

Pruning Minimal Reasoning Graphs for Efficient Retrieval-Augmented Generation

Ning Wang^{*1} Kuanyan Zhu^{*2} Daniel Yuehwoon Yee^{*3}
Yitang Gao⁴ Shiying Huang¹ Zirun Xu⁵ Sainyam Galhotra^{†1}

¹Cornell University ²University of Cambridge ³The University of Hong Kong ⁴HKUST
⁵University of British Columbia

nw366@cornell.edu kz345@cam.ac.uk u3636035@connect.hku.hk sg@cs.cornell.edu

Abstract

Retrieval-augmented generation (RAG) is now standard for knowledge-intensive LLM tasks, but most systems still treat every query as fresh, repeatedly re-retrieving long passages and re-reasoning from scratch, inflating tokens, latency, and cost. We present **Auto-PrunedRetriever**, a graph-style RAG system that *persists* the minimal reasoning subgraph built for earlier questions and *incrementally* extends it for later ones. AutoPrunedRetriever stores entities and relations in a compact, ID-indexed codebook and represents questions, facts, and answers as edge sequences, enabling retrieval and prompting over symbolic structure instead of raw text. To keep the graph compact, we apply a two-layer consolidation policy (fast ANN/KNN alias detection plus selective k -means once a memory threshold is reached) and prune low-value structure, while prompts retain only overlap representatives and genuinely new evidence. We instantiate two front ends: AUTOPRUNEDRETRIEVER-REBEL, which uses REBEL (Huguet Cabot and Navigli, 2021) as a triplet parser, and AUTOPRUNEDRETRIEVER-LLM, which swaps in an LLM extractor. On GraphRAG-Benchmark (Medical and Novel), both variants achieve **state-of-the-art complex reasoning accuracy**, improving over HippoRAG2 (Jiménez Gutiérrez et al., 2025) by roughly 9–11 points, and remain competitive on contextual summarize and generation. On our harder STEM and TV benchmarks, AutoPrunedRetriever again ranks first, while using up to two orders of magnitude fewer tokens than graph-heavy baselines, making it a practical substrate for long-running sessions, evolving corpora, and multi-agent pipelines.

1 Introduction

Retrieval-augmented generation (RAG) grounds LLMs in external knowledge, reducing hallucinations, enabling citation, and allowing updates without full retraining. Recent advances in dense retrieval and generation have yielded strong performance on open-domain question answering and tool-augmented assistants (Lewis et al., 2020; Karpukhin et al., 2020a; Izacard et al., 2022; Ram et al., 2023). However, moving from retrieving relevant text to solving complex, knowledge-intensive tasks still requires multi-hop reasoning: composing evidence across documents, enforcing temporal or structural constraints, and maintaining consistency across repeated or related queries.

In practice, most RAG systems treat each query independently. Even when multiple questions are closely related—or arise sequentially in agentic workflows, systems repeatedly re-retrieve overlapping passages and re-reason from scratch. This leads to substantial redundancy in retrieved context, inflated token usage, higher latency, and increased cost. These inefficiencies are especially pronounced in long-running sessions and multi-agent settings (e.g., planner–researcher–verifier pipelines), where similar reasoning chains are revisited many times (Yao et al., 2023; Wu et al., 2024)..

Graph-based RAG methods address some of these issues by lifting retrieval from flat text passages to structured representations over entities and relations (Han et al., 2024; Sun et al., 2022; Baek et al., 2023; Wang et al., 2023). By explicitly modeling compositional structure, GraphRAG systems improve multi-hop reasoning and disambiguation. Yet existing approaches still face three fundamental bottlenecks when deployed over evolving corpora and long reasoning chains. We demonstrate these challenges with an example.

Graph-based RAG methods address some of these issues by lifting retrieval from flat text passages to structured representations over entities and relations (Han et al., 2024; Sun et al., 2022; Baek et al., 2023; Wang et al., 2023). By explicitly modeling compositional structure, GraphRAG systems improve multi-hop reasoning and disambiguation. Yet existing approaches still face three fundamental bottlenecks when deployed over evolving corpora and long reasoning chains. We demonstrate these challenges with an example.

Example 1. Consider a small corpus containing five documents describing (i) a corporate acquisition, (ii) regulatory status under the EU Digital Services Act, (iii) GDPR incident histories, (iv)

^{*}Equal contribution.

[†]Corresponding author.

2024 vendor contracts, and (v) subsidiary relationships (Fig. 1). From this corpus, we ask three related questions:

Q1: Which subsidiaries acquired after Jan. 1, 2021 are subject to the EU Digital Services Act?

Q2: For those subsidiaries, did GDPR incident rates decrease post-acquisition?

Q3: Which 2024 vendor contracts involve those same subsidiaries?

A standard GraphRAG pipeline constructs an entity–relation graph over the corpus and answers each query via neighborhood expansion. Even in this minimal setting, three limitations emerge. **(M1) Graph construction and maintenance:** entity aliases and naming variants require global checks and relinking as new evidence arrives. **(M2) Reasoning granularity:** neighborhood-based expansion retrieves broad subgraphs around central entities (e.g., the parent company or regulator), rather than the few edges that realize each reasoning chain. **(M3) Redundant retrieval:** when Q1–Q3 are issued sequentially or by multiple agents, largely overlapping subgraphs are repeatedly retrieved and serialized, compounding token and latency costs.

Notably, the answers to Q2 and Q3 reuse much of the reasoning structure required for Q1. However, existing systems fail to exploit this overlap, repeatedly reconstructing context instead of persisting and extending prior reasoning.

These observations suggest a shift in perspective: retrieval should not aim to recover all potentially relevant context, but instead identify, cache, and reuse the minimal reasoning structure needed to answer a query, and incrementally extend that structure as new questions arrive. These limitations motivate three design principles: local incremental structure, path-centric retrieval, and exact symbolic reuse, which guide the design of AutoPrunedRetriever.

Our Approach. We introduce AutoPrunedRetriever, a structure-first RAG system that treats reasoning paths, rather than passages or neighborhoods, as the primary retrieval unit. AutoPrunedRetriever converts text into symbolic triples and represents questions, facts, and answers as compact sequences of entity–relation edges. The system persists only the minimal subgraphs that support successful reasoning and reuses them across later queries, avoiding repeated re-retrieval and re-prompting.

To keep memory compact and stable over time,

AutoPrunedRetriever applies a two-layer consolidation policy: a lightweight, continuous alias-detection pass using approximate nearest neighbors, and a periodic, budget-triggered consolidation step that merges aliases and prunes low-value structure. Retrieval is explicitly path-centric, scoring candidate reasoning chains rather than expanding broad neighborhoods. Prompts are constructed from a compact symbolic codebook that includes only novel or non-redundant evidence, substantially reducing token usage while preserving grounding in source text.

We instantiate AutoPrunedRetriever with two front ends: AutoPrunedRetriever-REBEL, which uses a fixed triplet parser, and AutoPrunedRetriever-LLM, which replaces it with an LLM-based extractor. Across the GraphRAG benchmark as well as harder STEM and TV reasoning datasets, both variants achieve state-of-the-art complex reasoning accuracy while using up to two orders of magnitude fewer tokens than graph-heavy baselines. These results demonstrate that pruned, persistent reasoning structure, not larger graphs or longer prompts, is the key substrate for efficient, long-running, and agentic RAG systems.

1.1 Design Principles

The design of AutoPrunedRetriever is guided by three principles, each directly addressing one of the limitations illustrated in Example 1.

P1. Local, incremental structure (addresses M1). To avoid the cost and brittleness of global graph maintenance, AutoPrunedRetriever builds reasoning structure locally and incrementally. Text is encoded into symbolic triples and grouped into small, coherent graphs that can be extended over time. Entity consolidation is applied selectively at the symbol level, allowing aliases to be merged without re-extracting text or relinking global structure.

P2. Path-centric retrieval (addresses M2). Reasoning is realized by short chains of entities and relations, not broad neighborhoods. AutoPrunedRetriever therefore treats edge sequences as the primary retrieval unit and scores candidate reasoning paths directly, avoiding the retrieval of large subgraphs that do not contribute to the required inference.

P3. Exact symbolic reuse (addresses M3). To prevent repeated serialization of overlapping context, reuse across queries is exact and symbolic

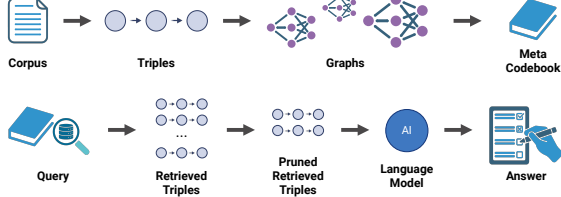


Figure 1: **AutoPrunedRetriever** pipeline: (1) encode into symbols and edges, (2) build chunked small graphs, (3) coarse→fine retrieval, (4) selector + compact prompt packing, (5) entity-only consolidation with a DPO wrapper to trade accuracy vs. tokens.

rather than textual. AutoPrunedRetriever caches reasoning subgraphs as compact sequences of entity–relation identifiers and constructs prompts that include only novel or non-redundant evidence, ensuring that token usage scales with new reasoning rather than repeated context.

We discuss more details about these principles in Section 2.

