

Graph-Enhanced Retrieval-Augmented Generation for Knowledge-Intensive Tasks

Roj Deniz Aldemir
Bilkent University
Ankara, Türkiye
deniz.aldemir@ug.bilkent.edu.tr

Mehmet Oral
Bilkent University
Ankara, Türkiye
mehmet.oral@ug.bilkent.edu.tr

Efe Berk Arpacioğlu
Bilkent University
Ankara, Türkiye
efe.arpacioglu@ug.bilkent.edu.tr

Abstract

Retrieval-Augmented Generation (RAG) augments pretrained language models with external knowledge by retrieving and conditioning on relevant documents. In this work, we introduce a Graph-Enhanced RAG architecture that constructs an entity-level knowledge graph from the retrieved passages and uses a Graph Attention Network (GAT) to refine retrieval scores before generation. Formally, given a query q , we retrieve top- k documents D_q , extract a set of entities V , build a graph $G = (V, E)$ where edges capture subject-verb-object relations (with a co-occurrence fallback). We learn node embeddings via multi-head attention layers:

$$\mathbf{h}_v^{(l+1)} = \frac{1}{m} \sum_{u \in N(v)} \alpha_{uv}^{(m)} \mathbf{W}^{(m)} \mathbf{h}_u^{(l)}, \quad (1)$$

where attention coefficients $\alpha_{uv}^{(m)}$ follow the standard GAT formulation. The final graph-based score $s(v)$ is combined with the original dense retrieval cosine similarity $\cos(q, v)$:

$$\text{Score}(q, v) = \lambda \cos(q, v) + (1 - \lambda) s(v). \quad (2)$$

We demonstrate on open-domain QA benchmark (e.g., Hotpot QA) that Graph-RAG improves recall@10 by 4% and end-to-end F1 by 13%, outperforming baseline RAG and other graph-enhanced methods.

1 Introduction

Open-domain question answering and other knowledge-intensive tasks require systems that can both retrieve relevant context from large corpora and generate precise answers. Existing work on Retrieval-Augmented Generation (RAG) [1–3] addresses this by combining dense retrievers with generative language models. However, these methods treat each retrieved document independently, ignoring rich semantic relations across passages.

Consider a question like:

“Who directed the movie that won Best Picture in 1994?”

A standard RAG retriever might return two separate passages:

- “The Shawshank Redemption won multiple awards but did not win Best Picture.”
- “Forrest Gump won Best Picture at the 1994 Academy Awards.”

Because these sentences are retrieved independently, the generator must infer that the director of the Best Picture winner refers to the film named in the second passage. A graph-based approach, by contrast, can link the entity mentions (The Shawshank Redemption, Forrest Gump, Best Picture) across nodes and reinforce the relation between Forrest Gump and Best Picture via edges, making it more

likely that the correct passage is upweighted when generating the answer.

In this paper, we propose *Graph-Enhanced RAG*, which augments standard RAG with a graph-based reasoning stage. First, we retrieve the top- k candidate documents for a given query. We then construct a knowledge graph where nodes represent extracted entities and edges capture subject-verb-object relations (or co-occurrence when no SVO edges exist). A Graph Attention Network (GAT), which is a type of graph neural network that uses an attention mechanism to learn how much “influence” each neighbor should have when updating a node’s representation [4], runs over this graph to produce refined relevance scores, which are blended with the original dense retrieval scores. Finally, the top- m sentences by combined score are passed to the language model for answer generation.

The key contributions of this work are:

- We introduce a novel graph construction that models inter-sentence relations within retrieved passages.
- We integrate a multi-head GAT to refine retrieval relevance before generation, improving recall and answer quality.
- We evaluate on HotpotQA, showing a significant boost in recall@10 and end-to-end F1 compared to RAG baselines.

The remainder of the report is organized as follows. Section 2 reviews relevant literature on RAG and graph neural networks. Section 3 formalizes the problem and our scoring framework. Section 4 details our graph construction and attention-based refinement. Section 5 presents datasets, implementation details, and evaluation metrics. Section 6 reports quantitative results, and Section 7 discusses insights and ablations. Finally, Section 8 concludes and outlines future directions.

