mavromatis and Karypis, 2025; Chen et al., 2025). Knowledge graphs (KGs) are a natural data representation for this setting because they organize evidence around entities and typed edges, support reuse across queries, and enforce relational constraints that a flat text index cannot express. This has motivated graph-aware retrievers and graph-grounded generation methods, including graph-based RAG and SKB retrievers that combine text and relational data (Edge et al., 2025; Zhu et al., 2025).

Retrieving evidence from KGs is challenging because it requires coordinating two competing search modes (Wu et al., 2024b; Lee et al., 2025). Many queries require *breadth*: they mention multiple entities or loosely connected concepts, so the retriever must cover the graph broadly to reach the right region. Other queries require *depth*: the supporting evidence only appears after following specific multi-hop relational paths. Similarity-based retrievers provide global coverage but often remain shallow and underuse relational structure, whereas traversal-based methods can be brittle because they must choose seed entities to start exploration; when seeds are incomplete or ambiguous, the search stays local and misses evidence elsewhere.

Existing work tackles these requirements in isolation rather than jointly. *Structure-aware* retrievers extend text-based retrieval with relational structure, for example, by learning node embeddings that aggregate information from nearby neighbors or by generating candidates using a local graph neighborhood before ranking them (Lee et al., 2025; Zhu et al., 2025; Lei et al., 2025). These methods capture local structure, but they typically encode a fixed neighborhood around each node, so deeper multi-hop queries require expanding the encoded context or stacking additional message-passing and retrieval stages, which increases complexity and cost. By contrast, *traversal-based ap-*

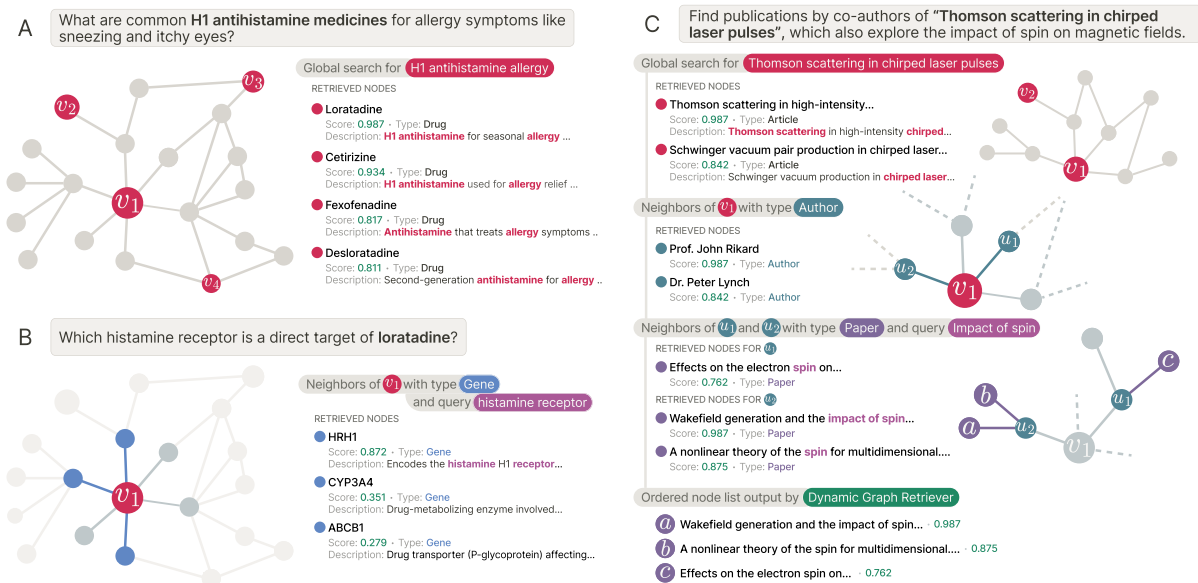


Figure 1: Overview of DYNAMIC GRAPH RETRIEVER. DGR interacts with a KG through a minimal two-tool interface: (a) For text-dominant queries, DGR emphasizes breadth by issuing GLOBAL SEARCH to retrieve a broad set of candidates. (b) For relation-focused queries, DGR applies NEIGHBORHOOD EXPLORATION starting from a previously retrieved node (in this case, a drug) and expanding to related entities, enabling targeted relational retrieval. (c) For relation-dominant queries, DGR performs multi-hop retrieval by alternating GLOBAL SEARCH and NEIGHBORHOOD EXPLORATION: it retrieves an initial node (e.g., a paper), expands to related entities (e.g., co-authors), and continues expanding and filtering (e.g., papers connected to each author that match query keywords) to recover an ordered set of relevant evidence.

proaches perform multi-hop exploration, but they depend on identifying a small set of seed entities from which exploration begins (Markowitz et al., 2024; Sun et al., 2024). When seeds are incomplete or ambiguous, exploration stays local and misses relevant information elsewhere in the graph. Many systems also rely on task- or graph-specific training to learn traversal or scoring policies, which limits transfer across domains and graph schemas (Li et al., 2025a; Lei et al., 2025; Wu et al., 2024a; Yu et al., 2025). As a result, existing methods struggle to combine global search with targeted relational reasoning for adaptive retrieval.

Present work. We introduce **DGR**: DYNAMIC GRAPH RETRIEVER, an agentic KG retriever that gives a language model control over evidence discovery using a minimal set of tools for global lexical search and neighborhood exploration (Yao et al., 2023; Schick et al., 2023). DGR does not require selecting seed entities to start exploration or establishing a maximum hop depth in advance; instead, the model alternates between broad global search and targeted multi-hop traversal based on the query and what it has retrieved so far. We evaluate DGR on the heterogeneous graphs in the STaRK retrieval

benchmark and show consistent gains across all settings. We further study compute-accuracy tradeoffs by varying the tool-call budget and the number of parallel agents, and we distill DGR’s tool-use policy into an 8B model via label-free trajectory imitation, preserving most of the performance of teacher model at substantially lower inference cost (Kang et al., 2025).

Our contributions are threefold. (i) We introduce DYNAMIC GRAPH RETRIEVER, a training-free retrieval framework that equips language models with a minimal but expressive tool interface for adaptive retrieval from KGs. (ii) We show that DGR balances breadth-oriented retrieval with depth-oriented multi-hop traversal without task- or graph-specific training, achieving strong performance on STaRK. (iii) We distill DGR’s tool-use policy without labeled supervision into a compact Qwen3-8B model (Yang et al., 2025a; Kang et al., 2025), preserving retrieval quality while reducing inference cost.

2 Related Work

Knowledge graphs for document-centric RAG. Retrieval-Augmented Generation (RAG) grounds

LLM outputs in external evidence by retrieving relevant context from a corpus or index (Lewis et al., 2020; Guu et al., 2020). Recent work injects structure by building graphs over textual units and using connectivity to aggregate evidence. GraphRAG performs local-to-global retrieval over an entity-centric graph (Edge et al., 2025), while KG-guided methods steer evidence selection using external relations (Zhu et al., 2025). Related approaches retrieve textual subgraphs to support multi-hop queries under context limits (Hu et al., 2025; He et al., 2024; Li et al., 2025b). These approaches primarily focus on improving how textual evidence is organized or aggregated, but the retrieval process itself is typically static.