2 Method

2.1 Overview

Free text is a noisy, length-biased substrate for reasoning: it repeats the same facts in many surface forms, conflates content with string realization, and penalizes reuse by charging tokens repeatedly for identical information. These properties make it ill-suited for persistent, multi-hop reasoning across related queries.

AutoPrunedRetriever addresses these issues with a symbol-first pipeline that implements the three design principles introduced in Section 1.1. Specifically, we: (1) encode questions, answers, and facts into a shared symbolic codebook of entities and relations (Sec. 2.2); (2) build local, coherent reasoning graphs that serve as retrieval atoms (Sec. 2.3); (3) retrieve paths rather than neighborhoods via coarse-to-fine scoring in symbol space (Sec. 2.4); (4) select and package only non-redundant structure into compact prompts (Secs. 2.5, 2.6); and (5) consolidate entities incrementally to keep the persistent graph compact over time (Sec. 2.7). Finally, we adapt compression behavior to accuracy–efficiency tradeoffs using a lightweight DPO-trained policy (Sec. 2.8).

Each component is designed to ensure that retrieval, reasoning, and prompting scale with new reasoning rather than repeated context.

2.2 Step 1: Symbolic Encoding

Intuition. We normalize free-form text into a small set of symbols and the edges they instantiate. This exposes shared structure across questions and facts and makes reuse exact via IDs rather than brittle string matches. Because natural language is heavy-tailed, most informational mass is carried by a small “core” vocabulary; encoding those into a codebook yields large compression gains (see Lemma 1). In App. A.2 (Lemma 14, Proposition 16) also shows that such E–R–E codebooks mitigate the “token dilution” effect of long sequence embeddings.

Formulation. For a text span y (question, answer, fact), a parser (REBEL or an LLM) produces triples

$$\tau(y) \subseteq \mathcal{U}_E \times \mathcal{U}_R \times \mathcal{U}_E.$$

We maintain a meta-codebook

$$\mathcal{C} = (E, R, M, Q, A, F, E_{\text{emb}}, R_{\text{emb}}),$$

where E, R are entity and relation dictionaries, $M \subseteq E \times R \times E$ is the sparse set of unique edges, and Q, A, F store sequences of edge IDs for questions, answers, and facts.

An INDEXIFY map

$$\mathbf{y} = \text{INDEXIFY}(\tau(y); E, R, M) \in M^*$$

serializes text into edge indices, extending (E, R, M) only when new symbols/edges appear. We append \mathbf{y} to the appropriate store in Q, A, F . The embedding tables $E_{\text{emb}} : E \rightarrow \mathbb{R}^d$ and $R_{\text{emb}} : R \rightarrow \mathbb{R}^d$ give us E–R–E embeddings for each edge $(e, r, e') \in M$.

2.3 Step 2: Chunked Small Graphs (Local-First Construction)

Intuition. Instead of inserting every triple into a global graph, we build *small, locally coherent graphs* in a single pass. The question is no longer “where in the global graph does this triple live?” but “does this triple *fit* the current small graph?” This keeps working sets small and updates local, while preserving global consistency via shared IDs in M .

Runs and modalities. For each modality $y \in \{Q, A, F\}$ we build a run repository

$$\text{Runs}(y) \in (M^*)^*,$$

a list of edge-ID sequences (runs) that correspond to small graphs.

Streaming construction. We maintain a current small graph $G_k = (V_k, E_k)$ with an embedding centroid. For each incoming triple, we compute a *fit score* that combines:

1. **Semantic cohesion:** cosine similarity between the triple embedding and the centroid of G_k .
2. **Structural continuity:** a bonus when the triple continues a path (e.g., $\text{tail}_{\text{prev}} = \text{head}_{\text{new}}$) or reuses nodes in V_k .

If the (bonus-adjusted) score exceeds a threshold τ , we append to G_k ; otherwise we close G_k , linearize its edges into $\mathbf{g}_k \in M^*$, store $\mathbf{g}_k \in \text{Runs}(y)$, and start G_{k+1} . This yields maximal locally coherent segments (Lemma 2).

Boundary refinement. A single pass can overcut near boundaries, so we run a local *merge-if-not-a-true-cut* test on adjacent runs: we re-embed the concatenation of the boundary region, re-run the same segmenter, and merge if the new segmentation does *not* place a cut at the original boundary. Surviving boundaries are self-consistent fixed points of the segmenter (Lemma 3).

The result is a sequence of compact, coherent runs that serve as retrieval atoms. We quantify their intra-run cohesion and the induced working-set reduction for retrieval in Lemmas 4 and Lemmas 5.

2.4 Step 3: Coarse to Fine Path Retrieval

Intuition. Naively embedding every query against every run is slow and favors long, noisy chunks. We instead *first* work in symbol space to get a small shortlist and *then* do a more detailed triple-level check on that shortlist. This two-layer scheme keeps cost near $O(k)$ rather than $O(n)$ while preserving precision (Lemma 7).

Coarse stage (symbol-space recall). Each run decodes to triples (h, ρ, t) . We pool entity embeddings into $E(\cdot)$ and relation embeddings into $R(\cdot)$. For a query q and candidate run f we compute a simple *max-pair* score over entities and relations,

$$s_{\text{coarse}}(q, f) = w_{\text{ent}} \max_{i,j} \cos(E(q)_i, E(f)_j) + w_{\text{rel}} \max_{p,r} \cos(R(q)_p, R(f)_r),$$

and keep the top- k runs as a high-recall shortlist I_k . This stage only touches small entity/relation sets and is therefore very fast.

Fine stage (triple-aware re-ranking). For each candidate in I_k , we linearize triples $(h \ \rho \ t)$ into short lines, embed them, and build a similarity matrix \mathbf{S} between query and candidate lines. The final score is a weighted sum of five simple terms computed on \mathbf{S} : a top- t mean over the best entries (relational strength), a coverage term counting how many query triples are well matched, a soft many-to-many overlap term, a greedy 1:1 match encouraging parsimonious alignment, and a small whole-chunk bonus gated by full-text similarity to avoid “longer is better”. We then take the global TopM runs across all queries and channels (answers / facts / prior questions). Exact formulas and hyperparameters are deferred to App. A.3.

2.5 Step 4: Knowledge Selection

Intuition. Real corpora have heavy-tailed reuse of motifs, so naive “pull everything similar” either misses paraphrased support or floods the prompt with near-duplicates. We therefore operate on runs and, for each channel (answers, facts, related questions), choose an action

$$s \in \{\text{include all}, \text{unique}, \text{not include}\}.$$

Here *include all* keeps all surviving runs (e.g., for safety/audit regimes), *unique* keeps one representative per semantic cluster to avoid echoing and token bloat, and *not include* drops the channel when it adds little beyond context. Semantic clusters are defined in run-embedding space; we pick a consensus representative per cluster, and show in Lemma 8 that this representative stays close to paraphrastic variants.

2.6 Step 5: Compact prompt Construction

Given the selected runs, we assemble the minimal symbolic state the LLM needs:

$$\begin{aligned} \mathbf{U} &= \mathbf{q} \cup \mathbf{a} \cup \mathbf{f}, \\ E' &= \{h : t : (h, \rho, t) \in \mathbf{U}\}, \\ R' &= \{\rho : (h, \rho, t) \in \mathbf{U}\}. \end{aligned}$$

We maintain two equivalent encodings and choose the cheaper at query time: (1) *word triples*, which list (h, ρ, t) directly for low-redundancy regimes; and (2) *compact indices*, which map E', R' to short IDs and represent query/answer/fact sequences as ID triples. The prompt payload is $\Pi = (E', R', \mathbf{q}, \mathbf{a}, \mathbf{f}, \text{rules})$, with a brief textual header explaining the ID format. Token cost scales with $|E'| + |R'| + |\mathbf{q}| + |\mathbf{a}| + |\mathbf{f}|$, typically far below concatenated passages.

2.7 Step 6: Entity-Only Consolidation

Intuition. Long-running deployments accumulate aliases, misspellings, and near-duplicates (*IBM* vs. *International Business Machines*). We consolidate *entities only*, then remap all edges and sequences.

- **Layer 1 (ANN+KNN).** Continuously build an ANN-backed k -NN graph over entity embeddings; connect pairs with cosine above a conservative threshold τ_E and form provisional alias groups.

- **Layer 2 (on-demand k -means).** When memory or $|E|$ exceeds a budget, refine groups with k -means and choose medoid representatives $m_E(\cdot)$, which minimize within-cluster distortion (Lemma 11). Each edge (u, r, v) is remapped to $(m_E(u), r, m_E(v))$, and duplicates are removed.

Because questions/answers/facts are stored as edge-ID sequences, this remap automatically cleans them as well. This quotienting can only reduce edge and sequence cardinalities (Lemma 9) and does not increase the number of sentence-level text encodings (Lemma 10).

2.8 Step 7: Adaptive Compression via DPO

The “right” selector choice depends on the query (ambiguity, hops), model (context length, robustness), domain (redundancy), and user goals (accuracy vs. latency/tokens). We learn a small categorical policy π_θ over the selector actions per channel, conditioned on features such as query length, ambiguity scores, model ID, and token budget.

Offline, for each query we evaluate several selector configurations y and compute a utility

$$U(x, y) = \alpha \cdot \text{Acc} + \delta \cdot \text{Faithfulness} - \beta \cdot \text{Tokens} - \gamma \cdot \text{Latency}.$$

We form preference pairs (x, y^+, y^-) whenever $U(x, y^+) > U(x, y^-)$, and train $\pi_\theta(y \mid x)$ with DPO against a fixed reference policy. Under a Bradley–Terry preference model, DPO aligns policy log-odds with utility differences up to a scaling and reference correction (Lemma 12), and a simple action lattice yields monotone token control under a budget constraint (Lemma 13).