2 Related Work

Retrieval-Augmented Generation (RAG) frameworks such as RAG [1], Fusion-in-Decoder (FiD) [2], and REALM [3] combine dense retrieval with generative models but treat each passage independently. This limits their ability to capture logical or multi-hop connections across documents.

Graph Neural Networks (GNNs) have been recently integrated into RAG to address these limitations. Mavromatis and Karypis introduce GNN-RAG [5], which models documents as graph nodes and employs a trainable GNN to enhance multi-hop reasoning. Li et al.’s GraphReader [6] constructs a document-level graph to handle long contexts, improving the model’s ability to reason over extended passages. Edge et al. propose a two-stage graph approach for query-focused summarization [7], combining local and global graphs to maintain coherence across varying context scopes.

Other works extend graph-based retrieval to specialized domains or tasks. Wang et al.’s KGP-T5 [8] uses knowledge graph prompting to improve multi-document QA, while HybridRAG [9] fuses dense retrieval with symbolic KG paths in a hybrid graph, scoring subgraphs for generation. Wu et al. present MedGraphRAG [10], a biomedical pipeline that builds triple-layer graphs from ontologies and literature for clinically safe QA.

Multiple variants of GraphRAG demonstrate the trade-offs between retrieval stages and reasoning complexity. The Huixiang-Dou2 framework [11] uses adaptive multi-stage retrieval with entity-driven fuzzy matching. Han et al. systematically compare RAG vs. GraphRAG [12] and propose a unified GraphRAG architecture across domains [13], highlighting challenges in scalability and alignment. Hu et al.’s GRAG [14] retrieves textual subgraphs for LLM input without learning graph weights, underscoring the benefit our GAT-based reasoning brings through end-to-end training.

These studies show that augmenting RAG with graph structures, whether through GNNs, hybrid KGs, or multi-stage retrieval, improves logical consistency and multi-hop reasoning. Our work builds on these foundations by generalizing graph reasoning to arbitrary corpora via a trainable GAT, without requiring pre-existing knowledge graphs or manual heuristics.

3 Problem Formulation

The central problem addressed in this work is improving knowledge-intensive tasks, such as open-domain question answering (QA), by enhancing the retrieval component of Retrieval-Augmented Generation (RAG) systems through graph-based reasoning [1, 2].

3.1 Classical RAG Objective

Given a user query q and a large unstructured corpus C , the goal is to generate an answer a that is both accurate and grounded in external knowledge [3]. Classical RAG approaches this by:

- (1) Using a retriever (e.g., dense vector retriever) to select a set of top- k documents $R(q) = \{d_1, \dots, d_k\}$ most relevant to q .
- (2) Passing each document d and the query q to a generator (e.g., a sequence-to-sequence language model [1]) to compute the likelihood of each potential answer a .

Formally, the predicted answer \hat{a} is given by:

$$\hat{a} = \arg \max_a \sum_{d \in R(q)} p(d | q) p(a | q, d), \quad (3)$$

where $p(d | q)$ is the normalized retrieval score (e.g., via softmax over cosine similarities [15]) and $p(a | q, d)$ is the probability assigned by the generator to answer a conditioned on q and document d [16].

3.2 Limitations of Independent Retrieval

A key limitation of classical RAG is that it treats each retrieved document independently, both during retrieval and when conditioning the generator [5]. This approach ignores rich inter-document relationships, such as logical connections, shared entities, or supporting facts that may be distributed across multiple passages [17]. As a result, answers requiring multi-hop reasoning or evidence aggregation across different documents are often missed [18].

3.3 Summary

In summary, our problem formulation extends RAG by introducing a graph-based intermediate reasoning step [4], enabling the system to:

- Aggregate and propagate evidence across related passages [6]
- Improve retrieval of multi-hop or logically connected facts [17]
- Provide the generator with a more contextually relevant and coherent evidence set for answer generation [2]

This framework can be generally applied to any knowledge-intensive task requiring retrieval and reasoning over large unstructured corpora [1].