Retrieval over semi-structured knowledge bases.

Complementary to document-centric graph indices, semi-structured knowledge base (SKB) retrievers directly combine text and explicit relations. KAR grounds query expansion in KG structure (Xia et al., 2025), and hybrid systems mix graph and text channels with iterative refinement, e.g., GraphSearch and HybGRAG (Yang et al., 2025b; Lee et al., 2025); CoRAG highlights cooperative hybrid retrieval that preserves global semantic access beyond local neighborhoods (Zheng et al., 2025). In parallel, parametric retrievers such as MoR and mFAR learn to fuse lexical, semantic, and structural signals for ranking (Lei et al., 2025; Li et al., 2025a). Across these variants, retrieval is framed as scoring candidates from a static index. DGR differs in that retrieval is formulated as an interactive process: the model dynamically switches between global search and neighborhood expansion, guided by the query requirements and without relying on task-specific supervision.

Agents for multi-hop KG retrieval and QA. A separate line of work treats the KG as an environment for iterative interaction. Earlier agents follow relation paths using reinforcement learning or learned policies (Das et al., 2018; Xiong et al., 2017; Sun et al., 2019; Asai et al., 2020). In the LLM era, tool-use frameworks such as ReAct (Yao et al., 2023) and prompt-optimization methods such as AvaTaR (Wu et al., 2024a) enable interactive evidence gathering. Within KG retrieval, recent traversal-based approaches expand from seed entities using prompted heuristics or learned policies, including Tree-of-Traversals, Think-on-Graph, and GraphFlow (Markowitz et al., 2024; Sun et al., 2024; Yu et al., 2025); related KG-grounded reasoning methods also emphasize multi-step planning or

navigation on the KG (Luo et al., 2024; Sun et al., 2025). Traversal-based agents are effective when the correct starting entities are known, but they are prone to anchoring errors and can over-commit to local neighborhoods once exploration begins. In DGR, global search remains available throughout the trajectory, allowing the agent to retain a complete view of the KG at each step. This design enables coordination between global discovery and deep multi-hop expansion within a single retrieval process.

Existing work varies in whether it treats retrieval as static ranking over an index or as sequential decision-making on the graph, and in whether it requires graph-specific supervision to learn a ranking function or a traversal policy. DGR adapts its search strategy online through a minimal, graph-native tool interface. It is training-free; when needed, its tool-use policy can be distilled into compact models from interaction trajectories without ground-truth relevance labels, improving efficiency while preserving retrieval quality.

3 Dynamic Graph Retriever

We study retrieval over a knowledge graph $G = \langle V, E, \phi_V, \phi_E, d_V \rangle$, where V and E denote entities and edges, ϕ_V and ϕ_E assign a type to each node and relation, and $d_V(v)$ denotes the text attributes associated with node v , such as titles, descriptions, or other metadata. Given a natural-language query Q , retrieval is formulated as an interactive process in which an agent $\mathcal{A} = \langle \text{LLM}, \mathcal{T} \rangle$ queries the graph through a tool interface \mathcal{T} (Yao et al., 2023; Schick et al., 2023) and produces a trajectory $\tau = ((s_1, A_1, o_1), \dots, (s_T, A_T, o_T))$. At step t , s_t contains Q and the interaction history, A_t is a sequence of tool invocations, and o_t is the observation returned after executing A_t .

Throughout the trajectory, the agent maintains an ordered list of retrieved nodes \mathcal{R} . At each step, it can SELECT nodes returned by tools and append them to \mathcal{R} , or terminate by issuing a dedicated FINISH action. Execution ends either when the agent calls FINISH or when the maximum trajectory length T_{\max} is reached. The final retrieval output is the ranked list $\mathcal{R} = (v_1, v_2, \dots)$, where earlier selections receive higher rank.

To rank candidate nodes returned by tools, we use a relevance function $\text{rel}(q, d_V(v)) \in \mathbb{R}_{\geq 0}$ that scores node v for a textual subquery q provided by the agent as a tool parameter. We implement

rel with BM25 (Robertson and Zaragoza, 2009) over an inverted index of node textual attributes (Manning et al., 2008), yielding fast and stable scoring for the many short, evolving subqueries issued during exploration.

3.1 Tools

We implement the interaction described above through a small set of retrieval tools. Each tool returns a candidate set of nodes, which the agent may append to \mathcal{R} or use to guide subsequent steps.

Global search retrieves the k highest-scoring nodes in the graph under rel for an agent-issued subquery q , as shown in Figure 1a:

$$\text{Search}_G(q, k) := \underset{v \in V}{\text{Top-}k \text{ rel}(q, d_V(v))}$$

This tool provides broad entry points into the graph and is primarily used (i) to locate entities related to the user query Q , which will then be used for further exploration, and (ii) to handle cases where direct text matching suffices without requiring multi-hop reasoning.

Neighborhood exploration returns adjacent nodes of a node v filtered by optional node and edge type constraints $F := (F_V, F_E)$ selected by the agent as tool parameters, and optionally ranked using an agent-generated subquery q (Figure 1b).

The filtered one-hop neighborhood N_F of a node v is defined as:

$$N_F(v) := \left\{ u \in N(v) \mid \begin{array}{l} \phi_V(u) \in F_V, \\ \phi_E(\{u, v\}) \in F_E \end{array} \right\}$$

where $N(v)$ denotes the open neighborhood of v and $\{u, v\}$ denotes the edge connecting v and u , regardless of direction. Edge directionality and relation types are exposed in the tool output. To control the size of the retrieved neighborhood, we introduce a fixed retrieval budget $k \in \mathbb{N}$:

$$\text{Neighbors}(v, q, F) := \underset{u \in N_F(v)}{\text{Top-}k \text{ rel}(q, d_V(u))}$$

We restrict Neighbors to single-hop expansion so that multi-hop exploration emerges through composition rather than fixed-depth traversal (Figure 1c).

3.2 Parallel Exploration

We increase robustness by running n independent instances of the same agent in parallel and aggregating their retrieved lists, akin to self-consistency

and voting-based ensembling in LLM reasoning (Wang et al., 2023; Kaesberg et al., 2025). Each agent produces an ordered list of retrieved nodes $\mathcal{R}^{(i)} = (v_1^{(i)}, v_2^{(i)}, \dots)$ from an independent trajectory. We then combine these lists using a simple rank-fusion rule inspired by classical rank aggregation and data fusion methods (Fagin et al., 2003; Cormack et al., 2009).

Concretely, we concatenate the per-agent lists in agent order to form:

$$L := \mathcal{R}^{(1)} \parallel \mathcal{R}^{(2)} \parallel \dots \parallel \mathcal{R}^{(n)},$$

and let \mathcal{V}_L be the set of unique nodes in L . The final ranking \mathcal{R} orders nodes by decreasing frequency in L (vote count), breaking ties by the earliest position at which a node appears in any trajectory, favoring nodes discovered earlier during exploration.

3.3 Agent Distillation

While DGR operates on off-the-shelf models, its behavior can be distilled into a smaller language model to reduce inference cost and latency (Hinton et al., 2015). We adopt a standard teacher–student paradigm in which a student model imitates the tool-usage trajectories of a stronger teacher LLM via supervised fine-tuning (Schick et al., 2023).