At inference time, the policy picks (include all, unique, or not include) per channel, steering AutoPrunedRetriever to the appropriate operating point—e.g., overlap-heavy scaffolds for ambiguous multi-hop questions vs. aggressive dedu-

plication under tight budgets—while retaining the same symbolic infrastructure.

3 Experiments

We study:

- **RQ1 (§3.2):** complex reasoning performance on Medical, Novel, STEM, TV.
- **RQ2 (§3.3):** efficiency (tokens, latency, workspace) on STEM/TV.
- **RQ3 (§3.4):** overall performance on the full GraphRAG benchmark.

3.1 Experimental Settings

Devices and models. GraphRAG experiments (§3.4) run on an Intel i9-13900KF, 64 GB RAM, RTX 4090 (24 GB); STEM/TV experiments (§3.2, §3.3) use an A100 (40 GB). All LLM-backed parsing/judging uses the gpt-4o-mini API.

Datasets. **Medical** and **Novel** are from the GraphRAG-Benchmark (Xiang et al., 2025) with Fact Retrieval, Complex Reasoning, Contextual Summarize, and Creative Generation. We additionally construct **TV** and **STEM** from the HotpotQA Wikipedia pipeline (HotpotQA Team, 2019), grouped into 47 TV and 32 STEM micro-corpora; TV questions focus on character/episode relations, STEM on cross-sentence scientific inference.

Evaluation and parsers. We follow the LLM-judge protocol of Xiang et al. (Xiang et al., 2025). AUTOPRUNEDRETRIEVER is evaluated with two front ends: a REBEL-based triplet parser (Huguet Cabot and Navigli, 2021) and an LLM-based parser (gpt-4o-mini); both use the same pruning and indices-only prompting pipeline.

3.2 Full Complex-Reasoning Evaluation

We aggregate all complex-reasoning sources: **Medical-CR**, **Novel-CR**, **STEM**, **TV**, to test a single pruned, symbolic pipeline across technical, narrative, and pop-culture reasoning.

Quantitative results. Across all four sets, AUTOPRUNEDRETRIEVER is consistently strongest (Fig. 2). On **Medical-CR** and **Novel-CR**, the REBEL variant reaches **72.49%** and **63.02%** ACC vs. HIPBORAG2 (61.98%, 53.38%), i.e., +10.51/+9.64 points. The LLM-parser variant is close (71.59%, 62.80%), indicating that the gain mainly comes from symbolic pruning

and retrieval. On **STEM**, REBEL/LLM obtain **81.4%/78.1%** vs. HIPPO-RAG2 (69.9%); on **TV**, **68.2%/65.2%** vs. HIPPO-RAG2 (59.5%). The ordering is the same on all four: APR-REBEL > APR-llm > HippoRAG2 > LightRAG.

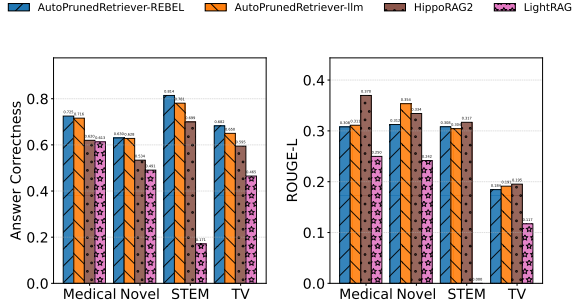


Figure 2: Average answer correctness on all complex-reasoning sets (Medical-CR, Novel-CR, STEM, TV) for HippoRAG2, LightRAG, AutoPrunedRetriever-REBEL, AutoPrunedRetriever-llm.

Case study (STEM): reasoning over ecological chains. Consider the STEM question “How does the variability in the size of brown bears across different regions serve as evidence for understanding their adaptability in various environments?” In our run this was a question where AUTOPRUNEDRETRIEVER beat both HIPPO-RAG2 and LIGHTRAG. What our retriever actually surfaced was a *compact symbolic subgraph* centered on three functional edges: (1) region → resource availability, (2) resources → body size, and (3) body size → environmental adaptability. Because those three edges were present together, the LLM could reconstruct the full causal chain:

region ⇒ food ⇒ size ⇒ adaptability.

The generated answer therefore explained that large coastal/Kodiak bears reflect high salmon (high calories), whereas smaller inland bears reflect limited resources, and that this *variation itself* is evidence of species-level adaptability. By contrast, HIPPO-RAG2 retrieved a broader, taxonomy-oriented context about brown-bear subspecies (“the taxonomy of brown bears remains somewhat bewildering”), which led the model to produce a *descriptive* answer (“there are many varieties, so sizes differ”) but not a *mechanistic* one (no resource → size link). LIGHTRAG retrieved topic-level chunks like “Brown bears vary greatly in size depending on where they live,” which was enough for correlation but not for causation. This illustrates the core advantage: our pruning keeps

only the minimal but *functional* intermediates, so the model can follow the hops in order instead of guessing them.

Case study (TV): retrieving the exact two pieces. A similar pattern appears on TV-style, entangled questions, e.g., “In The Simpsons minor-character descriptions, what in-universe line explains the Yes Guy’s stretched-out ‘Ye-e-e-s?!’, and what does the entry say about Wiseguy not actually having a single fixed name?” This question is hard not because the language is long, but because the answer lives in *two* separate mentions: one that gives the in-universe justification (“I had a stro-o-o-o-oke”) and another that clarifies the meta-labeling of the Wiseguy character. AUTOPRUNEDRETRIEVER retrieved precisely those two pieces as separate edges/nodes and presented them together, so the LLM could output both the in-universe gag *and* the meta-level note about the character not having a fixed proper name. HIPPO-RAG2, which tends to over-expand its graph, pulled a broader “recurring jokes / minor characters” context and produced a generic “it’s a running gag” answer that failed to name the stroke line. LIGHTRAG, which collapses to topic-level chunks, also stayed at the descriptive level (“recurring jokes create humor”) and missed the exact line. This shows that for entangled narrative questions, the benefit is not “more graph,” but “the *right* two edges at once.”

Overall, complex reasoning is where the pruned, indices-only pipeline shows the clearest advantage over prior GraphRAG variants.

3.3 Efficiency on Instrumented Corpora (STEM, TV)

We measure efficiency on **STEM** and **TV**, where we fully control build-time and storage logging.

Retrieval prompt tokens and latency. Figure 3 shows average query-time input tokens. AUTOPRUNEDRETRIEVER-REBEL is most compact (about 1,090 tokens on STEM and 523 on TV), followed by AUTOPRUNEDRETRIEVER-LLM (3,027 / 592). HIPPO-RAG2 and LIGHTRAG send much larger contexts (1,898/1,589 and 8,846/2,964). End-to-end latency (Fig. 4) roughly tracks this token ordering: methods with longer prompts are slower, while APR-REBEL remains competitive.

Build-time graph/prompt tokens and workspace. Figure 5 reports serialized graph-side payloads and workspace size. APR-REBEL is

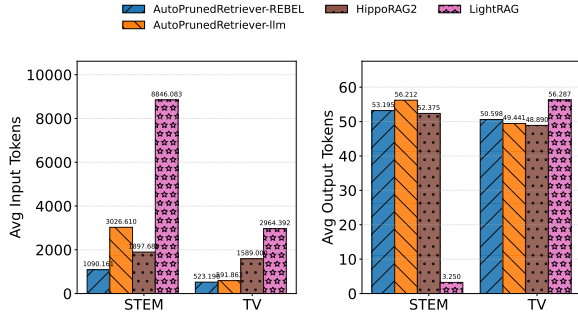


Figure 3: Input and output token usage on STEM and TV.

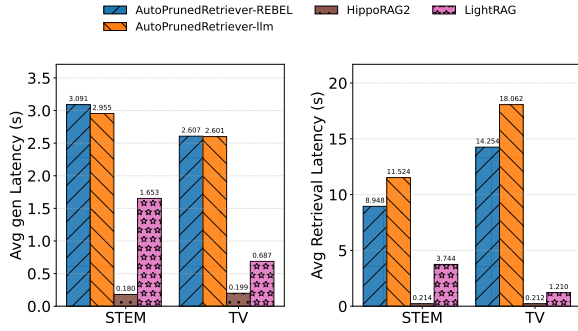


Figure 4: End-to-end latency on STEM and TV.

smallest on both corpora; APR-llm stores more LLM-extracted triples ($\approx 1.2 \times 10^6$ graph/prompt tokens on STEM and 2.84×10^6 on TV), but still below HIPPORAG2 and LIGHTRAG. Figure 6 shows that APR’s codebook / graph-size growth has plateaus where new items merge into existing entities, reflecting the two-layer entity-pruning step.

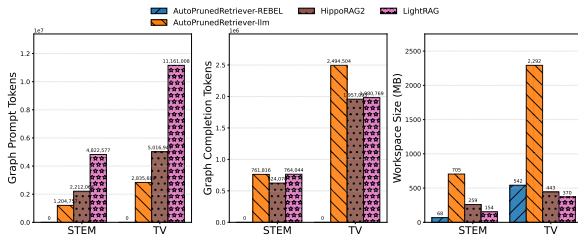


Figure 5: Build-time graph/prompt tokens and workspace usage on STEM and TV.