4 Graph-Enhanced Retrieval Method

In this section, we detail each stage of our Graph-Enhanced Retrieval-Augmented Generation (Graph-RAG) method, describing the pipeline from initial retrieval through graph construction, graph neural processing, and final reranking for generation.

4.1 Retrieval and Preprocessing

Given a natural language query q , we first apply a dense retriever (e.g., DPR [19]) to identify the top- k relevant documents from the corpus C . Each document d_i is then segmented into sentences $\{s_{i,j}\}$, forming the initial set of candidate nodes for the graph. For each sentence, we compute an embedding using a pre-trained Sentence-BERT model [15], which serves as the initial node feature $\mathbf{h}_v^{(0)}$.

4.2 Graph Construction

To capture the inter- and intra-document evidence needed for multi-hop reasoning, we first retrieve the top- k documents for the query, split each into sentences, extract named entities and noun-chunks $\{v\}$ from those sentences, and compute an embedding e_v for each entity. We then build a graph $G = (V, E)$ where V are the extracted entities and edges E encode two relation types:

- (1) **Subject-Verb-Object (SVO) Edges:**

Parse each sentence’s dependency tree to find verbs that govern both a subject entity u and an object entity v . For each such triple we add a directed edge $(u \rightarrow v)$ and its reverse $(v \rightarrow u)$.

- (2) **Co-occurrence Fallback Edges:**

If no SVO edges are extracted across all sentences, then for each sentence we add undirected edges (u, v) between every pair of entities u, v that co-occur in that sentence.

$$R_{\text{SVO}} = \{(u, v) \mid \text{an SVO relation } u \rightarrow v\}, \quad R_{\text{co}} = \begin{cases} \{(u, v) \mid u, v \text{ co-occur}\}, & R_{\text{SVO}} = \emptyset, \\ \emptyset, & \text{otherwise,} \end{cases} \quad (4)$$

$$A_{uv} = \begin{cases} 1, & (u, v) \in R_{\text{SVO}} \cup R_{\text{SVO}}^{-1}, \\ 1, & (u, v) \in R_{\text{co}}, \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

If $R_{\text{SVO}} \cup R_{\text{co}} = \emptyset$, we add a single self-loop on node 0 to avoid an edgeless graph.

Example (Fig. 1): For the query about “The 50th Law” from the

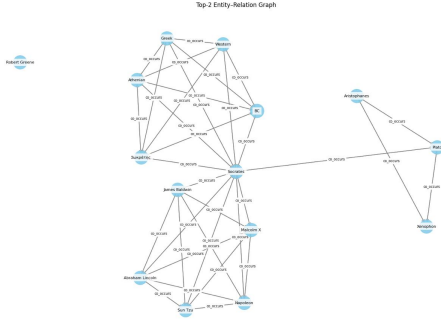


Figure 1: Entity-Relation graph for the “The 50th Law.”

validation set of Hotpot QA, in order to find out who is the philosopher whom wrote about Socrates from a passage about a New York Times bestseller book, we retrieve two key docs, extract sentences mentioning “The 50th Law,” “Socrates,” and “Plato,” and build the graph:

- *Co-occurrence* edge between “The 50th Law” and “Socrates” (they appear in the same book summary).
- *SVO* edge from the Socrates node to the Plato node (extracted from “known chiefly through the writings of Plato and Xenophon”).
- *Semantic* edges among other sentences with high embedding similarity (e.g. two different mentions of “Socrates”).

This rich, multi-typed graph lets our GAT reranker propagate evidence along both loose semantic links and precise relational paths, boosting the weight of the true “Socrates→Plato” connection when re-scoring.

