

Towards a Learnable Retriever for Textual Graph Question Answering

Anonymous ACL submission

Abstract

Textual graph question answering (GraphQA) and graph-based retrieval-augmented generation (GraphRAG) have gained increasing attention as a way to ground large language models (LLMs) on structured knowledge. Despite extensive efforts, existing methods suffer from two fundamental limitations. First (*coverage*), retrievers may fail to recall the evidence nodes required for reasoning, resulting in insufficient coverage of the gold answer set. Second (*compactness*), naive attempts to boost recall typically lead to an explosion of retrieved subgraph size, introducing excessive irrelevant information that compromises compactness and overwhelms the LLM. The gold answer set provides critical supervision in retrieving relevant information while filtering out irrelevant ones, which however, is rarely exploited by existing methods. In this work, we propose TALENT, a learnable retriever that utilizes gold answer signals to enhance both retrieval coverage and compactness. Specifically, we stratify graph nodes into three distinct levels: gold answers, answer-related nodes, and irrelevant noise. Afterwards, we adopt a weighted pseudo-label loss to prioritize gold answers and preserve answer-related nodes for coverage, while discarding irrelevant noises for compactness. Empirical results across two GraphQA datasets demonstrate that TALENT consistently improves downstream performance.

1 Introduction

Large language models (LLMs) have achieved impressive progress in general-purpose generation and reasoning (Thoppilan et al., 2022; Bang, 2023), but their outputs are ultimately bounded by what is encoded in the model parameters. Particularly, trained on large-scale open-domain corpora collected before a fixed cutoff date, LLMs may hallucinate or omit critical evidence when a query requires precise, up-to-date, or domain-specific knowledge (Zhao et al., 2023b; Wang et al., 2023).

To mitigate these limitations, Retrieval-Augmented Generation (RAG) has emerged as a dominant paradigm that grounds LLM outputs in external evidence by dynamically retrieving relevant knowledge snippets for context-aware generation (Lewis et al., 2020; Shi et al., 2024; Gao et al., 2023). While conventional RAG methods (Lewis et al., 2020; Shi et al., 2024) mainly retrieve evidence from unstructured text corpora, recent work has increasingly shifted toward textual graphs as a more expressive and structured knowledge source. This trend has spurred the development of Graph-based RAG (GraphRAG) and Graph Question Answering (GraphQA). By exploiting the explicit semantics and relational structure of textual graphs, graph-based RAG methods support more grounded reasoning and complex information synthesis (Logan IV et al., 2019; Luo et al., 2025; He et al., 2024).

Conventional retrieval methods in GraphRAG can be categorized into non-parametric, Graph neural network (GNN)-based, and LLM-based approaches. Non-parametric retrievers construct subgraphs using classical graph search algorithms or predefined rules (Taunk et al., 2023; Yasunaga et al., 2021). GNN-based methods encode graph structure with GNNs and identify task-relevant components to form the retrieved subgraph (Mavromatis and Karypis, 2024). LM-based retrievers utilized LM’s ability in their retrieval process (Jiang et al., 2023).

However, these three categories rely primarily on query similarity, but largely ignored the critical signals provided by the gold answers. Hence, existing retrievers have two fundamental limitations that hinder their effectiveness, including *coverage* and *compactness*. An illustrative example is shown in Figure 1, where we partition the nodes into three groups: (A) *Gold answers* which refer to all nodes appear in the gold answer list provided by the dataset, (B) *answer-related nodes* that are

085 semantically close to the gold answers but are not
086 themselves gold answers, and (C) *irrelevant noise*
087 that is similar to the query, e.g., containing “New-
088 ton” or “Isaac Newton”, yet is not answer-related.

089 First, *coverage* refers to the limitation that an
090 existing query-based retriever may fail to cover
091 the gold answer set in the retrieved subgraph (blue
092 circle). As shown in Figure 1, the retriever only
093 recovers two gold professions, *physicist* and *scien-
094 tist*, while missing other correct professions such
095 as *mathematician*, *astronomer*, *philosopher*, and
096 *chemist*. This happens because similarity-based
097 ranking tends to prioritize nodes that tightly match
098 dominant query tokens (“Isaac Newton”), rather
099 than nodes that directly express the professions-
100 type answers. As a result, the retrieved nodes have
101 low coverage of the gold answer set.

102 Second, *compactness* refers to the limitation that
103 query-based retriever can introduce excessive irrel-
104 evant noise. This happens because query similarity
105 is a weak proxy for *answer-related* evidence: es-
106 pecially in large graphs, many nodes can exhibit
107 strong lexical or semantic overlap with the query.
108 As shown in Figure 1, the retriever additionally
109 pulls in a substantial amount of Newton-named
110 but non-profession nodes, e.g., "Isaac Newton:
111 The Last Sorcerer", "Sir Isaac Newton", and other
112 Newton-centric mentions. These nodes inflate the
113 retrieved subgraph, consume limited context bud-
114 get, and can distract or mislead the downstream
115 generator from the few signals that actually deter-
116 mine the correct professions.

117 A key observation is that the *gold answer set*
118 provides precisely the missing supervision for re-
119 trieval. Unlike query-relevant retrieval methods,
120 gold answers specify what should be covered, and
121 can therefore help alleviate the two limitations. For
122 *coverage*, gold answers identify the evidence nodes
123 that are truly required for correct reasoning. Train-
124 ing the retriever with gold-answer signals encour-
125 ages it to consistently include these answer-related
126 nodes, rather than relying solely on naive query
127 resemblance that may miss them. For *compactness*,
128 gold answers also help distinguish *answer-related*
129 *evidence* from other weakly related or even irrel-
130 evant nodes. Specifically, gold supervision can
131 guide the model to prioritize nodes that explain or
132 connect to the correct answers, while suppressing
133 those that are unlikely to contribute to the correct
134 answer. As a result, the retriever can remain recall-
135 oriented on the right signals without inflating the
136 retrieved subgraph. To date, such gold signals are

137 rarely exploited to guide or train the retriever.

138 In this work, we propose TALENT, a learnable
139 retriever that leverages gold answer signals to im-
140 prove both retrieval coverage and compactness. We
141 start from a rule-based candidate subgraph that
142 provides a high-recall pool of nodes. To improve
143 retrieval coverage, we train a retriever based on
144 a weighted pseudo-label objective, incorporating
145 gold answer signals to guide the retrieval process.
146 To promote answer-critical evidence, we prioritize
147 gold answer and down-weight other rest signals
148 to improve retrieval compactness. Intuitively, the
149 loss assigns the highest weight to gold answers to
150 ensure coverage, assigns a much smaller weight to
151 answer-related nodes to retain useful context, and
152 ignores irrelevant noise to enforce compactness.
153 This design injects gold signal supervision into
154 retrieval process. After retriever training, the gen-
155 erator is fine-tuned on contexts constructed from
156 the retriever outputs, translating improved retrieval
157 into better downstream QA performance. Our con-
158 tributions are summarized as follows:

- 159 • **Idea.** Gold answer signals effectively improve
160 both coverage and compactness, outperforming
161 retrieval based solely on query similarity.
- 162 • **Method.** We propose a learnable retriever
163 trained with a weighted pseudo-label objective
164 that leverages gold answers to guide retrieval.
165 By prioritizing gold answer nodes and down-
166 weighting weakly related nodes, the retriever
167 leads to better downstream QA performance.
- 168 • **Evaluation.** We conduct extensive experiments
169 on GraphQA benchmarks to validate both re-
170 trieval effectiveness and end-task QA perfor-
171 mance. Our method outperform best competitor
172 by at least 1.13%.