3.4 Full GraphRAG Benchmark

We now evaluate on the full GraphRAG benchmark (Xiang et al., 2025), i.e., **Medical** and **Novel** across Fact Retrieval, Complex Reasoning, Contextual Summarize, and Creative Generation (Table 1).

On **Contextual Summarize**, APR transfers well: on Medical it reaches **68.78%/70.14%** ACC (REBEL/LLM), slightly above FAST-GRAPH-RAG (67.88%), and on Novel it reaches

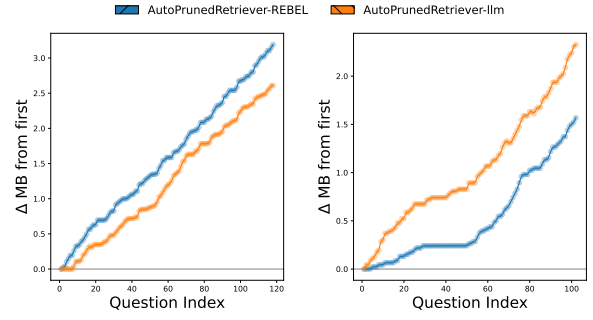


Figure 6: APR codebook / graph-size evolution (left: STEM, right: TV).

82.55%/83.10%, far above MS-GRAPH-RAG (LOCAL) (64.40%). On **Fact Retrieval**, APR-REBEL matches or exceeds strong baselines in ROUGE-L (e.g., **38.02** on Novel) while keeping ACC competitive with classic RAG and Hippo-style systems. For **Creative Generation**, APR-llm attains **62.97%** ACC on Novel and 65.02% on Medical, near or above other graph-based methods.

Token usage on GraphRAG (Fig. 7) remains low: APR-REBEL uses **1,110** tokens on Novel and **1,341** on Medical; APR-llm uses **956** and **2,243**, placing both among the most compact graph-aware systems while staying competitive or SOTA on complex reasoning and summarization.

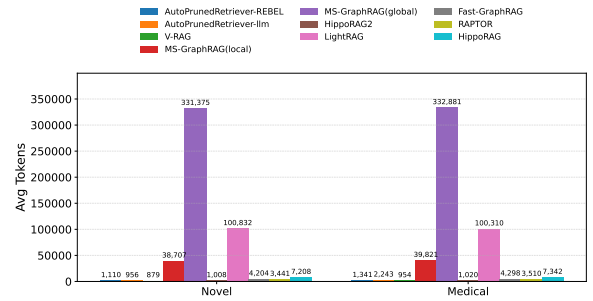


Figure 7: Average input token usage on the GraphRAG benchmark (Medical, Novel).

4 Related Work

Retrieval-Augmented Generation (RAG). RAG grounds LLMs in external knowledge via retriever-reader pipelines such as DrQA (Chen et al., 2017), DPR (Karpukhin et al., 2020b), and the original RAG model (Lewis et al., 2020). REALM integrates retrieval into pre-training with a non-parametric memory (Guu et al., 2020), while FiD performs passage-wise encoding and fusion for strong open-domain QA at higher decoding cost (Izacard and Grave, 2021). Atlas and follow-ups show that retrieval-augmented

Table 1: GraphRAG benchmark results on *Novel* and *Medical*.

Category	Model	Fact Retrieval		Complex Reasoning		Contextual Summarize		Creative Generation		
		ACC	ROUGE-L	ACC	ROUGE-L	ACC	Cov	ACC	Cov	FS Cov
Novel Dataset	RAG (w/o rerank)	58.76	<u>37.35</u>	41.35	15.12	50.08	82.53	41.52	47.46	37.84
	RAG (w/ rerank)	60.92	36.08	42.93	15.39	51.30	83.64	38.26	49.21	40.04
	MS-GraphRAG (local) (Edge et al., 2025)	49.29	26.11	50.93	24.09	64.40	75.58	39.10	55.44	35.65
	HippoRAG (Gutiérrez et al., 2025)	52.93	26.65	38.52	11.16	48.70	85.55	38.85	71.53	38.97
	HippoRAG2 (Jiménez Gutiérrez et al., 2025)	<u>60.14</u>	31.35	53.38	<u>33.42</u>	64.10	70.84	48.28	49.84	30.95
	LightRAG (Guo et al., 2025a)	58.62	35.72	49.07	24.16	48.85	63.05	23.80	57.28	25.01
	Fast-GraphRAG (CircleMind AI, 2024)	56.95	35.90	48.55	21.12	56.41	80.82	46.18	57.19	36.99
	RAPTOR (Sarathi et al., 2024)	49.25	23.74	38.59	11.66	47.10	82.33	38.01	<u>70.85</u>	35.88
	Lazy-GraphRAG (Edge et al., 2024a)	51.65	36.97	49.22	23.48	58.29	76.94	43.23	50.69	<u>39.74</u>
	AutoPrunedRetriever-REBEL	49.25	38.02	63.02	31.25	<u>82.55</u>	<u>83.95</u>	<u>59.94</u>	25.78	21.21
	AutoPrunedRetriever-llm	45.99	26.99	<u>62.80</u>	35.35	83.10	83.86	62.97	34.40	22.13
Medical Dataset	RAG (w/o rerank)	63.72	29.21	57.61	13.98	63.72	77.34	58.94	35.88	57.87
	RAG (w/ rerank)	<u>64.73</u>	30.75	58.64	15.57	65.75	78.54	60.61	36.74	<u>58.72</u>
	MS-GraphRAG (local) (Edge et al., 2025)	38.63	26.80	47.04	21.99	41.87	22.98	53.11	32.65	39.42
	HippoRAG (Gutiérrez et al., 2025)	56.14	20.95	55.87	13.57	59.86	62.73	64.43	69.21	65.56
	HippoRAG2 (Jiménez Gutiérrez et al., 2025)	66.28	<u>36.69</u>	61.98	36.97	63.08	46.13	68.05	<u>58.78</u>	51.54
	LightRAG (Guo et al., 2025a)	63.32	37.19	61.32	24.98	63.14	51.16	-	-	-
	Fast-GraphRAG (CircleMind AI, 2024)	60.93	31.04	61.73	21.37	<u>67.88</u>	52.07	<u>65.93</u>	56.07	44.73
	RAPTOR (Sarathi et al., 2024)	54.07	17.93	53.20	11.73	58.73	<u>78.28</u>	-	-	-
	Lazy-GraphRAG (Edge et al., 2024a)	60.25	31.66	47.82	22.68	57.28	55.92	62.22	30.95	43.79
	AutoPrunedRetriever-REBEL	61.28	32.96	72.49	30.79	<u>68.78</u>	40.15	64.04	32.19	11.12
	AutoPrunedRetriever-llm	61.25	34.69	<u>71.59</u>	<u>31.11</u>	70.14	40.59	65.02	33.06	28.62

LMs with updatable indices and reranking can reach strong few-shot performance (Izacard et al., 2022), but these text-first systems still treat queries independently and remain passage-centric and token-heavy for multi-hop or long-context reasoning.

Vector Retrieval and Re-ranking. Dense retrieval embeds queries and documents into a shared space (Lee et al., 2019; Karpukhin et al., 2020b) and is typically served by ANN indexes such as FAISS (Johnson et al., 2017) or HNSW (Malkov and Yashunin, 2016). Two-stage pipelines then apply stronger rerankers: BERT-based cross-encoders (Nogueira and Cho, 2019) and late-interaction models like ColBERT/ColBERTv2 (Khattab and Zaharia, 2020; Santhanam et al., 2022) improve ranking quality while trading off latency and indexability.

Graph-based RAG (GraphRAG). Graph-based RAG replaces flat passages with entity–relation structure. MS-GraphRAG builds graphs over private corpora and retrieves summarized subgraphs, with a local variant that builds per-session graphs (Edge et al., 2024c). Lazy-GraphRAG defers graph construction to query time (Edge et al., 2024b). HippoRAG and HippoRAG2 introduce neuro-inspired long-term memory that consolidates and reuses reasoning

paths (Gutiérrez et al., 2025a,b). LightRAG and Fast-GraphRAG simplify graphs into key–value forms for low-latency retrieval (Guo et al., 2025b; CircleMind AI, 2024), while RAPTOR uses a tree of recursive summaries instead of an explicit graph (Sarathi et al., 2024). Together, these illustrate the shift from passage-centric to structured retrieval, but still pay substantial graph-serialization and token cost.

Memory-Augmented Architectures.

External-memory LMs maintain persistent, updatable stores alongside parameters, exploring scalable memory (Wu et al., 2022), efficient retrieval (Liu et al., 2024), and continual consolidation (Dao et al., 2023). These build on differentiable memory architectures such as memory networks (Weston et al., 2014) and neural Turing machines (Graves et al., 2014).

Efficient Model Deployment. LLM efficiency is improved by quantization and low-rank adaptation (e.g., QLoRA) (Dettmers et al., 2023), mixed-precision training (Micikevicius et al., 2018), and hardware-aware optimization (Ai et al., 2024). Adaptive computation and dynamic routing (Schuster et al., 2022; Li et al., 2023) further trade depth and routing complexity against accuracy, complementing retrieval- and memory-side efforts to cut tokens and latency.