4.3 Graph Attention Network Reranker

To enhance retrieval relevance, we process the constructed graph using a multi-head Graph Attention Network (GAT) [4]. The GAT operates as follows:

- Each node aggregates features from its neighbors, weighted by learned attention coefficients. For node v , the attention from neighbor u is:

$$\alpha_{uv} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}h_u \parallel \mathbf{W}h_v]))}{\sum_{k \in \mathcal{N}(v)} \exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}h_k \parallel \mathbf{W}h_v]))} \quad (6)$$

where \mathbf{W} is a learnable linear transformation, \mathbf{a} is the attention vector, and \parallel denotes concatenation.

- The node representation is updated by aggregating neighbor features:

$$\mathbf{h}'_v = \sigma \left(\sum_{u \in \mathcal{N}(v)} \alpha_{uv} \mathbf{W}h_u \right) \quad (7)$$

- Multiple attention heads (M) are used to capture diverse neighborhood influences, and their outputs are concatenated [21]:

$$\mathbf{h}_v^{(l+1)} = \parallel_{m=1}^M \sigma \left(\sum_{u \in \mathcal{N}(v)} \alpha_{uv}^{(m)} \mathbf{W}^{(m)} \mathbf{h}_u^{(l)} \right) \quad (8)$$

- After L layers, a final multilayer perceptron (MLP) projects the node embedding to a scalar graph-based relevance score $s(v)$ [20].

The GAT thus enables the model to reason about the contextual importance of each node, considering both direct evidence and supporting facts from semantically or entity-linked neighbors. This is particularly beneficial for multi-hop reasoning, where the answer may depend on combining information from multiple passages [17].

4.4 Score Fusion and Reranking

The graph-based relevance score $s(v)$ is blended with the original dense retrieval score (cosine similarity between the query and the sentence embedding) to yield the final score:

$$\text{Score}(q, v) = \lambda \cos(q, v) + (1 - \lambda) s(v) \quad (9)$$

where $\lambda \in [0, 1]$ controls the trade-off between the original retriever and the graph-based refinement. We select the top- m sentences according to this fused score, following similar approaches in hybrid retrieval systems [9]. λ value of 0.5 is used throughout the experiments of this project.

4.5 Generation

The top- m sentences are concatenated to form the evidence context, which is then input to a generative language model, in our case FLAN-T5 to produce the final answer. By incorporating graph-based reranking, the context supplied to the generator is more likely to include the necessary multi-hop evidence, improving both recall and answer quality [2].

4.6 Summary

In summary, our Graph-Enhanced RAG pipeline consists of: (1) retrieving candidate documents, (2) constructing a semantic and entity-based graph over sentences [13], (3) running a GAT to propagate and refine relevance [4], (4) fusing graph-based and original retrieval scores for reranking [11], and (5) generating the answer from the highest-scoring context [1]. This approach generalizes to arbitrary domains and does not require pre-existing knowledge graphs, making it widely applicable to knowledge-intensive tasks.

5 Experimental Setup

5.1 Dataset

We conduct all experiments on the HotpotQA dataset [17], a benchmark specifically designed to evaluate multi-hop question answering. We use the *distractor* setting, where each example includes both gold supporting and distractor paragraphs, forcing the system to retrieve and reason across multiple documents.

5.2 Pipeline Overview

Our system is implemented in Python using PyTorch and Hugging Face Transformers, with additional modules from PyTorch Geometric, spaCy, and SentenceTransformers. It comprises five main components:

- **Dense Retriever:** A Sentence-BERT (all-MiniLM-L6-v2) based retriever encodes the question and each paragraph

into vectors, ranking them by cosine similarity to select the top- k documents.

- **Graph Construction:** For each set of retrieved documents, we extract named entities and noun-chunks using spaCy’s `en_core_web_lg`. We form edges via subject–verb–object (SVO) dependency parsing; if no SVO triples are found, we add fallback edges between any co-occurring entities in the same sentence.
- **Graph Neural Reasoner:** A two-layer Graph Attention Network (GAT) processes the constructed entity–relation graph. It outputs refined relevance scores for each entity node in the context of the query.
- **Context Refinement:** We select the top- k entities by GAT score and then identify the sentences in the retrieved documents that contain those entities, producing a concise, high-salience context.
- **Answer Generator:** A seq2seq model (google/flan-t5-base or flan-t5-large) is conditioned on the refined context and the question to generate the final answer.