173 2 Preliminaries

174 **Notation.** We use calligraphic letters to denote
175 sets, e.g., \mathcal{V} and \mathcal{E} for node and edge sets. A textual
176 graph is represented as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where each
177 node $v \in \mathcal{V}$ has a textual attribute $\text{text}(v)$, e.g.,
178 an entity name or concept, and each edge $e =$
179 $(u, r, v) \in \mathcal{E}$ denotes a relation r connecting the
180 source node u to the destination node v . The textual
181 attribute of an edge is denoted as $\text{text}(e)$. We use \mathcal{S}
182 to denote the subgraph of \mathcal{G} , where $\mathcal{S} \subseteq \mathcal{G}$ indicates
183 that its node set $\mathcal{V}_{\mathcal{S}} \subseteq \mathcal{V}$ and its edge set $\mathcal{E}_{\mathcal{S}} \subseteq \mathcal{E}$.
184 We use q for the query and a for the final answer
185 from the generator. We use a^* to denotes the token
186 sequence of the gold answer in the dataset and use

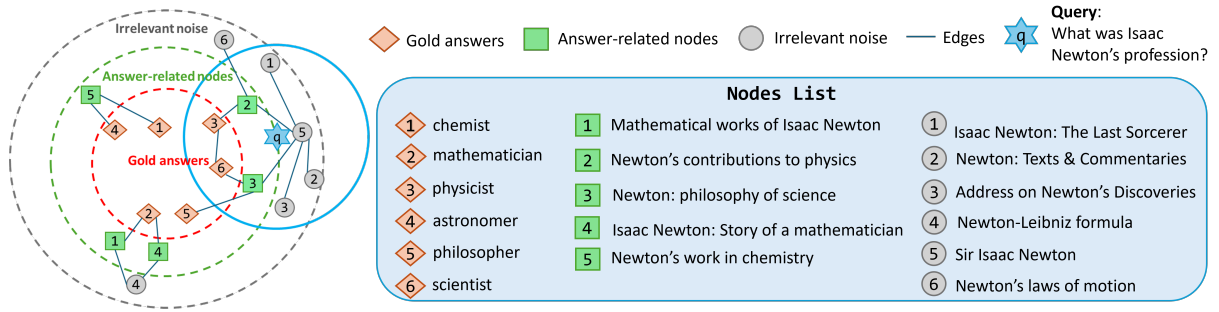


Figure 1: **An illustrative example of two limitations of query-based graph retrieval.** For the query “What was Isaac Newton’s profession?”, a query-similarity retriever returns a Newton-centric subgraph (blue circle). We partition the nodes into three groups: *Gold answers* (only a small subset of the gold answer set is retrieved by the subgraph in the blue circle, showing limited *coverage*), *Answer-related nodes* (topically relevant but not exact gold answers), and *Irrelevant noise* (nodes that match the query tokens such as “Newton” yet are not answer-related). The presence of substantial noise alongside partial gold coverage highlights the *compactness* issue: improving recall by expanding query-similar candidates can inflate the subgraph with distractors.

\mathcal{V}^* as the set of nodes appearing in the gold answer, i.e., the gold answer set. We denote the retriever and generator by R_θ and G_ϕ , respectively.

2.1 Problem Setup

We consider a general retrieve-then-generate setup (He et al., 2024) for textual graph question answering (GraphQA). Each instance in the dataset consists of a natural-language query q and a textual graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ as the knowledge base. Due to the large scale of \mathcal{G} , it is impractical to directly serialized \mathcal{G} into an LLM prompt, hence the system first retrieves a relevant *subgraph*, conditioned on which, an answer is further generated.

Formally, a retriever R_θ extracts a subgraph \mathcal{S} from the textual graph \mathcal{G} based on the query q :

$$\mathcal{S} = R_\theta(q, \mathcal{G}). \quad (1)$$

A generator G_ϕ then predicts the answer a based on the query q and the retrieved subgraph \mathcal{S} :

$$a = G_\phi(q, \mathcal{S}). \quad (2)$$

2.2 Candidate Connected Subgraph

For large textual graphs, directly serializing the full graph often exceeds the LLM context budget. Following G-Retriever (He et al., 2024), we first construct a *candidate connected subgraph* \mathcal{S} using the Prize-Collecting Steiner Tree (PCST) algorithm (Bienstock et al., 1993), and use \mathcal{S} as the input context for downstream modules.

Indexing and similarity retrieval. Given a query q and a textual graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, we embed the query and textual attributes of node \mathcal{V} and edges \mathcal{E} using a sentence encoder such as Sentence-BERT (Reimers and Gurevych, 2019):

$$\mathbf{z}_q = \text{LM}(q), \quad (3) \quad 219$$

$$\mathbf{z}_v = \text{LM}(\text{text}(v)), \quad v \in \mathcal{V}, \quad (4) \quad 220$$

$$\mathbf{z}_e = \text{LM}(\text{text}(e)), \quad e \in \mathcal{E}. \quad (5) \quad 221$$

We compute cosine similarity scores and retrieve the top- k_v nodes and top- k_e edges:

$$\mathcal{V}_{\text{sim}} = \arg \text{topk}_{v \in \mathcal{V}, k=k_v} \cos(\mathbf{z}_q, \mathbf{z}_v), \quad (6) \quad 224$$

$$\mathcal{E}_{\text{sim}} = \arg \text{topk}_{e \in \mathcal{E}, k=k_e} \cos(\mathbf{z}_q, \mathbf{z}_e), \quad (7) \quad 225$$

where k_v and k_e control the number of retrieved nodes and edges, respectively.

Subgraph construction. We assign *prizes* to retrieved candidates based on their ranks, so higher-ranked items receive larger prizes. Concretely, for $v \in \mathcal{V}_{\text{sim}}$ with rank r_v we set $p(v) = k_v - r_v$ to prioritize those with higher similarity, and analogously define $p(e) = k_e - r_e$ for $e \in \mathcal{E}_{\text{sim}}$; items outside these retrieved lists receive zero prize. PCST then selects a connected subgraph $\mathcal{S} = (\mathcal{V}_{\mathcal{S}}, \mathcal{E}_{\mathcal{S}})$ that trades off prize coverage and compactness:

$$\mathcal{S} = \arg \max_{\substack{\mathcal{S} \subseteq \mathcal{G} \\ \text{connected}}} \left(\sum_{v \in \mathcal{V}_{\mathcal{S}}} p(v) + \sum_{e \in \mathcal{E}_{\mathcal{S}}} p(e) - c |\mathcal{E}_{\mathcal{S}}| \right). \quad (8) \quad 238$$

where c controls the size-coverage trade-off.

3 Method

Overview. Our proposed TALENT mainly consists of two parts: (i) a learnable nodes retriever R_θ trained with weighted pseudo-label supervision (Section 3.1), and (ii) an answer generator G_ϕ fine-tuned with standard Supervised Fine-Tuning (SFT) (Section 3.2). The framework for TALENT is shown in Figure 2. Given a query q and a full textual graph \mathcal{G} , we first apply a rule-based graph

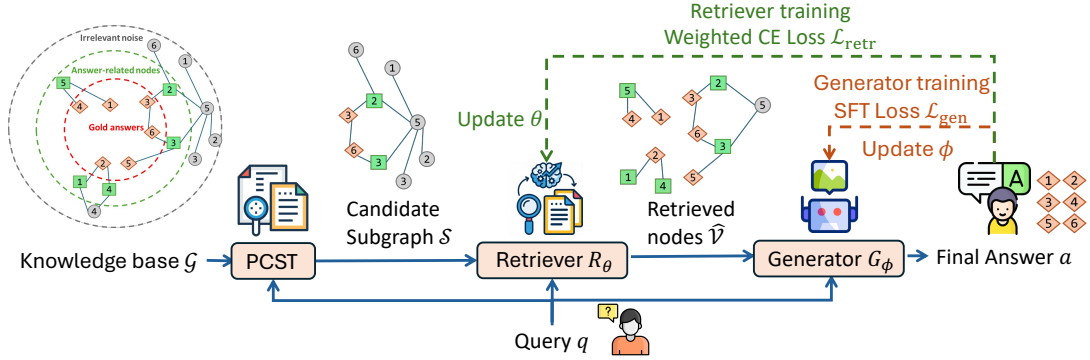


Figure 2: Overview of TALENT. The proposed TALENT includes the following steps: (1) Given Knowledge base \mathcal{G} and query q , PCST retrieves a candidate subgraph \mathcal{S} . (2) Given the candidate subgraph \mathcal{S} and query q , we train a learnable retriever R_θ with weighted pseudo-label supervision, producing $\hat{\mathcal{V}}$. (3) Given retrieved node set $\hat{\mathcal{V}}$ and q , we train a generator G_ϕ with standard SFT. (4) G_ϕ outputs the final answer a .