5 Limitations

Our study has several limitations. First, we evaluate AUTOPRUNEDRETRIEVER primarily on English, knowledge-intensive QA benchmarks (GraphRAG-Benchmark, STEM, TV); it is unclear how well the same pruning and symbolization scheme transfers to other languages, domains, or noisy user logs. Second, our pipeline depends on upstream triple extractors (REBEL or an LLM); systematic extraction errors or missing relations can still harm downstream reasoning, and we do not jointly train extraction and retrieval. Finally, we focus on text-only corpora and single-turn question answering in agentic settings, leaving multimodal inputs, tool-use workflows, and human-in-the-loop updates to future work.

References

- Qinbin Ai, Mingxuan Chen, Wenhui Chen, Tri Dao, Daniel Fu, Albert Gu, and 1 others. 2024. Efficient inference for large language models: A survey. *Foundations and Trends in Machine Learning*, 17(2):145–387.
- Jinheon Baek, Myeongho Jeong, Sung Ju Hwang, and Jungwoo Park. 2023. Knowledge graph grounded question answering with graph neural networks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, EMNLP ’23.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading wikipedia to answer open-domain questions. In *ACL*.
- CircleMind AI. 2024. [Fast graphrag: High-speed graph-based retrieval-augmented generation](#). Analytics Vidhya Blog.
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2023. Hungry hungry hippos: Towards language modeling with state space models. In *International Conference on Learning Representations*, ICLR ’23.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. In *Advances in Neural Information Processing Systems*, volume 36 of *NeurIPS ’23*.
- 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](#). *Preprint*, arXiv:2404.16130.
- Darren Edge, Ha Trinh, and Jonathan Larson. 2024a. [LazyGraphRAG: Setting a New Standard for Quality and Cost](#). Microsoft Research Blog.
- Darren Edge, Ha Trinh, Jonathan Larson, Alex Chao, and Robert Osazuwa Ness. 2024c. [From local to global: A graph rag approach to query-focused summarization](#). ArXiv:2404.16130 [cs.CL].
- Alex Graves, Greg Wayne, and Ivo Danihelka. 2014. [Neural Turing machines](#). *Preprint*, arXiv:1410.5401.
- Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. 2025a. [Lightrag: Simple and fast retrieval-augmented generation](#). *Preprint*, arXiv:2410.05779.
- Zirui Guo, Lianghao Xia, Yanhua Yu, Tu Ao, and Chao Huang. 2025b. [Lightrag: Simple and fast retrieval-augmented generation](#). ArXiv:2410.05779 [cs.IR].
- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2025a. [HippoRAG: Neurobiologically Inspired Long-Term Memory for Large Language Models](#). ArXiv:2405.14831 [cs.CL].
- 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:2502.14802 [cs.CL].
- Bernal Jiménez Gutiérrez, Yiheng Shu, Yu Gu, Michihiro Yasunaga, and Yu Su. 2025. [Hipporag: Neurobiologically inspired long-term memory for large language models](#). *Preprint*, arXiv:2405.14831.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Papat, and Ming-Wei Chang. 2020. REALM: Retrieval-augmented language model pre-training. In *ICML*.
- Haoyu Han, Yue Wang, Hagai Shomer, Kai Guo, Jiawei Ding, Yang Lei, Bo Li, and Jie Tang. 2024. [Retrieval-augmented generation with graphs \(graphrag\)](#). *Preprint*, arXiv:2404.16130.
- HotpotQA Team. 2019. Preprocessed wikipedia for hotpotqa. <https://hotpotqa.github.io/wiki-readme.html>.
- Pere-Lluís Hugué Cabot and Roberto Navigli. 2021. [REBEL: Relation extraction by end-to-end language generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 2370–2381, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Gautier Izacard and Edouard Grave. 2021. Leveraging passage retrieval with generative models for open-domain question answering. In *ICLR*.
- 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. *Journal of Machine Learning Research*, 23(251):1–43.

- Bernal Jiménez Gutiérrez, Yiheng Shu, Weijian Qi, Sizhe Zhou, and Yu Su. 2025. [From rag to memory: Non-parametric continual learning for large language models](#). *arXiv preprint arXiv:2502.14802*. Poster at ICML 2025.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with GPUs. *arXiv:1702.08734*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020a. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020b. [Dense passage retrieval for open-domain question answering](#). *arXiv preprint arXiv:2004.04906*.
- Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and effective passage search via contextualized late interaction over BERT. In *SIGIR*.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent retrieval for weakly supervised open-domain question answering. In *ACL*.
- 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. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33 of *NeurIPS '20*, pages 9459–9474.
- Yuhong Li, Daniel Fu, Wenhui Chen, Manoj Kumar, Noah Smith, Ming-Wei Chen, and Christopher Ré. 2023. Efficient streaming language models with attention sinks. In *Proceedings of the 40th International Conference on Machine Learning, ICML '23*.
- Zhenghao Liu, Jianfeng Wang, Jie Tang, Jie Zhou, and Ashwin Kalyan. 2024. Memory-efficient large language models via dynamic memory allocation. In *International Conference on Learning Representations, ICLR '24*.
- Yury A. Malkov and Dmitry A. Yashunin. 2016. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *arXiv:1603.09320*.
- Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. 2018. Mixed precision training. In *International Conference on Learning Representations, ICLR '18*.
- Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage re-ranking with BERT. In *arXiv:1901.04085*.
- Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. *Transactions of the Association for Computational Linguistics*, 11:1316–1331.
- Keshav Santhanam, Omar Khattab, and et al. 2022. Colbertv2: Effective and efficient retrieval via lightweight late interaction. In *NAACL*.
- Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. 2024. [Raptor: Recursive abstractive processing for tree-organized retrieval](#). In *Proceedings of the International Conference on Learning Representations (ICLR)*. ArXiv:2401.18059 [cs.CL].
- Tal Schuster, Adam Fisch, Tommi Jaakkola, and Regina Barzilay. 2022. Confident adaptive language modeling. In *Advances in Neural Information Processing Systems*, volume 35 of *NeurIPS '22*, pages 17456–17472.
- Zeyu Sun, Hongyang Yang, Rui Zhou, Chen Wang, Xiaodong Liu, and Xuanjing Huang. 2022. Graphene: Retrieval-augmented generation for efficient knowledge-intensive reasoning. In *Advances in Neural Information Processing Systems*, volume 35 of *NeurIPS '22*.
- Xiaorui Wang, Zhiyuan Zhang, Yizheng Hao, Lei Li, Lei Hou, Zhiyuan Liu, Xuan Song, and Maosong Sun. 2023. Graph-based memory for large language models. In *Proceedings of the 40th International Conference on Machine Learning, ICML '23*.
- Jason Weston, Sumit Chopra, and Antoine Bordes. 2014. Memory networks. In *International Conference on Learning Representations, ICLR '15*.
- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. 2024. [Autogen: Enabling next-gen LLM applications via multi-agent conversation](#).
- Yuhuai Wu, Markus Rabe, DeLesley Hutchins, and Christian Szegedy. 2022. Memorizing transformer. In *International Conference on Learning Representations, ICLR '22*.
- 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](#). *Preprint*, arXiv:2506.05690.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2023. [React: Synergizing reasoning and acting in language models](#). *Preprint*, arXiv:2210.03629.

A appendix

A.1 Theoretical Properties

A.1.1 Encoding and Heavy-Tailed Repetition

Lemma 1 (Concentration of repetition). *Let \mathcal{V} be the corpus vocabulary and $f : \mathcal{V} \rightarrow \mathbb{N}$ token frequencies. In typical language corpora f is heavy-tailed: there is a core set $\mathcal{V}_{\text{core}} \subset \mathcal{V}$ such that*

$$\sum_{v \in \mathcal{V}_{\text{core}}} f(v) \approx \Theta(|\mathcal{C}|),$$

where $|\mathcal{C}|$ is the total token count. Thus, indexing recurrent entities/relations captures most informational mass while reducing redundancy.