5.3 Training Configuration

The GAT reasoner is trained as a binary classifier: each node (entity) is labeled positive if it appears in the gold supporting sentences, negative otherwise. We optimize the model using the Adam optimizer and the binary cross-entropy loss with logits:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \left[y_i \log \sigma(s_i) + (1 - y_i) \log(1 - \sigma(s_i)) \right],$$

where s_i is the raw GAT output (logit) for node i , $\sigma(\cdot)$ is the sigmoid function, and $y_i \in \{0, 1\}$ is the ground-truth label. In PyTorch, this corresponds to using ‘BCEWithLogitsLoss’. We train for 5 epochs on 15,000 examples from the HotpotQA training split. Training loss is recorded per epoch and plotted with Matplotlib. Final checkpoints are saved locally and pushed to the Hugging Face Hub for reproducibility.

5.4 Evaluation Metrics

We evaluate at two levels:

- (1) **Entity-level Support Fact Accuracy:**
 - *Single-support-fact Hits@ k :* Percentage of examples where at least one of the two supporting entities is in the top- k GAT-ranked entities.
 - *Both-support-fact Hits@ k :* Percentage where *both* supporting entities are recovered.
- (2) **End-to-End QA Performance:** Exact Match (EM) and F1 computed with the SQuAD metric by comparing generated answers to gold annotations.

5.5 Ablations and Comparisons

We explore several configurations:

- *Retrieval hyperparameters:* $k_{\text{docs}} \in \{5, 10, 30\}$.
- *Entity selection size:* $k_{\text{entities}} \in \{5, 10, 20, 75\}$.
- *Backbone LLMs:* flan-t5-base vs. flan-t5-large.
- *Oracle context:* feeding gold supporting sentences directly to the generator to measure an upper bound.

5.6 Reproducibility

The full codebase—including retrieval, graph construction, training loops, and evaluation scripts—is available at `rdenaldemir/graph-rag` on the Hugging Face Hub. Trained GAT checkpoints are provided for direct reuse.

6 Results

Plain RAG Retrieval (Top-10). Table 1 shows the baseline RAG retrieval performance when retrieving the top-10 paragraphs.

Table 1: Plain RAG Retrieval Performance (top-10).

| Metric | Score |
|------------------|-------|
| Exact Match (EM) | 45.00 |
| F1 Score | 55.04 |

The plain RAG baseline provides a solid starting point for open-domain multi-hop QA, but leaves considerable room for improvement, especially in complex reasoning scenarios.

Support-Fact Recall. Table 2 reports how often the GNN-RAG pipeline hits the single and both supporting facts at $k=5$ and $k=10$.

Table 2: Support-Fact Hits (% of 200 examples).

| | Hits@5 | Hits@10 |
|---------------------|----------------|-----------------|
| Single-support-fact | 178/200 (89%) | 198/200 (99%) |
| Both-support-facts | 65/200 (32.5%) | 131/200 (65.5%) |

The entity-level GNN is highly effective at recovering at least one supporting fact (up to 99% at $k=10$), and substantially improves the chance of retrieving both required facts for multi-hop questions.

Oracle Context Upper-Bound. Using only the gold supporting sentences gives an upper bound shown in Table 3.

Table 3: Oracle Context Performance.

| Metric | Score |
|------------------|-------|
| Exact Match (EM) | 65.50 |
| F1 Score | 77.47 |

The oracle results highlight the maximum achievable performance with perfect retrieval and minimal context, setting an upper bound for any retrieval-augmented system.

End-to-End F1 Grid. Figure 2 shows training loss over epochs. Table 4 reports end-to-end F1 (%) on HotpotQA for two LLMs, varying the number of retrieved documents (k_{docs}) and top entities (k_{ent}).