builder (PCST) to obtain an initial candidate connected subgraph \mathcal{S} that fits the context budget. We then train a learnable retriever R_θ with weighted pseudo-label supervision to select answer-relevant nodes $\hat{\mathcal{V}}$ from \mathcal{S} (Section 3.1). Using answer-relevant nodes $\hat{\mathcal{V}}$, we construct textual contexts and fine-tune a generator G_ϕ with standard supervised learning (Section 3.2).

3.1 Retriever Training

As discussed earlier, PCST and other query-similarity based retrievers often fail to cover the answer-related nodes and include irrelevant noise, resulting in limited *coverage* and poor *compactness*. We therefore use gold answer signals to guide retriever training, encouraging prioritization of gold answers while down-weighting other signals. We first describe TALENT’s retriever prompting and then introduce the weighted pseudo-label objective. **Retriever prompting.** We implement our learnable retriever R_θ as an instruction-following LLM (with LoRA) that performs node selection based on a provided candidate subgraph. Given a query q and the PCST subgraph \mathcal{S} , we prompt the LLM to select the candidate nodes that are most helpful for answering q . The LLM outputs should be exactly one line in the format of “name1 | name2 | ...”. Each name such as name1 and name2 here represents the names of the nodes R_θ finally selects. This constrained format makes the output directly interpretable as a retrieved nodes set, which we denote by $\hat{\mathcal{V}} = R_\theta(q, \mathcal{S})$. It is worth noting that since our retriever R_θ is an LLM, it leverages its inherent parametric knowledge during the selection process. Consequently, the output node set $\hat{\mathcal{V}}$ may include nodes that were not present in the initial candidate subgraph \mathcal{S} (i.e., $\hat{\mathcal{V}} \not\subseteq \mathcal{V}_\mathcal{S}$). In our method, $\hat{\mathcal{V}}$ encompasses the complete set of nodes actually

generated by R_θ , reflecting both the information from \mathcal{S} and the LLM’s internal knowledge.

Fine-tuning with weighted pseudo-label supervision. Recall that existing works in GraphRAG use retrievers that are largely query-similarity driven, which leads to two limitations: coverage and compactness. Our goal is to train the retriever to output the node set $\hat{\mathcal{V}}$ that achieves both high *coverage* and *compactness*: it should reliably include nodes that appears in the gold answer set \mathcal{V}^* (coverage) while avoiding redundant or weakly related candidates (compactness). The gold answer provides precisely the missing supervision, i.e., the gold signal, for retrieval. To inject gold signals into retrieval, we leverage the dataset-provided gold answers and stratify candidate nodes into three levels: (i) *gold answers*, which are candidate nodes that are overlapped with the gold answer set \mathcal{V}^* , and provides evidence that the retriever should prioritize to support correct reasoning, (ii) *answer-related nodes*, which are candidate nodes within \mathcal{S} but are not gold answers, and may still helpful for providing context, though they are less directly tied to the final answer, and (iii) *irrelevant noise*, which are nodes not contained in \mathcal{S} , and are completely irrelevant and should be suppressed to maintain compactness. Given the candidate subgraph \mathcal{S} with node set $\mathcal{V}_\mathcal{S}$, and the gold answer set \mathcal{V}^* , we utilize these nodes as weighted pseudo label supervision by defining

$$\mathcal{P} = \mathcal{V}^* \cap \mathcal{V}_\mathcal{S} \quad (\text{strong positives; level-1}), \quad (9)$$

$$\mathcal{W} = \mathcal{V}_\mathcal{S} \setminus \mathcal{V}^* \quad (\text{weak positives; level-2}). \quad (10)$$

All nodes outside $\mathcal{V}_\mathcal{S}$ are treated as irrelevant noise (level-3) during retriever training.

We train R_θ as conditional generation of a target string y that concatenates nodes in \mathcal{P} and \mathcal{W} (separated by the pipe symbol “|”; strong positives first,

then weak positives)

$$y = p_1 | \dots | p_{|\mathcal{P}|} | w_1 | \dots | w_{|\mathcal{W}|}, \quad (11)$$

where $p_i, \forall i = 1, \dots, |\mathcal{P}|$ are nodes in \mathcal{P} and $w_j, \forall j = 1, \dots, |\mathcal{W}|$ are nodes in \mathcal{W} . Let x denote the retriever input prompt, which includes the query q and the textualized candidate subgraph \mathcal{S} , and let $y_{1:T}$ be the tokenized target sequence, where T is its token length. To prioritize level-1 (gold answer) nodes while still preserving level-2 (answer-related) candidates as a weak learning signal, we minimize the token-weighted negative log-likelihood:

$$\mathcal{L}_{\text{retr}}(\theta) = - \sum_{t=1}^T \omega_t \log p_{\theta}(y_t | x, y_{<t}). \quad (12)$$

Since we want to encourage R_{θ} to prioritize gold answer nodes while still allowing potentially useful answer-related nodes to contribute as a weak learning signal, we adopt the following token weights to achieve this

$$\omega_t = \begin{cases} 1, & \text{if } y_t \text{ belongs to an node in } \mathcal{P}, \\ \lambda, & \text{if } y_t \text{ belongs to an node in } \mathcal{W}, \\ 0, & \text{otherwise,} \end{cases} \quad (13)$$

where λ is a hyperparameters that controls the contribution of weak positives relative to strong positives.

3.2 Generator Training

Generator Prompting Given the predicted node set $\hat{\mathcal{V}} = R_{\theta}(q, \mathcal{S})$ output by the trained R_{θ} , we form the generator input as a *prompt* that concatenates the query q with subgraphs derived from $\hat{\mathcal{V}}$. These subgraphs contains (i) the retrieved nodes in $\hat{\mathcal{V}}$ and (ii) their associated edges $\hat{\mathcal{E}}$ that appear in \mathcal{S} . Concretely, we include edges for which at least one endpoint is in $\hat{\mathcal{V}}$:

$$\hat{\mathcal{E}} = \{(u, r, v) \in \mathcal{E} : u \in \hat{\mathcal{V}} \text{ or } v \in \hat{\mathcal{V}}\}. \quad (14)$$

We then serialize the retrieved nodes $\hat{\mathcal{V}}$ together with the selected edges $\hat{\mathcal{E}}$ into a *serialized graph text* and prepend the query q to form the generator prompt. Importantly, this *serialized graph text* is not required to form one whole connected subgraph. This design avoids forcing unnecessary “bridge” nodes or edges solely to maintain connectivity, which often increases irrelevant context and hurts compactness. Instead, we allow the retrieved evidence to appear as multiple disconnected components, as long as it covers the nodes in $\hat{\mathcal{V}}$. Moreover, as we mentioned in Section 3.1, $\hat{\mathcal{V}}$ may include nodes that come from R_{θ} ’s parametric knowledge

and do not appear in the initial candidate subgraph \mathcal{S} , we retained these nodes as additional isolated nodes in the *serialized graph text*.