Sketch. Empirical token distributions in large corpora follow Zipf-like laws; a small set of types accounts for most occurrences. Mapping these to IDs and sharing them across queries/facts removes repeated surface forms without losing the dominant mass. \square

A.1.2 Chunked Small Graphs (Local-First Construction)

Lemma 2 (Maximal local-coherence partition). *Fix a threshold τ , a bounded continuity bonus $0 \leq b < \infty$, and a fit rule that is monotone nonincreasing as the centroid drifts away from a triple. The one-pass rule “append if fit-score $\geq \tau$, else cut” yields a partition $\mathcal{G} = (G_1, \dots, G_K)$ in which each G_k is maximal: no additional triple can be appended without violating the fit test.*

Sketch. Within a segment, appending triples can only decrease (or leave unchanged) the fit of future triples because the centroid moves away from any fixed candidate. Once a triple fails the test, any larger graph containing it would also fail. Thus, each cut point defines a maximal prefix w.r.t. the local rule. \square

Lemma 3 (Boundary-consistency merge). *Let Θ denote the segmenter parameters (threshold, continuity bonus, etc.). For adjacent runs (L, R) , define $\text{MERGE}(L, R)$ as true iff the segmenter applied to the concatenation $L \parallel R$ either (i) produces a single chunk or (ii) places its first cut away from the original boundary at $|L|$. Then a boundary is removed iff the segmenter, when given both sides at once, would not cut at that location. Surviving boundaries are fixed points of the segmenter under local re-evaluation.*

Sketch. The merge test re-runs exactly the same algorithm and hyperparameters on the local window. If the “true” segmentation prefers a different cut, we merge; otherwise we keep the boundary. This is equivalent to requiring boundary self-consistency under the same rule. \square

Lemma 4 (Intra-chunk cohesion bound). *Assume unit-normalized triple embeddings (v_i) and an acceptance rule $\cos(\bar{c}, v_i) + \delta_i \geq \tau$ with $0 \leq \delta_i \leq b$, where \bar{c} is the running centroid. For any completed small graph G_k with $|G_k| \geq 2$,*

$$\frac{2}{|G_k|(|G_k|-1)} \sum_{p < q} \cos(v_p, v_q) \geq \tau - b.$$

Sketch. The centroid always lies in the convex hull of the triple embeddings. The acceptance condition ensures each new triple is not too far (in cosine) from the current centroid, up to the continuity bonus b . Averaging pairwise cosines and using triangle-type inequalities yields the bound. \square

Lemma 5 (Working-set reduction for retrieval). *Let $M = |\mathcal{M}|$ be the number of distinct edges, and suppose runs have lengths ℓ_1, \dots, ℓ_K with $L = \max_k \ell_k$. If symbolic pre-filtering selects H candidate runs for re-ranking, then the fine stage touches at most $H \cdot L$ edges. Compared to scanning all edges, this yields an asymptotic reduction factor $\Omega(M/(H \cdot L))$.*

Sketch. Each candidate run contributes at most L edges to be scored. Bounding H and L independently of M (corpus size) gives sublinear dependence on M . \square

Corollary 6 (Precision–recall / segmentation tradeoff). *Increasing τ (or decreasing the continuity bonus b) shortens runs, improves cohesion, and reduces retrieval latency at the cost of recall. Decreasing τ lengthens runs, improves recall, but increases candidate sizes. The boundary merge step counteracts over-segmentation by removing unstable cuts.*

Sketch. Higher thresholds cause earlier cuts; lower thresholds allow more heterogeneous content in each run. The merge rule prunes cuts that the full-window segmenter would not reproduce. \square

A.1.3 Coarse retrieve

Lemma 7 (Efficiency of Coarse→Fine with Max–Pair Filtering). *Let n be corpus size, d the embedding dimension, and $k \ll n$ the coarse shortlist size. The coarse stage computes, per query–candidate pair, constants over small entity/relation sets; the fine stage evaluates only k items with sentence embeddings. Thus cost drops from $O(n \cdot d)$ to $O(k \cdot d)$ while preserving precision provided k retains high-overlap candidates.*

A.1.4 Semantic Selection and Consensus

Lemma 8 (Semantic consensus is close to its members). *Let runs r_1, \dots, r_m have embeddings $\psi(r_i)$ and cosine similarity $\text{sim}(r_a, r_b) = \cos(\psi(r_a), \psi(r_b))$. Let \bar{r} maximize the average similarity*

$$\bar{r} \in \arg \max_r \sum_{i=1}^m \text{sim}(r, r_i).$$

If all runs in the cluster are mutually similar, i.e., $\text{sim}(r_a, r_b) \geq \theta$ for some θ close to 1, then \bar{r} is also close to every member:

$$1 - \text{sim}(r_a, \bar{r}) \leq \varepsilon(\theta) \quad \text{for all } a,$$

for a function $\varepsilon(\theta) \rightarrow 0$ as $\theta \rightarrow 1$.

Sketch. On the unit sphere, cosine similarity defines a bounded distance. If \bar{r} were far from some member r_a while all runs are mutually close, moving \bar{r} toward r_a would increase its average similarity to the cluster, contradicting maximality. The bound $\varepsilon(\theta)$ follows from standard geometric arguments. \square

A.1.5 Consolidation of Entities and Edge Sequences

Lemma 9 (Quotient consolidation reduces edge cardinality). *Let $G = (V, R, E)$ be a directed multigraph with labeled edges $E \subseteq V \times R \times V$. Let \sim be the equivalence relation on V induced by entity consolidation and $\pi : V \rightarrow V/\sim$ the projection. Define $\phi : E \rightarrow E'$ by $\phi(u, r, v) = (\pi(u), r, \pi(v))$ and let $E' = \text{uniq}(\phi(E))$. Then:*

1. $|E'| \leq |E|$, with equality iff ϕ is injective on E .
2. For any edge sequence $\sigma = (e_{i_1}, \dots, e_{i_T})$, the remapped-and-deduped sequence $\sigma' = \text{uniq}(\phi(\sigma))$ satisfies $|\sigma'| \leq |\sigma|$.

Sketch. E' is the image of E under ϕ followed by deduplication, so it cannot have more elements. Sequences inherit this property pointwise; collapsing equal edges cannot increase length. \square

Lemma 10 (Vectorization cost is preserved). *Let v_s be a schema vectorizer depending only on symbolic indices and precomputed vectors ($E_{\text{emb}}, R_{\text{emb}}$). Then the number of sentence-level text encodings required by the pipeline is unchanged by applying ϕ and deduplicating edges.*

Sketch. Consolidation changes only which indices point to which vectors. All sentence encodings are done before consolidation (for triples and chunks), so later index manipulation reuses them. \square

Lemma 11 (Medoid representatives minimize within-cluster distortion). *Let $C \subseteq V$ be a cluster of entities and define cosine dissimilarity $d(\mathbf{x}, \mathbf{y}) = 1 - \cos(\mathbf{x}, \mathbf{y})$. A medoid $r^* \in C$ satisfies*

$$r^* \in \arg \min_{r \in C} \sum_{i \in C} d(\mathbf{e}_i, \mathbf{e}_r),$$

and for any other representative $r \in C$,

$$\sum_{i \in C} d(\mathbf{e}_i, \mathbf{e}_{r^*}) \leq \sum_{i \in C} d(\mathbf{e}_i, \mathbf{e}_r).$$

Sketch. This is the standard property of medoids: by definition they minimize total dissimilarity within the cluster. \square

A.1.6 DPO Wrapper and Policy Behavior

Lemma 12 (DPO aligns policy log-odds with utilities). *Assume preferences follow a Bradley–Terry model:*

$$\Pr(y^+ \succ y^- \mid x) = \sigma(\lambda [U(x, y^+) - U(x, y^-)])$$

for some $\lambda > 0$, where U is a latent utility. Let π_{ref} be fixed. Then any minimizer π_θ of the DPO objective satisfies, up to a normalization constant C_x ,

$$\begin{aligned} \log \pi_\theta(y \mid x) - \log \pi_\theta(y' \mid x) \\ = \frac{\lambda}{\beta_{\text{dpo}}} [U(x, y) - U(x, y')] \\ + \Delta_{\text{ref}}(y, y' \mid x) + C_x. \end{aligned}$$

where Δ_{ref} depends only on π_{ref} .

Sketch. DPO maximizes the conditional log-likelihood of observed pairwise preferences under

a logistic link, with a reference correction. At optimum, gradient stationarity enforces proportionality between policy log-odds and utility differences, offset by the fixed reference. \square

Lemma 13 (Monotone token control via action lattice). *Suppose each channel’s action set $\{\text{include}, \text{all}, \text{unique}, \text{not include}\}$ forms a lattice under \succeq with*

$$\text{include} \succeq \text{all} \succeq \text{unique} \succeq \text{not include},$$

such that $\text{Tokens}(x, y)$ is monotone nonincreasing down the lattice and $\text{Acc}(x, y)$ is Lipschitz in a task metric. Then there exists a Lagrange multiplier $\eta^ \geq 0$ such that the DPO-trained policy with penalty $\eta^* \cdot \text{Tokens}$ attains a target budget B , and tighter budgets $B' < B$ can be achieved by increasing η (eventually collapsing to always not include).*

Sketch. On a finite action set, the mixed policy over actions yields a convex set of achievable (tokens, accuracy) pairs. A standard Lagrangian argument with a monotonically ordered action set gives existence of a multiplier realizing each feasible budget, and increasing the penalty pushes mass toward cheaper actions. \square

A.2 Effect of Token Length and Benefits of Entity–Relation Factorization

Sequence embeddings. For simplicity, we approximate the embedding of a text span S of length n tokens by the mean of its token embeddings:

$$z(S) \approx \frac{1}{n} \sum_{i=1}^n x_i \in \mathbb{R}^d.$$

Each token embedding decomposes into

$$x_i = s_i + \varepsilon_i,$$

where s_i is the semantic signal and ε_i is zero-mean noise with $\mathbb{E}[\varepsilon_i] = 0$ and $\text{Var}(\langle q, \varepsilon_i \rangle) = \sigma^2$ for any unit query vector q .

We partition the tokens into: (i) relevant tokens R (carry information the query cares about) and (ii) irrelevant tokens I (background, boilerplate, narration), with $|R| = m$ and $|I| = n - m$.