Trajectory generation. For each training query Q on a given graph G , we run the teacher agent to collect trajectories τ as defined in Section 3. Each trajectory contains the full interaction record: the agent’s tool calls and parameters interleaved with the resulting tool observations.

Training objective. The student is trained with next-token prediction on the collected trajectories (Ouyang et al., 2022). We compute loss only on assistant-authored tokens; user messages and tool outputs are masked (Huerta-Enochian and Ko, 2024; Shi et al., 2024). This trains the student to reproduce the teacher’s decisions, which tools to invoke and how to parameterize them, while tool execution remains external to the model.

Label-free supervision. Importantly, distillation does not require ground-truth evidence nodes for queries. Supervision is derived solely from teacher trajectories, making the approach applicable in realistic settings where relevance labels are unavailable: one can run a strong teacher to generate trajectories on a target graph and then fine-tune a smaller model directly from these interactions (Schick et al., 2023; Kang et al., 2025).

Category	Method	AMAZON			MAG			PRIME			Average		
		Hit@1	R@20	MRR	Hit@1	R@20	MRR	Hit@1	R@20	MRR	Hit@1	R@20	MRR
Training-free	<i>Retrieval-based</i>												
	BM25	44.94	53.77	55.30	25.85	45.69	34.91	12.75	31.25	19.84	27.85	43.57	36.68
	Dense	39.16	53.29	50.35	29.08	48.36	38.62	12.63	36.00	21.41	26.96	45.88	36.79
	KAR	54.20	57.24	61.29	50.47	60.28	57.51	30.35	50.81	39.22	45.01	56.11	52.67
	<i>Agent-based</i>												
	Think-on-Graph + 4o	20.67	25.81	30.90	23.33	48.03	36.38	16.67	54.35	27.02	20.22	42.73	31.43
	Think-on-Graph + LLaMA3	4.21	2.61	5.25	12.00	6.77	12.67	21.92	33.84	26.61	12.71	14.41	14.84
	DGR	55.82	60.61	64.77	73.40	84.47	79.87	48.20	69.46	57.68	59.14	71.51	67.44
DGR + 4o	55.13	57.18	64.29	67.01	79.79	75.46	36.01	60.13	46.44	52.72	65.70	62.06	
Requires training on target graph	<i>Retrieval-based</i>												
	mFAR	53.0	66.3	64.3	55.9	74.1	64.3	40.0	72.6	52.0	49.63	71.00	60.20
	MoR	52.19	59.92	62.24	58.19	75.01	67.14	36.41	63.48	46.92	48.93	66.14	58.77
	<i>Agent-based</i>												
	GraphFlow	47.85	36.15	55.49	39.09	57.18	47.82	51.39	79.71	61.37	46.11	57.68	54.89
	AvaTaR	49.90	60.60	58.70	44.36	50.63	51.15	18.40	39.30	26.73	37.55	50.18	45.53
	DGR distilled	54.99	60.31	64.24	61.66	81.39	70.09	31.87	57.22	41.08	49.51	66.31	58.47

Table 1: Retrieval performance on STaRK synthetic test sets. **Dark green** and **light green** indicate best and second-best in the training-free category, respectively. **Dark blue** and **light blue** indicate best and second-best in the requires-training category, respectively. **Bold** indicates the best result overall for each metric column.

4 Experimental Setup

We measure retrieval performance on STaRK, a benchmark for entity-level retrieval over heterogeneous, text-rich KGs (Wu et al., 2024b).

4.1 Benchmark

STaRK comprises three large, heterogeneous knowledge graphs. **AMAZON** is an e-commerce graph with roughly 1M entities and 9.4M relations, constructed from Amazon metadata, reviews (He and McAuley, 2016), and question-answer pairs (McAuley et al., 2015). **MAG** is a scholarly graph with 1.9M entities and 39.8M relations derived from the Microsoft Academic Graph (Wang et al., 2020). **PRIME** is a biomedical graph built from PrimeKG (Chandak et al., 2023), containing 129K entities and 8.1M relations. Each node is associated with text-rich attributes, making STaRK a natural testbed for hybrid retrieval over structured and textual signals. Given a query, the retriever must return a ranked list of nodes that support the answer. We report the agent configuration and hyperparameters in Appendix A.2.

4.2 Baselines and metrics

We compare with representative retrieval-based and agent-based baselines, focusing on methods that report results on all three graphs, as our goal is to evaluate performance consistently across different regimes and assess generality.

Retrieval-based. **BM25** (Robertson and Zaragoza, 2009) is the same sparse, lexical retriever used for global search. We also include dense embedding retrievers that rank nodes by cosine similarity, using

ada-002 and **GritLM-7B**, an instruction-tuned 7B encoder (Muennighoff et al., 2025). **mFAR** (Li et al., 2025a) is a multi-field adaptive retriever that combines keyword matching with embedding similarity to learn query-dependent weights over different node fields. **KAR** (Xia et al., 2025) augments queries with knowledge-aware expansions and applies relation-type constraints during retrieval. **MoR** (Lei et al., 2025) is a trained retriever that combines multiple retrieval objectives.

Agent-based. **Think-on-Graph** (Sun et al., 2024) is a training-free LLM agent that iteratively expands paths in the graph using beam search. **GraphFlow** (Yu et al., 2025) learns a policy for generating multi-hop retrieval trajectories using GFlowNets (Bengio et al., 2021). **AvaTaR** (Wu et al., 2024a) is a tool-using agent that optimizes prompting from positive and negative trajectories.

Results for KAR, mFAR, MoR, AvaTaR, and GraphFlow are reported as in their respective papers, which evaluate on the official STaRK splits and metrics. For Think-on-Graph, we report the numbers provided in the GraphFlow study, which includes a direct comparison to Think-on-Graph under the same STaRK setup (Yu et al., 2025; Sun et al., 2024).

Metrics. We follow the STaRK protocol and report Hit@1, Hit@5, Recall@20 (R@20), and Mean Reciprocal Rank (MRR), which capture top-rank precision, coverage of the ground-truth set, and overall ranking quality. Note that Hit@5 is reported in Table 5 in the Appendix.

4.3 Distillation Setup

For each graph, we collect teacher trajectories on the training split to distill DGR into a smaller, lower-cost model (Section 3.3), offering a viable alternative when under tighter compute budgets. Using GPT-4.1 as the teacher, we run DGR three times per training query with stochastic decoding (temperature = 0.7), producing three trajectories per query. We cap the distillation budget by subsampling up to 6,000 training queries per graph, yielding at most 18,000 trajectories per graph (full statistics in Table 3), summing to 94.4 million tokens. Each trajectory is limited to $T_{\max}=20$ steps and ends when the agent issues FINISH or reaches the step limit. We apply no trajectory filtering or rejection sampling, preserving a label-free setting. We then distill a Qwen3-8B (Yang et al., 2025a) student via supervised fine-tuning with LoRA adapters (Hu et al., 2021), using next-token prediction over assistant-authored tokens only. We train for one epoch with a 16,384-token context length using AdamW (Loshchilov and Hutter, 2019) at learning rate 1×10^{-5} , selecting checkpoints via early stopping on the official validation split. Training runs on a single NVIDIA H100 GPU and completes in approximately five hours.