Fine-tuning with standard SFT. We fine-tune an llm-based generator G_{ϕ} using standard supervised fine-tuning (SFT). For each training instance, we first run the trained retriever R_{θ} to obtain the generator prompt constructed in the previous paragraph (i.e., the *serialized graph text* formed from $\hat{\mathcal{V}}$ and the associated $\hat{\mathcal{E}}$, together with the query q). We then train G_{ϕ} to generate the gold answer a^* conditioned on this prompt:

$$\mathcal{L}_{\text{gen}}(\phi) = - \sum_{t=1}^{T_a} \log p_{\phi}(a_t^* | q, \hat{\mathcal{V}}, \hat{\mathcal{E}}, a_{<t}^*), \quad (15)$$

where $a^* = (a_1^*, \dots, a_{T_a}^*)$ denotes the *token* sequence of the gold answer, and T_a is its length. This SFT step transfers the improved evidence selection from retrieval into downstream QA.

3.3 Inference and Evaluation

At inference time, we apply the full retrieve-then-generate pipeline using the retriever and generator fine-tuned on the training (and validation) splits. For each test instance (q, \mathcal{G}) , we first construct the candidate subgraph \mathcal{S} by PCST, then run the learned retriever to obtain $\hat{\mathcal{V}} = R_{\theta}(q, \mathcal{S})$. We materialize $(\hat{\mathcal{V}}, \hat{\mathcal{E}})$ into a textual graph context and feed it to the generator to produce the final answer $a = G_{\phi}(q, \hat{\mathcal{V}}, \hat{\mathcal{E}})$. We report task metrics computed from the generator outputs on the test split.

4 Experiments

4.1 Experiment Setup

Datasets and Metrics. To evaluate the effectiveness of our method on textual QA tasks, we conduct experiments on two textual graph question answering datasets, ExplaGraphs and WebQSP from the GraphQA benchmark (He et al., 2024). Specifically, (1) ExplaGraphs’s primary task is to determine whether the arguments are supportive or contradictory. (2) For the knowledge graph question answering dataset WebQSP, the primary task is to answer questions that necessitate multi-hop reasoning. For evaluation, we adopt Accuracy for ExplaGraphs and Hit@1 for WebQSP. Dataset statistics are summarized in Table 5, and detailed experimental configurations are provided in Appendix A.

Baselines. We compare against two groups of baselines. (1) Inference-Only methods that include Zero-shot, Zero-CoT (Kojima et al., 2022), CoT-BAG (Wang et al., 2024), KAPING (Baek et al.,

2023), and Graph-based Inference. (2) Prompt-Tuning methods that include GraphToken (Perozzi et al., 2024), G-Retriever (He et al., 2024), and TAONA-A (Yan et al., 2025). For G-Retriever, we include one variant that fine-tunes the LLM-based generator with LoRA (G-Retriever w/ LoRA).

Experiment pipeline. We adopt the open-source LLaMA2-Chat-7B model (Touvron et al., 2023) as both the retriever and generator, and conduct all experiments on a single NVIDIA A100 GPU. We follow the same train/validation/test splits as G-Retriever and report results averaged over four random seeds 0, 1, 2, 3. For WebQSP, since the full KG exceeds the context budget, we first construct a PCST candidate subgraph with $k_v=k_e=10$ and then apply our learnable retriever. During retriever training, we use the weighted pseudo-label objective with strong-positive weight 1 and weak-positive weight $\lambda = 0.01$. Additional hyperparameter details are provided in Appendix B.

4.2 Benchmark results

Table 1 reports the benchmark results on WebQSP (Hit@1) and ExplaGraphs (Accuracy). Overall, TALENT achieves the best performance on both datasets, improving over the strongest competitors. We draw the following three observations.

(1) TALENT achieves state-of-the-art performance. On WebQSP, TALENT attains a Hit@1 of 74.92%, ranking first among all methods and outperforming the strongest prior baselines, G-Retriever w/ LoRA and TAONA-A by 1.13% and 3.69%, respectively. On ExplaGraphs, TALENT similarly achieves the highest accuracy of 91.79%, improving over G-Retriever w/ LoRA and TAONA-A by 4.74% and 4.78%, respectively. These consistent gains across WebQSP and ExplaGraphs indicate that TALENT generalizes well under different graph QA objectives and metrics.

(2) Comparison to inference-only methods. On WebQSP, inference-only baselines (Zero-shot, Zero-CoT, CoT-BAG, KAPING, and Graph-based Inference) remain substantially behind TALENT, with the strongest baseline achieving a Hit@1 of only 52.64% compared to 74.92% for TALENT. This performance gap stems from their reliance on test-time prompting and rule-based subgraph selection, without task-specific fine-tuning of either the retriever or generator, leaving the LLM poorly adapted to textual-graph reasoning. In contrast, TALENT fine-tunes both components, enabling more reliable evidence selection and more accu-

rate answers. Notably, Graph-based Inference performs particularly poorly, with a Hit@1 of 47.22% on WebQSP and an accuracy of 33.93% on ExplaGraphs, indicating that merely exposing graph information can even hinder reasoning when the context is noisy or misaligned. It further show that learning to select answer-relevant nodes is essential beyond simply providing the graph.

(3) Comparison to prompt-tuning methods. On WebQSP, prompt-tuning approaches form substantially stronger baselines than inference-only prompting methods. In particular, G-Retriever and its LoRA variant demonstrate that fine-tuning with a PCST candidate subgraph yields large gains. Nevertheless, TALENT further improves over G-Retriever w/ LoRA from 73.79% to 74.92%, indicating remaining headroom beyond rule-based, query-driven retrieval. Overall, TALENT’s superior performance supports our central claim: training retrieval with gold signals via weighted pseudo-label supervision better balances answer-node coverage and compactness, reducing redundancy and benefiting generation. On ExplaGraphs, prompt-tuning methods again outperform inference-only baselines (e.g., GraphToken 85.08%, G-Retriever w/ LoRA 87.05%, TAONA-A 87.01%), yet TALENT remains best with an accuracy of 91.79%.

4.3 Coverage-Compactness Tradeoff

In this section, we study how TALENT’s retriever achieves a better coverage-compactness tradeoff for its output node set $\hat{\mathcal{V}}$. We conduct experiments on WebQSP, with results shown in Figure 3, and make the following observations.

Compactness: similar coverage with substantially fewer nodes. We first compare the compactness of nodes retrieved by the initial PCST candidate subgraph ($k_v=k_e=10$) and TALENT, as shown in Figure 3. Under similar gold-node coverage, TALENT’s retriever output $\hat{\mathcal{V}}$ contains significantly fewer nodes. Both methods achieve comparable Hit@1 and Recall on WebQSP, yet $\hat{\mathcal{V}}$ includes only about half as many nodes: on average 16.77 nodes for TALENT versus 31.75 nodes for the PCST subgraph. Meanwhile, $\hat{\mathcal{V}}$ attains notably higher precision and F1, improving precision from 9.26% to 19.98% and F1 from 12.22% to 21.66%. These results suggest that the weighted pseudo-label supervision effectively incorporates gold signals, enabling the retriever to filter redundant or weakly related nodes and produce a more compact context without sacrificing coverage.

Table 1: Benchmark results (%).