Lemma 14 (Token dilution). *Let q be a unit query vector aligned with the average relevant signal*

$$\mu_{\text{rel}} := \frac{1}{m} \sum_{i \in R} s_i \quad \text{with} \quad \langle q, \mu_{\text{rel}} \rangle = \alpha > 0.$$

Assume irrelevant tokens have no systematic alignment with q , i.e., $\mathbb{E}[\langle q, s_i \rangle] = 0$ for $i \in I$. Then the signal-to-noise ratio (SNR) of the passage embedding in the direction q decays as

$$\text{SNR}(S) := \frac{(\mathbb{E}[\langle q, z(S) \rangle])^2}{\text{Var}(\langle q, z(S) \rangle)} \propto \frac{1}{n}.$$

Proof. We have

$$z(S) = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} \sum_{i \in R} (s_i + \varepsilon_i) + \frac{1}{n} \sum_{i \in I} (s_i + \varepsilon_i).$$

Taking inner product with q and expectation, and using $\mathbb{E}[\langle q, \varepsilon_i \rangle] = 0$ for all i and $\mathbb{E}[\langle q, s_i \rangle] = 0$ for $i \in I$,

$$\mathbb{E}[\langle q, z(S) \rangle] = \frac{1}{n} \sum_{i \in R} \langle q, s_i \rangle = \frac{m}{n} \alpha.$$

Thus, for fixed m and α , the expected signal scales as $m\alpha/n$.

For the variance, by independence and identical variance of the noise:

$$\begin{aligned} \text{Var}(\langle q, z(S) \rangle) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n \langle q, \varepsilon_i \rangle\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}(\langle q, \varepsilon_i \rangle) = \frac{\sigma^2}{n}. \end{aligned}$$

Therefore

$$\text{SNR}(S) = \frac{(m\alpha/n)^2}{\sigma^2/n} = \frac{m^2 \alpha^2}{\sigma^2} \cdot \frac{1}{n},$$

which shows $\text{SNR}(S) \propto 1/n$ as claimed. \square

Corollary 15 (Length bias). *Consider two passages S_{short} and S_{long} that contain the same m relevant tokens (same fact, same α) but have lengths n_s and n_ℓ with $n_\ell > n_s$. Then*

$$\mathbb{E}[\langle q, z(S_{\text{short}}) \rangle] = \frac{m}{n_s} \alpha > \frac{m}{n_\ell} \alpha = \mathbb{E}[\langle q, z(S_{\text{long}}) \rangle]$$

Thus, even for equally relevant content, longer passages tend to produce lower expected similarity to q and are systematically disadvantaged in retrieval.

Implication. Lemma 14 and Corollary 15 formalize a “token dilution” effect: when we embed entire passages, the representation of a fact is weakened by irrelevant tokens, and the SNR decreases as $1/n$ with passage length. Consequently,

retrieval quality depends not only on what is said, but also on how long and where it is written.

Entity-relation factorization. In our system, we instead represent knowledge as graph edges $(e, r, e') \in \mathcal{M} \subseteq E \times R \times E$ and store embeddings for entities and relations via the codebook

$$E_{\text{emb}} : E \rightarrow \mathbb{R}^d, \quad R_{\text{emb}} : R \rightarrow \mathbb{R}^d.$$

Thus each fact $f = (e, r, e')$ is encoded by the triple

$$E_{\text{emb}}(e), R_{\text{emb}}(r), E_{\text{emb}}(e'),$$

derived from short token sequences for entity names and relation labels, rather than from whole passages.

Proposition 16 (Advantages of E–R–E embeddings). *Let $E, R, \mathcal{M}, E_{\text{emb}}, R_{\text{emb}}$ be as in Section 2.2. Under the averaging+noise model of Lemma 14, the following hold:*

1. (High-SNR micro-embeddings) *There exist constants $c_1, c_2 > 0$, independent of the passage length n , such that for any fact $f = (e, r, e') \in \mathcal{M}$,*

$$c_1 \leq \text{SNR}(v) \leq c_2$$

$$\forall v \in \{E_{\text{emb}}(e), R_{\text{emb}}(r), E_{\text{emb}}(e')\}.$$

In particular, E–R–E embeddings do not suffer the $1/n$ decay of Lemma 14.

2. (Compositional query scoring) *Let a query q induce components $q_E, q_R \in \mathbb{R}^d$. For suitable nonnegative weights $\lambda_s, \lambda_r, \lambda_t$ we can score an edge $(e, r, e') \in \mathcal{M}$ via*

$$\text{score}(e, r, e' \mid q) = \lambda_s \langle q_E, E_{\text{emb}}(e) \rangle +$$

$$\lambda_r \langle q_R, R_{\text{emb}}(r) \rangle + \lambda_t \langle q_E, E_{\text{emb}}(e') \rangle,$$

i.e., as an inner product between a structured query and short, high-SNR E–R–E embeddings, instead of a single inner product with a noisy passage embedding $z(S)$.

3. (Localized interference and updates) *Each fact f has its own edge (e, r, e') ; interference between facts arises only through shared $E_{\text{emb}}(\cdot)$ or $R_{\text{emb}}(\cdot)$. Updating a single fact changes $O(1)$ vectors instead of re-embedding entire passages containing many unrelated facts.*

Proof sketch. (1) Let the surface string for $e \in E$ have n_E tokens, of which m_E are semantically relevant. By construction n_E is bounded by a small constant (e.g., 1–3), so $m_E \approx n_E = O(1)$. Applying the same calculation as in Lemma 14 with $n = n_E$ gives

$$\text{SNR}(E_{\text{emb}}(e)) \propto \frac{m_E^2}{n_E} = \Theta(1),$$

with constants independent of the passage length n in which f appears. The same argument applies to $R_{\text{emb}}(r)$ and $E_{\text{emb}}(e')$, yielding the claimed bounds c_1, c_2 .

(2) Because we store $E_{\text{emb}}(e), R_{\text{emb}}(r)$, and $E_{\text{emb}}(e')$ separately, any query q that decomposes into components (q_E, q_R) admits the factorized score above. Algebraically, this is a weighted sum of inner products between short E–R–E vectors and corresponding query components, rather than a single inner product $\langle q, z(S) \rangle$ with a length-dependent passage embedding.

(3) The representation of a fact f is the triple $(E_{\text{emb}}(e), R_{\text{emb}}(r), E_{\text{emb}}(e'))$. Adding, removing, or modifying f only affects these embeddings and other edges sharing e, r , or e' . No re-embedding of unrelated tokens is required, in contrast to passage-level embeddings that entangle many facts in the same $z(S)$. \square

A.3 Retrieval Details

For completeness we summarize the exact scoring terms used in the coarse and fine stages.

Coarse score. An indexed run \mathbf{y} decodes to triples $S(\mathbf{y}) = \{(h, \rho, t)\} \subseteq \mathcal{M}$. We collect entity and relation embeddings into matrices $E(\mathbf{y}) \in \mathbb{R}^{n_e \times d}$ and $R(\mathbf{y}) \in \mathbb{R}^{n_r \times d}$. For a query q and candidate run f ,

$$s_{\text{coarse}}(q, f) = w_{\text{ent}} \max_{i,j} \cos(E(q)_i, E(f)_j) \\ + w_{\text{rel}} \max_{p,r} \cos(R(q)_p, R(f)_r),$$

and we take $I_k = \text{TopK}_f s_{\text{coarse}}(q, f)$.

Fine score from triple lines. For each candidate $f \in I_k$, we linearize triples to short lines “ $h \rho t$ ”, embed query and candidate lines into $\mathbf{Q} \in \mathbb{R}^{n_q \times d}$ and $\mathbf{C} \in \mathbb{R}^{n_c \times d}$, and form the cosine matrix

$$\mathbf{S} = \widehat{\mathbf{Q}} \widehat{\mathbf{C}}^\top \in [-1, 1]^{n_q \times n_c}.$$

All fine-stage terms are computed on \mathbf{S} :

- **RelTopT**: flatten S , take the top- t entries, and average.
- **Coverage**: $\text{Cov}(\tau_{\text{cov}}) = \sum_i \mathbf{1}[\max_j S_{ij} \geq \tau_{\text{cov}}]$.
- **Many-to-many (MP)**: apply $\sigma((S_{ij} - \tau_{\text{pair}})/T_{\text{pair}})$ elementwise and normalize by $\sqrt{n_q n_c}$ or $\log(1 + n_q n_c)$.
- **Distinct 1:1**: greedily select the largest unused entries above τ_{dist} and average them with a $1/\sqrt{m}$ factor.
- **Whole-chunk gate**: compute a full-chunk cosine between concatenated query and candidate text, normalize by a length term, and gate with a sigmoid so very long but off-topic chunks do not get extra credit.

The final semantic score is a weighted sum

$$s_{\text{fine}} = \text{RelTopT} + \lambda_{\text{cov}} \text{Cov} + \lambda_{\text{mp}} \text{MP} + \lambda_{1:1} \text{Distinct} + \lambda_{\text{whole}} \text{WholeGate}.$$

with one set of weights and thresholds per dataset/model, reused across all experiments.

A.4 Prompt Format and Input Encoding

We encode each input as a compact, graph-structured prompt rather than a long text block. Concretely, a prompt consists of:

- A codebook of entities and relations, E' and R' (either as words or short IDs).
- Edge sequences for the query, prior knowledge, and facts: query edges q , knowledge edges k , and fact edges f .
- A short instruction block describing how to interpret each tuple (h, ρ, t) or ID triple.

Figure 8 contrasts a conventional raw-text prompt with our ID-based and word-based encodings. The ID variants use a JSON-style schema:

One of three formats depends on which one has fewer tokens.