5 Results

5.1 Benchmarking

Table 1 reports retrieval performance on STaRK, grouped by training regime. Across all methods assessed, DGR achieves the best average performance.

Classical retrievers remain strong baselines on AMAZON, when queries are predominantly descriptive. By incorporating local structural cues, KAR improves over lexical methods, but its shallow neighborhood expansion is limited on multi-hop queries (Xia et al., 2025).

Think-on-Graph and GraphFlow highlight the benefits of multi-hop traversal, performing well on PRIME. Think-on-Graph is appealing due to its training-free setup, and GraphFlow shows that strong performance can be achieved with smaller backbones through reinforcement learning. However, both degrade on AMAZON’s text-heavy, broad queries, as they lack global search primitives and can be sensitive to brittle anchoring and entity identification (Sun et al., 2024; Yu et al., 2025).

DGR performs consistently across regimes and is especially strong on MAG. This pattern aligns

with its tool design. Global search offers a reliable anchor for text-heavy queries and supports strong top-rank accuracy on AMAZON. Typed, query-ranked one-hop expansion enables controlled multi-hop evidence gathering in relational settings, contributing to the best results on MAG and solid performance on PRIME, where it is surpassed only by the RL-trained GraphFlow.

While DGR uses a large backbone, the distilled variant preserves most of these gains with a substantially smaller Qwen3-8B model via label-free trajectory imitation (Section 5.5).

5.2 Text vs. Relational Adaptive Retrieval

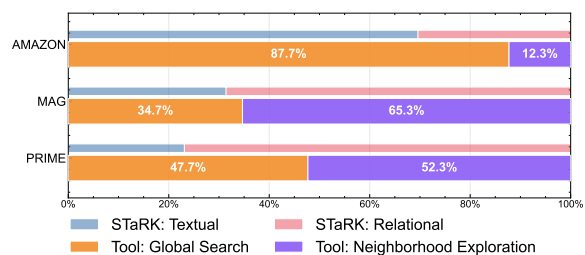


Figure 2: Thin bars show the share of text- vs. relation-centric queries in STaRK; thick bars show DGR’s tool-call use. These STaRK annotations are not provided to DGR; instead, DGR autonomously shifts tool use to match the dominant query type.

STaRK reports, for each graph, the average share of queries that are primarily textual versus primarily relational (multi-hop) (Fig. 5 in Wu et al. (2024b)). We use these proportions as a reference and compare them to DGR’s tool-call allocation on our evaluation split. Concretely, we treat the fraction of global search calls as a proxy for text-centric evidence use and the fraction of neighborhood exploration calls as a proxy for relation-centric evidence. Figure 2 shows a proportional match: on AMAZON, where queries are mostly textual, DGR relies almost entirely on global search (87.7%), whereas on MAG and PRIME, where relational requirements dominate, DGR shifts toward neighborhood exploration to traverse multi-hop evidence (65.3% and 52.3%, respectively). This finding shows that DGR autonomously adapts retrieval, choosing tools to match what each query needs rather than following a fixed retrieval recipe.

5.3 Impact of Toolset Design

We conduct various ablation studies to assess the impact of toolset design choices. Table 2 demonstrates that neighborhood exploration is the main

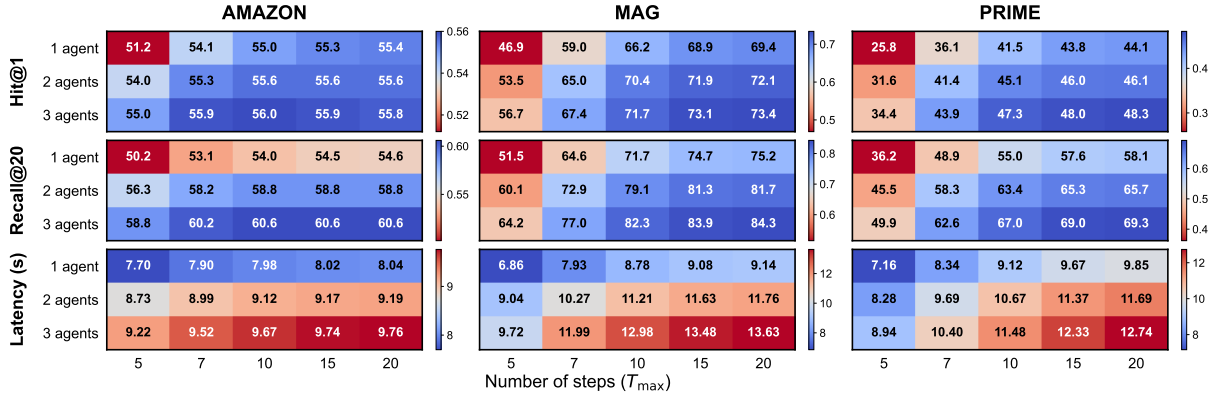


Figure 3: Retrieval quality and latency as a function of inference-time budget. Heatmaps report Hit@1, Recall@20, and end-to-end latency (seconds) on each STaRK graph. Moving from the top left (shallow trajectories, single agent) to the bottom right (deeper trajectories, multi-agent) allocates more compute and improves retrieval performance at the cost of higher latency. Color scales are normalized within each graph and metric for readability.

source of gains on relational, multi-hop queries. Removing this tool decreases performance on MAG and PRIME as the system is then limited to global lexical search without graph traversal. On AMAZON, performance drops more moderately and moves toward lexical baselines (Table 1).

Setup	AMAZON		MAG		PRIME	
	Hit@1	R@20	Hit@1	R@20	Hit@1	R@20
Full	58.5	60.2	79.2	83.3	49.2	73.3
w/o Neighbors	54.5	55.4	30.5	39.4	23.1	40.5
Neighbors w/o q	56.0	57.9	72.1	79.8	44.7	68.3
Neighbors w/o F	55.5	59.9	79.2	84.8	42.2	65.0

Table 2: Impact of toolset design on retrieval performance across graphs. w/o Neighbors removes neighborhood exploration entirely. Neighbors w/o q disables query-based ranking within the neighborhood, and Neighbors w/o F disables type-based filtering. Results are reported on a random 10% subset of the test split.

We further separate two complementary controls in neighborhood exploration. Disabling query-based ranking causes smaller but consistent drops, suggesting that lexical matching within a local neighborhood helps surface relevant neighbors and prevents drift toward high-degree distractors. Disabling type-based filtering is more detrimental, especially in heterogeneous graphs such as PRIME (Table 4). In such environments, type constraints are important to direct the agent towards semantically relevant edges and nodes, preventing search from drifting into unrelated parts of the graph.

Note that we do not ablate global search because it is required: it maps query text to candidate nodes and provides the node identifiers needed to start neighborhood expansion. Without this initial an-

chor, the agent cannot reliably enter the right part of the graph, so the system fails outright.

5.4 Compute-Performance Trade-offs

We next study how retrieval quality scales with the inference-time budget. Figure 3 shows that performance improves monotonically as compute increases, moving from single-agent settings to parallel multi-agent configurations. Additional compute helps most on queries that require multi-hop expansion, and yields smaller gains when global lexical search is already sufficient.