Dataset (Metrics)	ExplaGraphs (ACC)	WebQSP (Hit@1)
Zero-shot	56.50	41.06
Zero-CoT (Kojima et al., 2022)	57.04	51.30
CoT-BAG (Wang et al., 2024)	57.94	39.60
KAPING (Baek et al., 2023)	62.27	52.64
Graph-based Inference	33.93	47.22
Frozen LLM + Prompt Tuning (PT)	57.63	48.34
GraphToken (Perozzi et al., 2024)	85.08	57.05
G-Retriever	85.16	70.49
G-Retriever w/ LoRA	87.05	73.79
TAONA-A (Yan et al., 2025)	87.01	71.23
TALENT	91.79	74.92

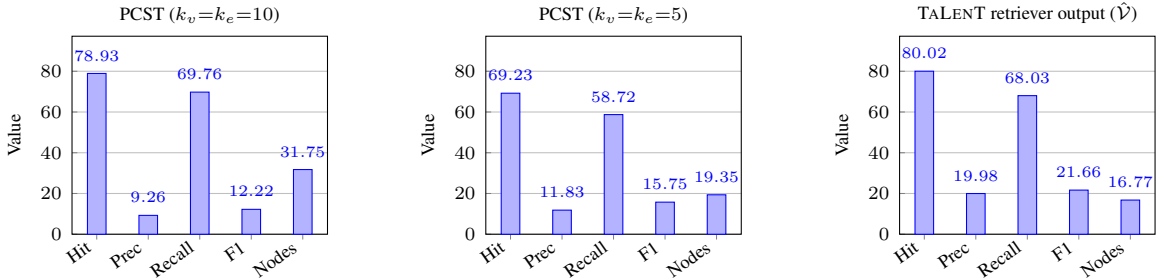


Figure 3: WebQSP test-set retrieval results (macro-averaged). We measure how well PCST subgraphs and TALENT’s learned retriever output $\hat{\mathcal{V}}$ cover the gold answer set \mathcal{V}^* . We compare: (1) the PCST candidate subgraph with $k_v=k_e=10$ (this is the retriever input), (2) a smaller PCST subgraph with $k_v=k_e=5$ (a size-matched reference), and (3) the nodes set $\hat{\mathcal{V}}$ selected by TALENT’s learned retriever. Our retriever keeps coverage close to PCST($k_v=k_e=10$) while using far fewer nodes, and it also covers more gold nodes than the size-matched PCST($k_v=k_e=5$).

Table 2: Breakdown of gold answer coverage contributed by different parts of the learned retriever output on WebQSP test. All metrics are macro-averaged over examples and reported as percentages.

Output subset	Hit@1	Recall
All retrieved nodes	80.02	68.03
Nodes within PCST	75.58	64.76
Nodes out of PCST	10.14	5.21

Coverage: better gold answer coverage under comparable node number. To further compare coverage, we report results for a smaller PCST subgraph with $k_v=k_e=5$, whose size is comparable to that of TALENT. As shown in Figure 3, under a similar node budget, TALENT’s retriever achieves substantially higher coverage. Specifically, TALENT improves Hit@1 from 69.23% to 80.02% and Recall from 58.72% to 68.03% compared to PCST. This demonstrates that, given a similar number of nodes, our gold-signal-involved weighted pseudo-label supervision enables the retriever to achieve significantly higher coverage than the rule-based, query-similarity-driven PCST algorithm.

4.4 Retrieval Beyond PCST

Fine-tuning on an LLM-based retriever may leverage parametric knowledge to retrieve nodes be-

yond the PCST candidate subgraph. We therefore examine whether TALENT’s retriever can retrieve such beyond-PCST nodes and how they contribute to performance. Specifically, we split the retriever output $\hat{\mathcal{V}}$ into nodes inside and outside the PCST subgraph. Shown in Table 2, nodes outside PCST alone recover a non-trivial portion of gold answer nodes, with a Hit@1 of 10.14% and a Recall of 5.21%. Moreover, incorporating beyond-PCST nodes improves Hit@1 from 75.58 to 80.02% and Recall from 64.76% to 68.03%. Results indicate that the trained LLM retriever effectively complements the PCST candidate subgraph with parametric-knowledge nodes that enhance gold-node coverage.

5 Ablation Study

We conduct two ablation studies on WebQSP to illustrate (i) the effect of pseudo-label retriever training (ii) the effect of learnable retriever.

5.1 On the Effect of Pseudo-label Training

Retrieval quality. Table 3 compares TALENT’s retriever against an ablated variant, which uses the same backbone LLM LLaMA2-Chat-7B model but keeps it *frozen* (no training) and directly prompts it to select helpful nodes from the PCST candi-

534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558

Table 3: Retriever ablation on WebQSP test set (macro-averaged over examples). “Ablate retriever training” replace TALENT’s learnable retriever by a frozen retriever without LoRA fine-tuning.

Retriever setting	Hit@1	Prec	Recall	F1
TALENT	80.02	19.98	68.03	21.66
Ablate retriever training	36.61	14.97	26.26	14.97

Table 4: QA ablations on WebQSP (Hit@1, %). All systems follow the same pipeline as TALENT except for the ablated component stated in each row.

System	Hit@1 (%)
TALENT w/o retriever training	68.37
TALENT w/o retriever	72.50
TALENT	74.92

date subgraph. Results show that the frozen retriever fails to recover gold nodes required for answering, with Hit@1, Precision, and Recall of 36.61%, 14.97%, and 26.26%, respectively. In contrast, TALENT’s retriever explicitly prioritizes gold answer nodes during retriever fine-tuning using the pseudo label training, substantially improving coverage with a Hit@1 of 80.02%, a precision of 19.98% and a recall of 68.03%, indicating that injecting gold signals into retriever training is crucial for selecting answer-relevant nodes.

Downstream QA performance. We next examine whether the retrieval improvement translates into better end-task QA. Table 4 summarizes QA Hit@1 under three systems: (i) TALENT with a *frozen* retriever (ablate training on retriever), (ii) TALENT *without* the retriever module (directly feeding the PCST candidate subgraph to the generator), and (iii) the full TALENT system. Switching from a frozen retriever to our pseudo-label trained retriever improves QA Hit@1 from 68.37% to 74.92%. This gain is consistent with Table 3: pseudo-label training recovers substantially more gold nodes, providing the generator with more reliable evidence.

5.2 On the Effect of Learnable Retriever

We further isolate the contribution of the retriever by removing it from the pipeline and directly feeding the PCST candidate subgraph into the generator, i.e., without R_θ or \hat{V} selection. As shown in Table 4, removing the learnable retriever reduces QA performance from 74.92% to 72.50%, indicating that an explicit node-selection retriever provides clear benefits for downstream QA. Notably, using a frozen retriever performs even worse (Hit@1 of 68.37%) than removing the retriever entirely, suggesting that retrieval is not inherently beneficial

when gold answer node coverage is low, as noisy context can mislead the generator. This result further underscores the importance of our gold-signal-based pseudo-label training, which substantially improves retrieval quality for end-task QA.

6 Related Work

Retrieval Augmented Generation. RAG grounds LLMs with external evidence to reduce hallucinations and improve faithfulness (Sun et al., 2024; Gao et al., 2023). Classic RAG systems follow a straightforward index-retrieve-generate workflow (Ma et al., 2023), while many recent variants strengthen retrieval with additional processing before and after retrieval. For example, pre-retrieval techniques such as query transformation, expansion, and rewriting are often adopted to better match the search space (Zheng et al., 2023; Gao et al., 2023), and post-retrieval reranking is commonly used to refine the retrieved evidence set (Blagojevi, 2023). More broadly, modular RAG systems incorporate heterogeneous sources and use LLMs to iteratively refine retrieval queries or evidence selection (Yu et al., 2022).

RAG over Graphs. Extensive works study how to combine graph-structured information with LLMs (Li et al., 2023b; Wang et al., 2024). Prior efforts span general graph modeling (Liu et al., 2023; Lei et al., 2023), multimodal architectures that incorporate graph-like structures (Yoon et al., 2023; Li et al., 2023a). Representative applications include fundamental graph reasoning (Zhao et al., 2023a), node prediction (Sun et al., 2023; Chen et al., 2024), graph-level prediction (Qian et al., 2023). A related line of work discusses GraphRAG / GraphQA, where the system selects graph evidence and conditions an LLM on graph context to answer questions (Logan IV et al., 2019; Luo et al., 2025).