(a) Raw-text prompt

Q: Which subsidiaries acquired since 2021 are exposed to new EU rules?
Context:
- In 2022, AlphaCorp acquired BetaLtd...
- EU Regulation 2024/12 applies to...
- Post-merger reports indicate ...
(plus additional retrieved passages ...)

(b) ID-referenced codebook with edge matrix

```
{
  "e": ["AlphaCorp", "BetaLtd",
        "EUReg2024_12", "2021+"],
  "r": ["acquired_in", "exposed_to",
        "subject_to"],
  "edge_matrix": [[0,0,3],
                  [1,1,2],
                  [0,2,2]],
  "questions(edges[i])": [0,1],
  "facts(edges[i])": [0,2],
  "rules": "<KB schema string>"
}
```

(c) ID-referenced compact triples

```
{
  "e": ["AlphaCorp", "BetaLtd",
        "EUReg2024_12", "2021+"],
  "r": ["acquired_in", "exposed_to",
        "subject_to"],
  "questions([e,r,e], ...)": [[0,0,3], [1,1,2]],
  "facts([e,r,e], ...)": [[0,0,3], [0,2,2]],
  "rules": "<KB schema string>"
}
```

(d) Word-level triples (no IDs)

```
{
  "questions(words)": [[AlphaCorp, acquired_in, 2021+],
                        [BetaLtd, exposed_to, EUReg2024_12]],
  "facts(words)": [[AlphaCorp, acquired_in, 2021+],
                    [AlphaCorp, subject_to, EUReg2024_12]],
  "rules": "<KB schema string>"
}
```

Figure 8: **AutoPrunedRetriever input encodings.** Panel (a) shows a conventional long-context prompt; (b) encodes the same information via an entity/relation codebook and an edge matrix; (c) uses explicit triple lists with IDs; (d) uses full-word triples.

(a) Edge-matrix JSON schema

```
---Knowledge Base---
[JSON format]
- e: list of entities (e[i] = entity string)
- r: list of relations (r[j] = relation string)
- edge_matrix: [[head_e_idx, r_idx, tail_e_idx]]
  * NOTE: edges[i] is just shorthand for edge_matrix[i]
- questions(edges[i]): questions linked by edge i
- given knowledge(edges[i]): prior answers linked by edge i
- facts(edges[i]): facts linked by edge i
```

(b) ID-based triple JSON schema

```
---Knowledge Base---
[JSON format]
- e: list of entities (e[i] = entity string)
- r: list of relations (r[j] = relation string)
- [e,r,e]: triple [head_e_idx, r_idx, tail_e_idx]
- questions([e,r,e], ...): question triples
- given knowledge([e,r,e], ...): prior answer triples
- facts([e,r,e], ...): fact triples
```

(c) Word-level triple schema

```
---Knowledge Base---
[JSON format]
- questions(words): question triples
- given knowledge(words): prior answer triples
- facts(words): fact triples
```

Figure 9: **Knowledge-base JSON specifications (“rules”) used by AutoPrunedRetriever.** The concrete encodings in Fig. 8 all instantiate one of these schemas.

- e: entity vocabulary (either strings or IDs).

- r : relation vocabulary.
- edge matrix or triple lists: $[h, r, t]$ indices into e/r or E/R .
- questions, knowledge, facts: subsets of edges tagged as questions, prior knowledge, or background facts.

A.5 Cross-domain qualitative comparison on STEM and TV

Domain / ID	Question (abridged)	AutoPrunedRetriever Analysis (abridged)	Result	HippoRAG2 Analysis (abridged)	Result	LightRAG Analysis (abridged)	Result	Analysis
STEM-5c755e96	Brown bear size variation and adaptability	<i>Context:</i> “Kodiak bears are largest due to high salmon availability; inland bears smaller with limited resources.” <i>Reasoning:</i> geography → food abundance → body-mass shift → adaptability. <i>Answer:</i> Larger coastal/Kodiak bears reflect rich caloric intake; smaller inland bears reflect scarcity ⇒ size variance evidences environmental adaptability. <i>Error:</i> — (correct)		<i>Context:</i> “The taxonomy of brown bears remains bewildering; multiple subspecies identified.” <i>Reasoning:</i> taxonomy → morphological variation (no environmental cause). <i>Answer:</i> Size variability indicates subspecies diversity. <i>Error:</i> Misses causal driver (resources) ⇒ descriptive, not mechanistic.		<i>Context:</i> “Brown bears vary greatly in size depending on where they live.” <i>Reasoning:</i> region → size → adaptation (shallow). <i>Answer:</i> Bears adapt to local conditions, so sizes differ. <i>Error:</i> Correlation only; lacks resource/metabolic link.		
STEM-4c26ae6d	Historical range → ecological role	<i>Context:</i> “Mexican grizzly / Kodiak / Himalayan subspecies; apex predators affecting vegetation and prey.” <i>Reasoning:</i> historical range → diversification → habitat adaptation → modern trophic role. <i>Answer:</i> Past range shaped regional lineages whose adaptations underwrite today’s grizzly apex role. <i>Error:</i> — (correct)		<i>Context:</i> “Pleistocene lineage prior to demise.” <i>Reasoning:</i> lineage timeline → extinction (no present ecology). <i>Answer:</i> Historical divergence explains current bears (vague). <i>Error:</i> Lacks link to present-day ecological function.		<i>Context:</i> “Ecological dynamics of predator–prey systems.” <i>Reasoning:</i> ecosystem complexity → generic role. <i>Answer:</i> Grizzlies play roles in ecosystems (generic). <i>Error:</i> No entity-level or causal path from range to role.		
STEM-1b8f5662	Physical adaptations → hunting success (moose)	<i>Context:</i> “Charge and scent-based ambush tactics; terrain affects prey choice.” <i>Reasoning:</i> morphology + habitat → tactic → success vs large prey. <i>Answer:</i> Bears’ strength/claws + terrain-leveraged tactics raise success on moose. <i>Error:</i> — (correct)		<i>Context:</i> “Brown bears as apex omnivores in ecosystems.” <i>Reasoning:</i> apex predator → survival (no tactics). <i>Answer:</i> As apex predators they can hunt large prey. <i>Error:</i> Omits behavioral mechanism/tactical link.		<i>Context:</i> “Bears interact with diverse ecosystems.” <i>Reasoning:</i> environment → adaptation (broad). <i>Answer:</i> Environmental adaptation enables hunting. <i>Error:</i> High-level summary; no tactic/terrain edge.		
TV-29d2f5b1	Caboose and Omega possession (Red vs Blue)	<i>Context:</i> “Caboose’s abnormal behavior linked to Omega possession and oxygen deprivation after suit reboot.” <i>Reasoning:</i> possession + hypoxia → erratic acts → friendly-fire. <i>Answer:</i> Caboose; accidents stem from AI control + hypoxia. <i>Error:</i> — (correct)		<i>Context:</i> “Carolina’s body taken by Omega; personality changes.” <i>Reasoning:</i> AI possession → behavior change (host misattributed). <i>Answer:</i> Carolina behaves abnormally due to Omega. <i>Error:</i> Entity confusion; temporal mismatch; misses hypoxia factor.		<i>Context:</i> “Omega AI causes aggression.” <i>Reasoning:</i> AI influence → abnormal behavior (partial). <i>Answer:</i> Omega explains erratic acts. <i>Error:</i> Omits oxygen-deprivation component; partial causality.		
TV-33a6bd74	Sarge’s Season 15 depression and redemption	<i>Context:</i> “Sarge creates fake enemies, betrays Reds/Blues, later saves them.” <i>Reasoning:</i> depression → betrayal → redemption (temporal). <i>Answer:</i> Depression triggers betrayal; later redemption by saving team. <i>Error:</i> — (correct)		<i>Context:</i> “Sarge experiences burnout.” <i>Reasoning:</i> depression → low morale (truncated). <i>Answer:</i> Sarge acts poorly due to burnout. <i>Error:</i> Misses betrayal–redemption arc.		<i>Context:</i> “Acts irrationally after long missions.” <i>Reasoning:</i> fatigue → misbehavior (truncated). <i>Answer:</i> Irrational actions due to fatigue. <i>Error:</i> Omits betrayal + later change-of-heart.		
TV-1e87f0a2	The Simpsons – Yes Guy / Wiseguy meta humor	<i>Context:</i> “Yes Guy’s ‘Ye-e-e-s?!’ explained by ‘I had a stro-o-o-oke’; Wiseguy labeled stereotype.” <i>Reasoning:</i> stroke gag → speech quirk → meta-reference. <i>Answer:</i> The gag is justified in-universe; Wiseguy isn’t a fixed name. <i>Error:</i> — (correct)		<i>Context:</i> “Running joke across episodes.” <i>Reasoning:</i> repetition → humor (no causal quote). <i>Answer:</i> It’s a recurring joke. <i>Error:</i> Lacks textual evidence explaining the quirk.		<i>Context:</i> “Minor characters recurring jokes.” <i>Reasoning:</i> trope repetition → humor (generic). <i>Answer:</i> Recurring jokes create humor. <i>Error:</i> Descriptive only; misses explicit line/second part.		

Table 2: Cross-domain qualitative comparison on STEM and TV questions. Each model column includes its retrieved context, reconstructed reasoning chain, final answer, and an error note (if any). AutoPrunedRetriever preserves minimal but functional causal paths; HippoRAG2 tends to over-expand associatively; LightRAG collapses to topic-level summaries.