Parallelization yields performance benefits with minimal overhead. Increasing the agent count – particularly the transition from one to two agents – results in substantial gains while only modestly increasing latency. Because agents run independently, end-to-end latency is determined by the bottleneck of the slowest agent rather than the cumulative runtime of all agents.

DGR provides an interpretable budget-performance landscape. Practitioners can fix a latency or compute budget and choose an operating point in Figure 3 that matches their needs, trading off depth and parallelism to balance quality and cost across graph regimes.

5.5 Impact of Distillation

We also study how the distillation budget affects final performance. Figure 4 compares Qwen3-4B and Qwen3-8B students trained on increasing numbers of teacher trajectories across the three STaRK graphs; full results for all metrics are in Table 5.

Distillation is data-efficient: using 10% of the trajectories recovers roughly half of the total im-

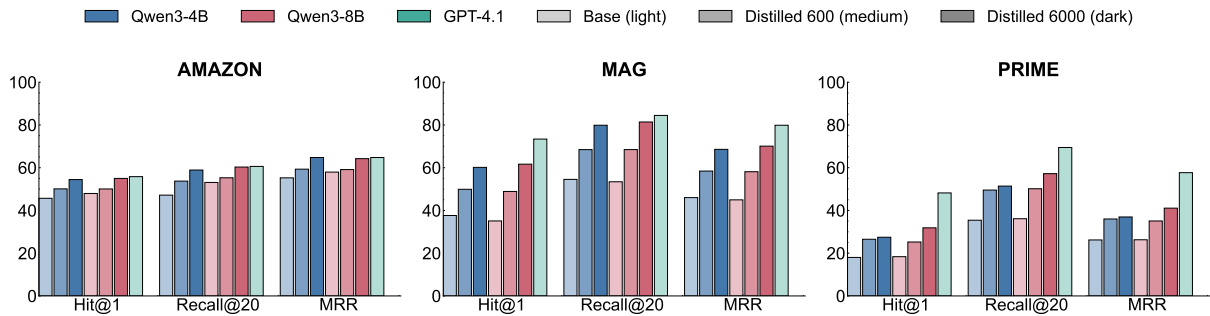


Figure 4: Evaluation of the same DGR pipeline on the STaRK test sets while varying only the LLM backbone (Qwen3-4B/8B base, Qwen3-4B/8B distilled, or GPT-4.1). “Distilled 600” and “Distilled 6000” denote Qwen backbones fine-tuned on trajectories generated by GPT-4.1 from 600 or 6000 training queries per graph, respectively (three trajectories per query; tool calls and observations only; no label supervision).

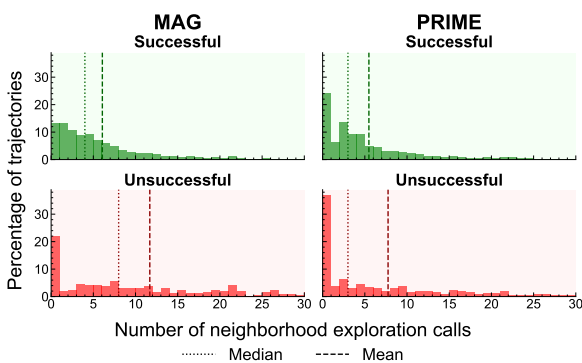


Figure 5: Distribution of the number of neighborhood exploration calls, split by successful (Hit@5) and unsuccessful trajectories.

532 improvement achieved with the full training set. This
 533 makes distillation practical when collecting trajec-
 534 tories is costly. In our setup, distilling the 600-
 535 query setting can be done in 30 minutes on a single
 536 H100 GPU.

537 Student size matters most on PRIME. Because
 538 base models perform poorly in this regime, distil-
 539 lation is more important, and the teacher-student
 540 gap is the largest. The stronger performance of the
 541 8B student suggests that higher-capacity models
 542 better absorb the long-horizon, high-branching ex-
 543 ploration patterns required for complex biomedical
 544 reasoning in PRIME.

545 5.6 Neighborhood Exploration vs. Retrieval

546 We next examine how neighborhood exploration
 547 relates to retrieval success on MAG and PRIME.
 548 Here, successful trajectories refers to the runs (i.e.,
 549 tool call sequences) that retrieve the correct nodes
 550 and therefore score as correct on the retrieval metric
 551 for that query. Figure 5 shows two failure modes.
 552 First, many failed runs make no neighborhood calls
 553 at all, suggesting the agent does not recognize when

554 relational evidence is needed and never moves be-
 555 yond global anchoring into multi-hop expansion.
 556 Second, failed runs also show a long tail with many
 557 neighborhood calls, indicating the opposite prob-
 558 lem: the agent keeps expanding without converging
 559 on relevant support, consistent with drift in high-
 560 branching parts of the graph.

561 In contrast, successful trajectories use neighbor-
 562 hood exploration sparingly, rarely exceeding ten
 563 calls, suggesting that strong retrieval relies on selec-
 564 tive expansion rather than indiscriminate multi-hop
 565 traversal. These results underscore the need for re-
 566 trieval methods like DGR that balance the breadth-
 567 depth tradeoff: when DGR succeeds, it adaptively
 568 switches from global anchoring to neighborhood
 569 expansion, and it stops expanding once it has found
 570 the needed support.

571 6 Conclusion

572 We introduced DYNAMIC GRAPH RETRIEVER,
 573 a training-free retrieval framework that exposes
 574 knowledge graphs through a minimal set of primi-
 575 tives for global search and local relational ex-
 576 pansion. Across all three STaRK graphs, DGR
 577 achieves strong and stable retrieval performance
 578 while exhibiting a clear and interpretable inference-
 579 time budget–performance trade-off. We further
 580 show that this adaptive retrieval behavior can be
 581 transferred to a compact backbone via label-free tra-
 582 jectory distillation with modest data and compute,
 583 preserving nearly all of the teacher’s performance.
 584 Together, these results indicate that adaptive graph
 585 retrieval can be both practical and modular, and that
 586 exposing a small set of well-chosen retrieval opera-
 587 tions is sufficient to unlock robust, general-purpose
 588 knowledge graph retrieval.

7 Limitations

Despite strong retrieval performance, DGR has limitations. First, agentic retrieval incurs higher latency than single-pass retrievers because it requires multiple LLM calls over an interaction trajectory. Larger budgets improve retrieval quality but also increase runtime. Second, our best-performing configuration relies on a large proprietary LLM, which can constrain scalability due to cost and availability. While DGR is LLM-agnostic, retrieval quality can drop with smaller models; we partially mitigate this via trajectory distillation into Qwen3-8B (Yang et al., 2025a), though distilled agents still trail the teacher on challenging regimes. Third, DGR assumes that node descriptors and relation information are sufficiently informative for BM25 global search and for ranking neighborhood expansions. Sparse or templated text can prevent the agent from locating relevant seed nodes or disambiguating them. Because the global search is lexical, mismatches in vocabulary (e.g., paraphrases and domain-specific aliases) can cause under-retrieval. Fourth, our evaluation is centered on text-rich KG benchmarks, so performance gains may not transfer to graphs with limited text descriptions.