7 Conclusion

We address two fundamental limitations in existing GraphRAG methods: *coverage*, where retrievers fail to recall required evidence nodes, and *compactness*, where excessive irrelevant information overwhelms the LLM. Both issues may stem from the absence of gold answer signals to guide evidence selection. To address this, we propose TALENT, a learnable retriever trained with a weighted pseudo-label loss that prioritizes answer-related signals to improve coverage while suppressing noise for compactness. Experimental results on two GraphQA benchmarks demonstrate that TALENT consistently improves downstream performance.

646 **Limitations**

647 While TALENT shows strong performance, we ac-
648 knowledge several limitations. First, retriever train-
649 ing focuses on selecting nodes only. Edges are
650 added by a deterministic rule based on the selected
651 node set, instead of being directly optimized with
652 learnable edge selection. Second, our weighted
653 pseudo-label supervision reduces node selection to
654 sequence generation by serializing the target nodes
655 into a pipe-separated string in a fixed order, where
656 strong positives precede weak positives, but we
657 do not study whether alternative orderings, would
658 further affect or improve learning. Future work
659 could explore joint node-and-edge selection and
660 order-robust training objectives.

661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715

References

Jinheon Baek, Alham Fikri Aji, and Amir Saffari. 2023. Knowledge-augmented language model prompting for zero-shot knowledge graph question answering. *arXiv preprint arXiv:2306.04136*.

Yejin Bang. 2023. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*.

Daniel Bienstock, Michel X Goemans, David Simchi-Levi, and David Williamson. 1993. A note on the prize collecting traveling salesman problem. *Mathematical programming*, 59(1):413–420.

Vladimir Blagojevi. 2023. Enhancing rag pipelines in haystack: Introducing diversityranker and lostinthemiddleranker.

Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250.

Ziwei Chai, Tianjie Zhang, Liang Wu, Kaiqiang Han, Xiaohai Hu, Xuanwen Huang, and Yang Yang. 2025. Graphllm: Boosting graph reasoning ability of large language model. *IEEE Transactions on Big Data*.

Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei Yin, Wenqi Fan, Hui Liu, and 1 others. 2024. Exploring the potential of large language models (llms) in learning on graphs. *ACM SIGKDD Explorations Newsletter*, 25(2):42–61.

Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. Precise zero-shot dense retrieval without relevance labels. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1762–1777.

Xiaoxin He, Yijun Tian, Yifei Sun, Nitesh Chawla, Thomas Laurent, Yann LeCun, Xavier Bresson, and Bryan Hooi. 2024. G-retriever: Retrieval-augmented generation for textual graph understanding and question answering. *Advances in Neural Information Processing Systems*, 37:132876–132907.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, and 1 others. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.

Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Wayne Xin Zhao, and Ji-Rong Wen. 2023. Structgpt: A general framework for large language model to reason over structured data. *arXiv preprint arXiv:2305.09645*.

Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in*

neural information processing systems, 35:22199–22213. 716
717

Bin Lei, Chunhua Liao, Caiwen Ding, and 1 others. 2023. Boosting logical reasoning in large language models through a new framework: The graph of thought. *arXiv preprint arXiv:2308.08614*. 718
719
720
721

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*. 722
723
724

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474. 725
726
727
728
729
730
731

Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*. 732
733
734

Xin Li, Dongze Lian, Zhihe Lu, Jiawang Bai, Zhibo Chen, and Xinchao Wang. 2023a. Graphadapter: Tuning vision-language models with dual knowledge graph. *Advances in Neural Information Processing Systems*, 36:13448–13466. 735
736
737
738
739

Yuhan Li, Zhixun Li, Peisong Wang, Jia Li, Xiangguo Sun, Hong Cheng, and Jeffrey Xu Yu. 2023b. A survey of graph meets large language model: Progress and future directions. *arXiv preprint arXiv:2311.12399*. 740
741
742
743
744

Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan Zhang. 2023. One for all: Towards training one graph model for all classification tasks. *arXiv preprint arXiv:2310.00149*. 745
746
747
748
749

Robert L Logan IV, Nelson F Liu, Matthew E Peters, Matt Gardner, and Sameer Singh. 2019. Barack’s wife hillary: Using knowledge-graphs for fact-aware language modeling. *arXiv preprint arXiv:1906.07241*. 750
751
752
753
754

Linhao Luo, Jiabin Ju, Bo Xiong, Yuan-Fang Li, Ghulamreza Haffari, and Shirui Pan. 2025. Chatrule: Mining logical rules with large language models for knowledge graph reasoning. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pages 314–325. Springer. 755
756
757
758
759
760

Xinbei Ma, Yeyun Gong, Pengcheng He, Hai Zhao, and Nan Duan. 2023. Query rewriting in retrieval-augmented large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5303–5315. 761
762
763
764
765

Costas Mavromatis and George Karypis. 2024. Gnnrag: Graph neural retrieval for large language model reasoning. *arXiv preprint arXiv:2405.20139*. 766
767
768

769	Bryan Perozzi, Bahare Fatemi, Dustin Zelle, Anton Tsitsulin, Mehran Kazemi, Rami Al-Rfou, and Jonathan Halcrow. 2024. Let your graph do the talking: Encoding structured data for llms. <i>arXiv preprint arXiv:2402.05862</i> .	language? <i>Advances in Neural Information Processing Systems</i> , 36.	825
770			826
771			
772		Keheng Wang, Feiyu Duan, Sirui Wang, Peiguang Li, Yunsen Xian, Chuantao Yin, Wenge Rong, and Zhang Xiong. 2023. Knowledge-driven cot: Exploring faithful reasoning in llms for knowledge-intensive question answering. <i>arXiv preprint arXiv:2308.13259</i> .	827
773			828
774	Chen Qian, Huayi Tang, Zhirui Yang, Hong Liang, and Yong Liu. 2023. Can large language models empower molecular property prediction? <i>arXiv preprint arXiv:2307.07443</i> .		829
775			830
776			831
777			
778	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .	Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. 2024. Next-gpt: Any-to-any multimodal llm. In <i>Forty-first International Conference on Machine Learning</i> .	832
779			833
780			834
781	Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. Replug: Retrieval-augmented black-box language models. In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 8371–8384.		835
782			
783		Yuchen Yan, Aakash Kolekar, Sahika Genc, Wenju Xu, Edward W Huang, Anirudh Srinivasan, Mukesh Jain, Qi He, and Hanghang Tong. 2025. To answer or not to answer (taona): A robust textual graph understanding and question answering approach. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> , pages 6360–6376.	836
784			837
785			838
786			839
787			840
788			841
789	Andy Sun, Tianqi Zheng, Aakash Kolekar, Rohit Patki, Hossein Khazaee, Xuan Guo, George Cai, David Liu, Ruirui Li, Yupin Huang, and 1 others. 2024. A product-aware query auto-completion framework for e-commerce search via retrieval-augmented generation method.	Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. Qaggn: Reasoning with language models and knowledge graphs for question answering. <i>arXiv preprint arXiv:2104.06378</i> .	843
790			844
791			845
792			846
793			847
794		Minji Yoon, Jing Yu Koh, Bryan Hooi, and Ruslan Salakhutdinov. 2023. Multimodal graph learning for generative tasks. <i>arXiv preprint arXiv:2310.07478</i> .	848
795	Shengyin Sun, Yuxiang Ren, Jiehao Chen, and Chen Ma. 2023. Large language models as topological structure enhancers for text-attributed graphs. <i>arXiv preprint arXiv:2311.14324</i> .		849
796			850
797		Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2022. Generate rather than retrieve: Large language models are strong context generators. <i>arXiv preprint arXiv:2209.10063</i> .	851
798			852
799	Dhaval Taunk, Lakshya Khanna, Siri Venkata Pavan Kumar Kandru, Vasudeva Varma, Charu Sharma, and Makarand Tapaswi. 2023. Grapeqa: Graph augmentation and pruning to enhance question-answering. In <i>Companion Proceedings of the ACM Web Conference 2023</i> , pages 1138–1144.		853
800			854
801			855
802			856
803		Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. 2023. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. <i>arXiv preprint arXiv:2303.16199</i> .	857
804			858
805	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, and 1 others. 2022. Lamda: Language models for dialog applications. <i>arXiv preprint arXiv:2201.08239</i> .		859
806			860
807			861
808		Jianan Zhao, Le Zhuo, Yikang Shen, Meng Qu, Kai Liu, Michael Bronstein, Zhaocheng Zhu, and Jian Tang. 2023a. Graphtext: Graph reasoning in text space. <i>arXiv preprint arXiv:2310.01089</i> .	862
809			863
810	Yijun Tian, Huan Song, Zichen Wang, Haozhu Wang, Ziqing Hu, Fang Wang, Nitesh V Chawla, and Panpan Xu. 2024. Graph neural prompting with large language models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pages 19080–19088.		864
811			865
812		Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023b. A survey of large language models. <i>arXiv preprint arXiv:2303.18223</i> , 1(2).	866
813			867
814			868
815			869
816	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, and 1 others. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .		870
817			
818		Huaxiu Steven Zheng, Swaroop Mishra, Xinyun Chen, Heng-Tze Cheng, Ed H Chi, Quoc V Le, and Denny Zhou. 2023. Take a step back: Evoking reasoning via abstraction in large language models. <i>arXiv preprint arXiv:2310.06117</i> .	871
819			872
820			873
821			874
822	Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov. 2024. Can language models solve graph problems in natural	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. <i>arXiv preprint arXiv:2304.10592</i> .	876
823			877
824			878
			879