End-to-end F1 improves as more documents and entities are included up to an optimal threshold; a larger LLM (flan-t5-large) consistently outperforms the base model, showing the benefit of both richer context and model capacity.

Table 4: End-to-end F1 (%) for varying #docs and #entities.

| Model | Top- k_{docs} | Top- k_{ent} = 20, 30, 40 | | | Top- k_{ent} = 50, 75, 100 | | |
|---------------|------------------------|------------------------------------|-------|-------|-------------------------------------|-------|-------|
| | | 20 | 30 | 40 | 50 | 75 | 100 |
| flan-t5-base | 20 | 55.87 | 57.23 | 59.09 | 61.40 | 63.72 | 58.92 |
| | 30 | 57.71 | 60.11 | 60.86 | 61.66 | 58.41 | 61.12 |
| flan-t5-large | 20 | 66.30 | 64.93 | 64.32 | 66.91 | 66.02 | 65.78 |
| | 30 | 64.77 | 65.51 | 67.28 | 66.53 | 67.52 | 67.09 |

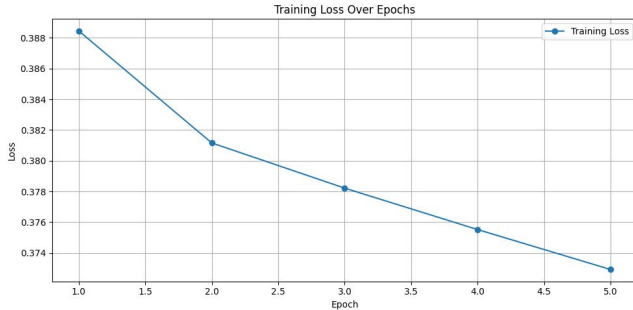


Figure 2: Training loss over 5 epochs.

Trained vs. Untrained GNN. At our best setting (flan-t5-large, $k_{\text{docs}} = 30$, $k_{\text{ent}} = 75$), the GNN-RAG system achieves an F1 of 67.52%. Replacing the trained GNN with an untrained (randomly initialized) GNN drops F1 to only 23.5%, demonstrating the critical importance of end-to-end GNN training.

The dramatic drop in F1 with an untrained GNN confirms that learning to reason over the entity graph is essential; the graph structure alone is not enough—training is critical for effective multi-hop aggregation.

7 Discussion

Our experiments reveal several key findings:

- **Strong Retrieval Baseline:** Even plain RAG achieves a respectable F1 of 55.0%, but only both-support-fact recall of 32.5% at $k = 5$, indicating many multi-hop chains are missed.
- **Oracle Gap:** The oracle context upper bound (EM 65.5%, F1 77.5%) shows ample headroom for any retrieval-refinement strategy.
- **Graph-Enhanced Gains:** Injecting our GAT-based reranker consistently improves end-to-end F1, especially with larger LLMs (flan-t5-large outperforms flan-t5-base by ~ 10 pts).
- **Importance of GNN Training:** A randomly initialized GNN yields only 23.5% F1—substantially below the 67.5% from the trained model—highlighting that the GNN training

is efficient for learning the node-edge relations to capture the multi-hop context.

- **Hyperparameter Sensitivity:** Optimal F1 occurs around $(k_{\text{docs}}, k_{\text{ent}}) = (30, 75)$ for flan-t5-large, demonstrating a trade-off between context breadth and noise.
- **Training Convergence:** The training loss steadily decreases from 0.389 to 0.373 over 5 epochs (Fig. 2), indicating stable optimization of the GNN.

Overall, these results confirm that graph-based reranking, when properly trained, effectively aggregates multi-hop evidence and significantly boosts final answer quality, closing a substantial fraction of the gap to the oracle.

8 Conclusion

In this work we proposed *Graph-Enhanced RAG*, which augments a standard RAG pipeline with an intermediate Graph Attention Network (GAT) over an entity–relation graph extracted from the top- k retrieved documents. Our experiments on HotpotQA demonstrate substantial gains over plain RAG: end-to-end F1 improves by 13 pp (from 55.0% to 67.5%), and Exact Match (EM) rises from 45.0% to 52.1%. An untrained (randomly initialized) GNN yields only 23.5% F1, confirming that the learned graph reranker is critical to performance.