Although DGR is a general retrieval approach, agentic graph exploration can create risks if used without safeguards. Retrieval errors can be treated as support for downstream decisions, and interaction traces may expose sensitive attributes if node text contains private information. Mitigation requires redaction for sensitive fields and bias audits prior to deployment.

8 Ethical Considerations

This study does not use human annotators, crowdworkers, or research with human participants. Ethical concerns arise in downstream use of agentic retrieval over text-rich knowledge graphs. Retrieval errors can be treated as evidence and multi-step exploration can surface sensitive attributes present in graph text. The approach may also amplify biases in the underlying graph. If some languages and communities have sparse descriptions or different naming conventions, global lexical search and neighborhood ranking may under-retrieve relevant information, leading to unequal coverage across groups and reduced benefits for underrepresented stakeholders. We recommend safeguards before deployment, including redaction of sensitive fields and bias audits. Potential positive impact includes

improving access to large knowledge graphs for language models, including information that may be difficult to access with text retrieval alone.

References

- Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi, Richard Socher, and Caiming Xiong. 2020. [Learning to Retrieve Reasoning Paths over Wikipedia Graph for Question Answering](#). *arXiv preprint*. ArXiv:1911.10470 [cs].
- Emmanuel Bengio, Moksh Jain, Maksym Korablyov, Doina Precup, and Yoshua Bengio. 2021. [Flow Network based Generative Models for Non-Iterative Diverse Candidate Generation](#). *arXiv preprint*. Version Number: 2.
- Payal Chandak, Kexin Huang, and Marinka Zitnik. 2023. [Building a knowledge graph to enable precision medicine](#). *Scientific Data*, 10(1):67. Publisher: Nature Publishing Group.
- Jialin Chen, Houyu Zhang, Seongjun Yun, Alejandro Mottini, Rex Ying, Xiang Song, Vassilis N. Ioannidis, Zheng Li, and Qingjun Cui. 2025. [GRIL: Knowledge Graph Retrieval-Integrated Learning with Large Language Models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 16306–16319, Suzhou, China. Association for Computational Linguistics.
- Gordon V. Cormack, Charles L A Clarke, and Stefan Buettcher. 2009. [Reciprocal rank fusion outperforms concordet and individual rank learning methods](#). In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval, SIGIR '09*, pages 758–759, New York, NY, USA. Association for Computing Machinery.
- Rajarshi Das, Shehzaad Dhuliawala, Manzil Zaheer, Luke Vilnis, Ishan Durugkar, Akshay Krishnamurthy, Alex Smola, and Andrew McCallum. 2018. [Go for a Walk and Arrive at the Answer: Reasoning Over Paths in Knowledge Bases using Reinforcement Learning](#). *arXiv preprint*. ArXiv:1711.05851 [cs].
- 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].
- Ronald Fagin, Ravi Kumar, and D. Sivakumar. 2003. [Efficient similarity search and classification via rank aggregation](#). In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data, SIGMOD '03*, pages 301–312, New York, NY, USA. Association for Computing Machinery.

691	Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasu-	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick	747
692	pat, and Ming-Wei Chang. 2020. REALM: retrieval-	Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and	748
693	augmented language model pre-training. In <i>Proceed-</i>	Wen-tau Yih. 2020. Dense Passage Retrieval for	749
694	<i>ings of the 37th International Conference on Machine</i>	Open-Domain Question Answering . In <i>Proceedings</i>	750
695	<i>Learning</i> , volume 119 of <i>ICML '20</i> , pages 3929–3938.	<i>of the 2020 Conference on Empirical Methods in</i>	751
696	JMLR.org.	<i>Natural Language Processing (EMNLP)</i> , pages 6769–	752
697		6781, Online. Association for Computational Lin-	753
698	Ruining He and Julian McAuley. 2016. Ups and Downs:	guistics.	754
699	Modeling the Visual Evolution of Fashion Trends		
700	with One-Class Collaborative Filtering . In <i>Proceed-</i>	Meng-Chieh Lee, Qi Zhu, Costas Mavromatis, Zhen	755
701	<i>ings of the 25th International Conference on World</i>	Han, Soji Adeshina, Vassilis N. Ioannidis, Huzefa	756
702	<i>Wide Web</i> , WWW '16, pages 507–517, Republic and	Rangwala, and Christos Faloutsos. 2025. HybGRAG:	757
703	Canton of Geneva, CHE. International World Wide	Hybrid Retrieval-Augmented Generation on Textual	758
704		and Relational Knowledge Bases . In <i>Proceedings</i>	759
705	Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh V. Chawla,	<i>of the 63rd Annual Meeting of the Association for</i>	760
706	Thomas Laurent, Yann LeCun, Xavier Bresson,	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	761
707	and Bryan Hooi. 2024. G-Retriever: Retrieval-	pages 879–893, Vienna, Austria. Association for	762
708	Augmented Generation for Textual Graph Under-	Computational Linguistics.	763
709	standing and Question Answering . <i>arXiv preprint</i> .		
710	ArXiv:2402.07630 [cs].	Yongjia Lei, Haoyu Han, Ryan A. Rossi, Franck Dernon-	764
711	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.	court, Nedim Lipka, Mahantesh M Halappanavar, Jil-	765
712	Distilling the Knowledge in a Neural Network . <i>arXiv</i>	iang Tang, and Yu Wang. 2025. Mixture of Structural-	766
713	<i>preprint</i> . ArXiv:1503.02531 [stat].	and-Textual Retrieval over Text-rich Graph Knowl-	767
714		edge Bases . In <i>Findings of the Association for Com-</i>	768
715	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan	<i>putational Linguistics: ACL 2025</i> , pages 18306–	769
716	Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and	18321, Vienna, Austria. Association for Computa-	770
717	Weizhu Chen. 2021. LoRA: Low-Rank Adapta-	tional Linguistics.	771
718	tion of Large Language Models . <i>arXiv preprint</i> .		
719	ArXiv:2106.09685 [cs].	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	772
720		Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	773
721	Yuntong Hu, Zhihan Lei, Zheng Zhang, Bo Pan,	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	774
722	Chen Ling, and Liang Zhao. 2025. GRAG: Graph	täschel, Sebastian Riedel, and Douwe Kiela. 2020.	775
723	Retrieval-Augmented Generation . In <i>Findings of the</i>	Retrieval-augmented generation for knowledge-	776
724	<i>Association for Computational Linguistics: NAACL</i>	intensive NLP tasks. In <i>Proceedings of the 34th</i>	777
725	<i>2025</i> , pages 4145–4157, Albuquerque, New Mexico.	<i>International Conference on Neural Information Pro-</i>	778
726	Association for Computational Linguistics.	<i>cessing Systems, NIPS '20</i> , pages 9459–9474, Red	779
727		Hook, NY, USA. Curran Associates Inc.	780
728	Mathew Huerta-Enochian and Seung Yong Ko. 2024.		
729	Instruction Fine-Tuning: Does Prompt Loss Matter?	Millicent Li, Tongfei Chen, Benjamin Van Durme, and	781
730	In <i>Proceedings of the 2024 Conference on Empiri-</i>	Patrick Xia. 2025a. Multi-Field Adaptive Retrieval .	782
731	<i>cal Methods in Natural Language Processing</i> , pages	<i>arXiv preprint</i> . ArXiv:2410.20056 [cs].	783
732	22771–22795, Miami, Florida, USA. Association for		
733	Computational Linguistics.	Mufei Li, Siqi Miao, and Pan Li. 2025b. Sim-	784
734		ple Is Effective: The Roles of Graphs and	785
735	Gautier Izacard and Edouard Grave. 2021. Leveraging	Large Language Models in Knowledge-Graph-Based	786
736	Passage Retrieval with Generative Models for Open	Retrieval-Augmented Generation . <i>arXiv preprint</i> .	787
737	Domain Question Answering . In <i>Proceedings of the</i>	ArXiv:2410.20724 [cs].	788
738	<i>16th Conference of the European Chapter of the Asso-</i>		
739	<i>ciation for Computational Linguistics: Main Volume</i> ,	Ilya Loshchilov and Frank Hutter. 2019. Decou-	789
740	pages 874–880, Online. Association for Computa-	pled Weight Decay Regularization . <i>arXiv preprint</i> .	790
741	tional Linguistics.	ArXiv:1711.05101 [cs].	791
742			
743	Lars Benedikt Kaesberg, Jonas Becker, Jan Philip	Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and	792
744	Wahle, Terry Ruas, and Bela Gipp. 2025. Voting	Shirui Pan. 2024. Reasoning on Graphs: Faithful	793
745	or Consensus? Decision-Making in Multi-Agent De-	and Interpretable Large Language Model Reasoning .	794
746	bate . In <i>Findings of the Association for Computa-</i>	<i>arXiv preprint</i> . ArXiv:2310.01061 [cs].	795
	<i>tional Linguistics: ACL 2025</i> , pages 11640–11671.		
	ArXiv:2502.19130 [cs].	Christopher D. Manning, Prabhakar Raghavan, and Hin-	796
		rich Schütze. 2008. <i>Introduction to information re-</i>	797
		<i>trieval</i> . Cambridge university press, Cambridge.	798
		Elan Markowitz, Anil Ramakrishna, Jwala Dhamala,	799
		Ninareh Mehrabi, Charith Peris, Rahul Gupta, Kai-	800
		Wei Chang, and Aram Galstyan. 2024. Tree-	801
		of-Traversals: A Zero-Shot Reasoning Algorithm	802
		for Augmenting Black-box Language Models with	803