A Additional Dataset Details

A.1 Datasets

Datasets and Metrics. To evaluate the effectiveness of our method on textual QA tasks, we conduct experiments on two textual graph question answering datasets, ExplaGraphs and WebQSP from the GraphQA benchmark (He et al., 2024). Specifically, (1) ExplaGraphs is developed for generative commonsense reasoning, with a specific focus on the construction of explanation graphs for stance prediction in debates. By providing detailed and unambiguous commonsense-augmented graphs, it enables an evaluation of whether arguments either support or refute a specific belief. ExplaGraphs’s primary task is to determine whether the arguments are supportive or contradictory. Its node denotes commonsense concepts and the edge denotes Commonsense relations. In this dataset, the gold answer is a binary label (i.e., ‘support’ or ‘counter’), which does not provide explicit supervision for the retriever. Therefore, we leverage the LLaMA2-Chat-7B model (Touvron et al., 2023) to generate pseudo gold reasoning node sets, and treat these annotated nodes as pseudo-label supervision \mathcal{V}^* for training the learnable retriever. Due to the small graph size of the graphs in ExplaGraphs, we directly use the entire graph as the candidate context for retrieval. Following GraphQA (He et al., 2024), we evaluate the end task with **Accuracy**. (2) WebQSP is a large-scale multi-hop knowledge graph QA dataset that consists of 4,737 questions. The primary task is to answer questions that necessitate multi-hop reasoning. It employs a subset of Freebase (Bollacker et al., 2008), with a specific focus on facts within 2 hops of the entities identified in the questions. The nodes denote Entities in Freebase and the edges denote Relations in Freebase. The gold answers are lists of Freebase entities, which naturally correspond to node lists and can directly serve as \mathcal{V}^* without requiring additional annotation. For metrics, we report Hit@1, counting a prediction as correct if the generated answer matches any of the answers in the gold answer list. Dataset statistics are summarized in Table 5. See Appendix A.2 for Pseudo-Gold Nodes Annotation for ExplaGraphs.

A.2 Pseudo-Gold Nodes Annotation for ExplaGraphs

ExplaGraphs provides gold supervision as a stance label (support/counter), rather than node-level evidence annotations. To make our learnable retriever

pseudo label supervision applicable, we introduce a lightweight preprocessing step on the training and validation splits: we use Llama-2-7b-chat to annotate a set of pseudo-gold reasoning nodes for each example. Concretely, we provide the model with (i) the full textual explanation graph, (ii) the arguments, and (iii) the gold label, and prompt it to select a set of nodes from the graph that form a critical reasoning path leading to the correct stance decision. The selected nodes are used as \mathcal{V}^* in our weighted pseudo-label objective for training the retriever. This annotation is used only to train our retriever: the generator is still fine-tuned with the original task supervision (predicting support vs. counter), and at inference time the pipeline does not require access to any node annotations.

B Hyperparameter Configuration

Retriever fine-tuning. We fine-tune the retriever with LoRA under our token-weighted cross-entropy objective. Given a question and a candidate subgraph, we construct weighted pseudo labels over candidate nodes: nodes overlapping with gold answer nodes are treated as strong positives with weight 1.0, while the remaining candidate nodes are treated as weak positives with weight $\lambda = 0.01$. The training target is a single-line, pipe-separated node list, formed by concatenating strong positives and weak positives using the delimiter “|” (strong positives first). We deduplicate node names in both retriever outputs and the contexts later used for generator training.

We use LoRA with rank $r=8$, scaling factor $\alpha=16$, and dropout 0.05. Retriever training uses 3 epochs, learning rate 5×10^{-5} , warmup ratio 0.03, per-device batch size 1 with gradient accumulation 8.

Generator fine-tuning. We fine-tune the generator with LoRA using standard supervised fine-tuning on retrieved textual graph contexts and gold targets. We use the same LoRA configuration ($r=8$, $\alpha=16$, dropout 0.05) for generator fine-tuning.

C More Related Works

Graphs and Large Language Models A growing body of work studies how to combine graph-structured information with LLMs (Li et al., 2023b; Wang et al., 2024). Prior efforts span general graph modeling (Liu et al., 2023; Lei et al., 2023), multimodal architectures that incorporate graph-like structures (Yoon et al., 2023; Li et al., 2023a). Rep-

Table 5: Statistics of datasets.

Dataset	ExplaGraphs	WebQSP
#Graphs	2,766	4,737
Average #Nodes	5.17	1370.89
Average #Edges	4.25	4252.37
Node Attribute	Commonsense concepts	Entities in Freebase
Edge Attribute	Commonsense relations	Relations in Freebase
Task	Commonsense reasoning	KGQA
Evaluation metrics	Accuracy	Hit@1

representative applications include fundamental graph reasoning (Zhao et al., 2023a), node prediction (Sun et al., 2023; Chen et al., 2024), graph-level prediction (Qian et al., 2023). A related line of work discusses GraphRAG / GraphQA, where the system selects graph evidence and conditions an LLM on graph context to answer questions (Logan IV et al., 2019; Luo et al., 2025).

Parameter-Efficient Fine-Tuning Parameter-efficient fine-tuning (PEFT) methods enable adapting large pretrained models with limited trainable parameters. Representative techniques include prompt tuning (Lester et al., 2021), prefix tuning (Li and Liang, 2021), and low-rank adaptation such as LoRA (Hu et al., 2022). Related adapter-style methods (e.g., LLaMA-adapter (Zhang et al., 2023)) have also supported rapid adaptation and inspired multimodal instruction-following systems such as MiniGPT-4 (Zhu et al., 2023) and NExT-Chat (Wu et al., 2024). PEFT has recently been explored in graph+LLM settings as well, including GraphLLM (Chai et al., 2025) and GraphToken (Perozzi et al., 2024) for graph reasoning, and GNP (Tian et al., 2024) for knowledge-graph QA.