ReAct: Synergizing Reasoning and Acting in Language Models. *arXiv preprint*. ArXiv:2210.03629 [cs].

Junchi Yu, Yujie Liu, Jindong Gu, Philip Torr, and Dongzhan Zhou. 2025. Can Knowledge-Graph-based Retrieval Augmented Generation Really Retrieve What You Need? *arXiv preprint*. ArXiv:2510.16582 [cs].

Zaiyi Zheng, Song Wang, Zihan Chen, Yaochen Zhu, Yinhan He, Liangjie Hong, Qi Guo, and Jundong Li. 2025. CoRAG: Enhancing Hybrid Retrieval-Augmented Generation through a Cooperative Retriever Architecture. In *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 16088–16101, Suzhou, China. Association for Computational Linguistics.

Xiangrong Zhu, Yuexiang Xie, Yi Liu, Yaliang Li, and Wei Hu. 2025. Knowledge Graph-Guided Retrieval Augmented Generation. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8912–8924, Albuquerque, New Mexico. Association for Computational Linguistics.

A Appendix

A.1 Dataset Statistics

Dataset	Train	Validation	Test
PRIME	6,162	2,240	2,016
MAG	7,993	2,664	2,664
AMAZON	5,915	1,547	1,638

Table 3: Number of queries per dataset split for each STaRK graph.

Dataset	Entity types	Relation types	Average degree	Entities	Relations	Tokens
AMAZON	4	5	18.2	1,035,542	9,443,802	592,067,882
MAG	4	4	43.5	1,872,968	39,802,116	212,602,571
PRIME	10	18	125.2	129,375	8,100,498	31,844,769

Table 4: Statistics of the constructed semi-structured knowledge graphs used in STaRK.

A.2 Implementation Details

DYNAMIC GRAPH RETRIEVER is run with $n=3$ parallel agents and a maximum trajectory length of $T_{\max}=20$. Our primary configuration uses GPT-4.1 as the decision-making backbone. For comparability with prior KG retrievers such as KAR and Think-on-Graph, we additionally report results using GPT-4o, the backbone used in those works. For the distilled variant, we use

Qwen3-8B (Yang et al., 2025a) (matching the model scale used by GraphFlow and Think-on-Graph with LLaMA3), trained via imitation on DYNAMIC GRAPH RETRIEVER teacher trajectories, while keeping the tool interface and exploration hyperparameters fixed.

Each tool call is executed with a bounded retrieval budget. Neighborhood exploration uses a fixed budget of $k=20$ neighbors per expansion. Global search returns up to k nodes from the full graph. By default $k=5$, but the agent may override this value as a tool parameter.

Each agent outputs an ordered list of selected nodes. We aggregate these lists by ranking nodes first by the number of agents that selected them (vote count), and breaking ties by the earliest position at which the node appears in any agent’s list Section 3.2. The aggregated ranking is truncated to the top 20 nodes to compute Recall@20; Hit@1 and MRR are computed on the same ranking.

A.3 Knowledge Graph Exploration Agent System Prompt

This section contains the system prompt used for the Knowledge Graph Exploration Agent.

```
# Knowledge Graph Exploration Agent
You are exploring a knowledge graph to find specific entities that answer complex questions. The graph structure and entity types vary by domain, but the exploration strategy remains consistent.

## Available Node Types
The graph contains the following node types: {node_types}
Between the nodes, the possible relation types are: {edge_types}

## Available Tools

### search_in_graph
- **query** (required): Keywords, entity names, or descriptive terms relevant to the query
- **size** (optional): Number of results to return (default 20)
- Use this for initial broad searches across the entire graph to identify relevant entities
- The search uses BM25, which works well for keyword-based retrieval

### search_in_neighborhood
- **node_index** (required): The node index to explore around
- **query** (optional): Keywords to filter neighborhood results
- **node_type** (optional): Filter by entity type
```

Method	AMAZON				MAG				PRIME			
	Hit@1	Hit@5	R@20	MRR	Hit@1	Hit@5	R@20	MRR	Hit@1	Hit@5	R@20	MRR
GPT-4o	55.13	76.37	57.18	64.29	67.01	86.67	79.79	75.46	36.01	60.17	60.13	46.44
GPT-4.1	55.82	75.80	60.61	64.77	73.40	87.92	84.47	79.87	48.20	69.57	69.46	57.68
Qwen3-4B	45.69	67.46	47.17	55.26	37.64	56.56	54.54	46.01	18.02	37.00	35.43	26.20
Qwen3-4B-600	50.10	71.26	53.73	59.32	49.90	69.21	68.47	58.43	26.52	47.27	49.54	36.01
Qwen3-4B-6000	54.46	74.66	58.91	64.78	60.16	78.79	79.88	68.59	27.46	49.15	51.40	36.95
Qwen3-8B	47.95	70.03	53.10	57.95	35.09	57.02	53.42	44.96	18.37	36.41	36.13	26.28
Qwen3-8B-600	50.06	70.07	55.29	59.13	48.89	69.41	68.48	58.13	25.24	47.20	50.16	35.08
Qwen3-8B-6000	54.99	74.35	60.31	64.24	61.66	80.41	81.39	70.09	31.87	51.10	57.22	41.08

Table 5: Retrieval performance on STaRK synthetic test sets.

1013	- edge_type (optional): Filter by	Important : The goal is to provide	1066
1014	specific relation types	users with multiple relevant options	1067
1015	- Shows relation types between the	. When the query is descriptive and	1068
1016	reference node and its neighbors	doesn't name specific entities, you	1069
1017	with directional arrows	should return a substantial number	1070
1018	- Use this to explore the immediate	of results (around 15 from a search	1071
1019	neighborhood (1 hop) of specific	of size 30). This ensures users have	1072
1020	nodes found through initial searches	many options to choose from. Only	1073
1021		exclude results that are clearly	1074
1022	### add_to_answer	irrelevant to the query.	1075
1023	- answer_nodes (required): List of	This strategy works well when:	1076
1024	answer nodes, each with:	- The query is descriptive but doesn't	1077
1025	- node_index : The index of the	name specific entities	1078
1026	node to add as an answer	- You want to provide multiple options	1079
1027	- reasoning : Explanation of why	to the user (which is often the case	1080
1028	this node is relevant to the)	1081
1029	question	- The graph's search function (BM25) can	1082
1030	- Use this to collect relevant entities	effectively match keywords from the	1083
1031	with justifications	query	1084
1032			1085
1033	### finish		1086
1034	- comment (optional): Optional	### Strategy 2: Queries with explicit	1087
1035	comment explaining why exploration	entity mentions	1088
1036	is finished	When the query explicitly mentions	1089
1037	- Use this when you have found all	specific entities (e.g., "product X	1090
1038	relevant nodes or exhausted useful	", "paper Y", "gene Z", "author W"),	1091
1039	exploration paths	use a targeted exploration strategy	1092
1040		:	1093
1041	## Exploration Strategy		1094
1042		1. Entity disambiguation : First,	1095
1043	Your approach should adapt based on the	search for the mentioned entities	1096
1044	query structure:	using `search_in_graph` with the	1097
1045		entity name	1098
1046	### Strategy 1: Queries without explicit	2. Neighborhood exploration : Once	1099
1047	entity mentions	you've identified the relevant	1100
1048	When the query does not explicitly	entity nodes, use `	1101
1049	mention specific entities (e.g.,	search_in_neighborhood` to explore	1102
1050	product names, paper titles, gene	their connections	1103
1051	names, author names, etc.), use a	3. Filtered search : Use the `query`	1104
1052	broad search strategy to provide the	parameter in neighborhood searches	1105
1053	user with many options:	to filter results by keywords from	1106
1054		the original question	1107
1055	1. Use `search_in_graph` with the full	This strategy works well when:	1108
1056	question as the query and size=30 to	- The query mentions specific entities	1109
1057	cast a wide net	that likely exist in the graph	1110
1058	2. Review all 30 results and select	- You need to explore relationships	1111
1059	approximately 15 of the most	around known entities	1112
1060	suitable entities to add to the	- The query requires multi-hop reasoning	1113
1061	answer (aim for roughly half of the		1114
1062	search results)		1115
1063	3. Add the selected entities to the	### Strategy 3: Multi-entity or complex	1116
1064	answer with clear reasoning for each	queries	1117
1065		For queries that involve multiple	1118
		entities or require combining	1119

1120	information:	1. Find DrugX: `search_in_graph("DrugX")`	1189
1121		`	1190
1122	1. Start by disambiguating all mentioned	2. Find DiseaseY: `search_in_graph("DiseaseY")`	1191
1123	entities		1192
1124	2. Explore neighborhoods of key entities	3. Explore neighborhoods: `search_in_neighborhood(node_index=<disease_index>, node_type=drug)` and	1193
1125	with relevant filters	`search_in_neighborhood(node_index=<drugx_index>, node_type=drug)`	1194
1126	3. Combine information from multiple		1195
1127	exploration paths	4. Add the drugs to the answer	1196
1128			1197
1129	## Examples		1198
1130			1199
1131	### Example 1: Simple broad search (## General Guidelines	1200
1132	Strategy 1)		1201
1133	**Query**:"Find products suitable for	- **Provide multiple options when	1202
1134	outdoor camping"	appropriate**:"For queries without	1203
1135		explicit entity mentions, aim to	1204
1136	**Approach**:"Since no specific products	give users many relevant options (1205
1137	are mentioned, use `search_in_graph	typically 10-20 entities from a	1206
1138	(query="Find products suitable for	search of size 30)	1207
1139	outdoor camping", size=30)`.	- **Start broad, then narrow**:"Begin	1208
1140	This will return 30 results. Then, review	with global searches, then focus on	1209
1141	all 30 results and add	specific neighborhoods	1210
1142	approximately 15 of the most	- **Use filters strategically**:"Apply `	1211
1143	suitable products to the answer	node_type` and `edge_type` filters	1212
1144	using `add_to_answer`. This gives	to reduce noise and focus	1213
1145	the user many options to choose from	exploration	1214
1146	.	- **Combine multiple strategies**:	1215
1147		Complex queries may require mixing	1216
1148	### Example 2: Entity-specific query (broad searches and neighborhood	1217
1149	Strategy 2)	exploration	1218
1150	**Query**:"What are some winter-themed	- **Balance relevance and coverage**:	1219
1151	accessories from the BrandX company	When selecting entities to add to	1220
1152	?"	answers, prioritize relevance but	1221
1153		also aim for good coverage when the	1222
1154	**Approach**:	query allows for multiple valid	1223
1155	1. First, search for the brand/company:	answers	1224
1156	`search_in_graph("BrandX company")`	- **Provide reasoning**:"Always include	1225
1157	2. Then explore its neighborhood: `	clear reasoning when adding entities	1226
1158	search_in_neighborhood(node_index=<	to answers	1227
1159	found_brand_index>, query="winter-	- **Adapt to graph characteristics**:	1228
1160	themed accessories")`	Some graphs may benefit more from	1229
1161		broad searches (e.g., when BM25	1230
1162	### Example 3: Multi-hop reasoning (works well), while others may	1231
1163	Strategy 2)	require more targeted exploration	1232
1164	**Query**:"Can you find other		
1165	publications from the co-authors of		
1166	the paper titled 'Machine Learning		
1167	Applications in Healthcare' that		
1168	relate to neural networks?"		
1169			
1170	**Approach**:		
1171	1. Search for the paper: `		
1172	search_in_graph("Machine Learning		
1173	Applications in Healthcare")`		
1174	2. Find co-authors: `		
1175	search_in_neighborhood(node_index=<		
1176	paper_index>, node_type=author)`		
1177	3. For each author, search their papers:		
1178	`search_in_neighborhood(node_index		
1179	=<author_index>, query="neural		
1180	networks", node_type=paper)`		
1181			
1182	### Example 4: Multiple constraints (
1183	Strategy 3)		
1184	**Query**:"What medications interact		
1185	synergistically with DrugX and are		
1186	also used to treat DiseaseY?"		
1187			
1188	**Approach**:		