D Potential Risks

The proposed TALENT improves textual GraphQA/GraphRAG by training an LLM-based retriever with gold-answer signals. While effective, it introduces several potential risks and limitations.

Dependence on candidate-subgraph quality.

For large graphs (e.g., WebQSP), our pipeline first constructs a candidate subgraph via a rule-based method (PCST) before applying the learnable retriever. If this candidate subgraph fails to include sufficient answer-related evidence, downstream retrieval and generation may be constrained accordingly.

Reliance on pseudo-gold node annotations for ExplaGraphs. ExplaGraphs provides only stance labels (support/counter) and does not include node-level evidence. To enable retriever training under our weighted pseudo-label objective, we use an LLM (LLaMA2-Chat-7B) to annotate pseudo-gold reasoning node sets on the training/validation splits. This preprocessing may introduce annotation noise or bias, which can propagate into retriever behavior, although inference does not require node annotations.

Potential ungrounded evidence from parametric knowledge. Because the retriever is an instruction-following LLM, it may leverage parametric knowledge and output nodes not present in the initial candidate subgraph. We retain such nodes as isolated nodes in the serialized context; however, these nodes may not be verifiable within the provided graph evidence, potentially affecting faithfulness.

Structured-output and optimization limitations.

As discussed in our paper limitations, retriever training focuses on selecting nodes (edges are added deterministically), and the weighted pseudo-label supervision reduces node selection to sequence generation with a fixed serialization order (strong positives precede weak positives). We do not study whether alternative orderings or permutation-invariant objectives could further improve learning.

Black-box constraints. Our approach assumes access to model weights (for LoRA fine-tuning) and to intermediate text/graph representations during retrieval and prompting. This may limit applicability in strictly black-box API scenarios where model adaptation or constrained structured decoding is not feasible.

1053	E Use Or Create Scientific Artifacts		
1054	Our work builds on established public bench-		
1055	marks and pre-trained language models for textual		
1056	graph understanding. We evaluate on GraphQA		
1057	benchmark datasets, including ExplaGraphs and		
1058	WebQSP, and use a subset of Freebase as the un-		
1059	derlying knowledge graph for WebQSP. For both		
1060	retriever and generator, we use LLaMA2-Chat-7B		
1061	as the backbone LLM, and follow prior work in first		
1062	constructing a candidate subgraph using PCST.		
1063	We do not modify the original dataset content.		
1064	For ExplaGraphs only, we create <i>pseudo-gold</i> node		
1065	annotations on the training/validation splits to train		
1066	the retriever; these annotations are not used at in-		
1067	ference time.		
1068	E.1 Cite Creators Of Artifacts		
1069	All external artifacts are credited to their original		
1070	publications and repositories. Specifically, we cite:		
1071	(i) the GraphQA benchmark and baselines (e.g.,		
1072	G-Retriever), (ii) ExplaGraphs and WebQSP, (iii)		
1073	the backbone LLM (LLaMA2-Chat-7B), and (iv)		
1074	PCST / sentence-encoder based indexing compo-		
1075	nents used for candidate subgraph construction.		
1076	E.2 Discuss The License For Artifacts		
1077	We comply with the licenses and terms of use of		
1078	all artifacts used in this work.		
1079	• ExplaGraphs. Used under its dataset li-		
1080	cence/terms as provided by the dataset au-		
1081	thors.		
1082	• WebQSP. Used for research purposes under		
1083	their respective licenses/terms as distributed		
1084	by the original providers.		
1085	• Backbone LLM. LLaMA2-Chat-7B is used		
1086	under the corresponding Meta Llama 2 license		
1087	and acceptable use policy.		
1088	E.3 Artifact Use Consistent With Intended		
1089	Use		
1090	We confirm that our use of datasets and pre-trained		
1091	models is consistent with their intended research		
1092	purposes. ExplaGraphs is designed for common-		
1093	sense reasoning with explanation graphs (stance		
1094	prediction), and WebQSP is designed for multi-hop		
1095	KGQA; both align directly with our goal of eval-		
1096	uating retrieval-augmented reasoning over textual		
1097	graphs. The backbone LLM is used for instruction-		
1098	following retrieval and generation, consistent with		
1099	common research usage.		
	E.4 Data Contains Personally Identifying Info		1100
	Or Offensive Content		1101
	The datasets we use are standard public bench-		1102
	marks. WebQSP is derived from Freebase and may		1103
	contain entities corresponding to real people (as		1104
	part of a public knowledge graph). We do not		1105
	add new personally identifying information beyond		1106
	what exists in the original datasets, and we use the		1107
	data solely for research evaluation. ExplaGraphs		1108
	consists of debate-style claims and explanations;		1109
	while not intended to contain Personally Identifi-		1110
	able Information, it may include sensitive topics		1111
	typical of argumentation datasets.		1112
	E.5 Documentation Of Artifacts		1113
	We provide dataset descriptions and usage details		1114
	in Appendix A, including task definitions, graph		1115
	semantics (node/edge attributes), and evaluation		1116
	metrics. We also provide detailed hyperparame-		1117
	ter configurations for retriever and generator fine-		1118
	tuning (including LoRA settings and optimization		1119
	hyperparameters) in Appendix B.		1120
	F Statistics For Data		1121
	Dataset statistics are summarized below (see Table		1122
	5 and Appendix A for details).		1123
	F.1 ExplaGraphs		1124
	• #Graphs: 2,766.		1125
	• Average #Nodes: 5.17; Average #Edges: 4.25.		1126
	• Node attributes: commonsense concepts; edge		1127
	attributes: commonsense relations.		1128
	• Task: commonsense reasoning for stance pre-		1129
	dition; metric: Accuracy.		1130
	F.2 WebQSP		1131
	• #Questions/Graphs: 4,737.		1132
	• Average #Nodes: 1370.89; Average #Edges:		1133
	4252.37.		1134
	• Node attributes: entities in Freebase; edge		1135
	attributes: relations in Freebase.		1136
	• Task: multi-hop KGQA; metric: Hit@1.		1137
	G Computational Experiments		1138
	All computational experiments are reproducible		1139
	given the specifications in the main paper and Ap-		1140
	pendix B.		1141

G.1 Model Size And Budget

We use LLaMA2-Chat-7B as the backbone for both retriever and generator. All experiments are run on a single NVIDIA A100 GPU, and we repeat experiments over four random seeds, reporting the average.

G.2 Experimental Setup And Hyper-params

We follow the GraphQA benchmark protocol and use PCST to build a candidate subgraph for large graphs (WebQSP), then apply TALENT on top. Key hyperparameters include:

- PCST candidate-subgraph settings: top- k_v nodes and top- k_e edges for prize assignment (e.g., $k_v = k_e = 10$ for WebQSP).
- Weighted pseudo-label supervision: strong-positive weight 1.0 and weak-positive weight λ (e.g., $\lambda = 0.01$).
- Structured retriever targets: a single-line pipe-separated node list; strong positives precede weak positives.
- LoRA configurations and optimizer schedules for retriever/generator fine-tuning (see Appendix B).

G.3 Descriptive Statistics

We report downstream task performance on the test split: Accuracy for ExplaGraphs and Hit@1 for WebQSP, following the benchmark definitions.

G.4 Parameters For Packages

All implementation and dependency details (libraries and versions) will be documented in the released codebase and configuration files.

H Ai Assistants In Research Or Writing

We use Llama-2-7b-chat as a preprocessing tool to produce pseudo-gold reasoning node sets for ExplaGraphs training/validation, enabling retriever supervision. This annotation is used only for retriever training; the generator is fine-tuned with original task supervision, and inference does not require node annotations.

In addition, AI assistant tools are used for language editing (spelling/grammar/clarity) of the manuscript. They are not used to generate experimental results or to replace human scientific judgment.