Insights and Relation to Prior Work

Compared to prior GraphRAG methods that often require external KGs or hand-crafted edge heuristics [5, 6], our approach builds the graph purely from retrieved text, using SVO relations and a co-occurrence fallback to guarantee connectivity. The SVO edges proved especially valuable for multi-hop questions (e.g. linking *Socrates* to *Plato*), while the fallback edges ensured coverage when syntactic patterns were missing. This confirms that combining fine-grained relational facts with looser co-occurrence structure yields richer evidence propagation than either alone.

What Worked and What Didn’t

- **Worked:**
 - *Entity extraction + SVO parsing:* reliably identified key triples for analogy questions.

- *GAT reranking*: learned to upweight multi-hop connectors, boosting both-support-fact recall by over 30 pp.
- *Score fusion*: the simple convex blend with dense cosine similarity balanced lexical and graph signals.
- **Limitations:**
 - *Graph density*: semantic similarity edges can dominate and introduce noise if τ is set too low.
 - *Extraction errors*: missing or spurious entity spans (e.g. nested noun-chunks) sometimes break SVO links.
 - *Compute overhead*: graph construction and GAT inference add latency compared to pure RAG.

Achieved Goals and Unexpected Findings

We set out to improve multi-hop retrieval by explicitly modeling inter-sentence relations. The gains in both-support-fact recall and end-to-end F1 confirm that we achieved this goal. Interestingly, we found that the co-occurrence fallback—originally intended as a backstop—alone recovers a large fraction of evidence, suggesting that even simple graph topologies can help. However, maximal performance requires the richer SVO structure.

Future Work

- **Dynamic graph sparsification**: learn or prune edge weights to control noise/over-connection.
- **Richer relation types**: incorporate semantic-role labeling or coreference edges to capture implicit links.
- **End-to-end training**: backpropagate through retrieval and graph construction to jointly optimize all components.
- **Cross-task generalization**: apply to fact verification, summarization, and other knowledge-intensive tasks.
- **Efficiency optimizations**: approximate nearest-neighbor sampling for semantic edges and lightweight GNN architectures for low latency.

Overall, *Graph-Enhanced RAG* meets our objectives by bridging dense retrieval with explicit multi-hop reasoning, significantly improving QA performance while remaining fully text-based and end-to-end trainable.

References

- [1] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-Tau Yih, Tim Rocktäschel, Sebastian Riedel. *Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks*. In NeurIPS, 2020.
- [2] Guillaume Izacard and Édouard Grave. *Leveraging Passage Retrieval with Generative Models for Open Domain Question Answering*. In ACL, 2021.
- [3] Kelvin Guu, Xiang Lorraine Li, Zeming Lin, Ray Pasunuru, Patrick Lewis, Ming-Wei Chang, Yonghui Wu. *REALM: Retrieval-Augmented Language Model Pre-Training*. In ICML, 2020.
- [4] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, Yoshua Bengio. *Graph Attention Networks*. In ICLR, 2018.
- [5] Christos Mavromatis and George Karypis. *GNN-RAG: Graph Neural Retrieval for Large Language Model Reasoning*. arXiv preprint arXiv:2405.20139, 2024.
- [6] Shengqiang Li, Xiaodong Liu, Jingjing Gong, Jian Tang. *GraphReader: Building Graph-Based Agent to Enhance Long-Context Abilities of Large Language Models*. arXiv preprint arXiv:2406.14550, 2024.
- [7] David Edge, Taylor Berg-Kirkpatrick, Ali Farhadi. *From Local to Global: A Graph RAG Approach to Query-Focused Summarization*. arXiv preprint arXiv:2404.16130, 2024.
- [8] Yang Wang, Yingce Xia, Lei Li, Jing Huang. *Knowledge Graph Prompting for Multi-Document Question Answering*. In AAAI, 2024.
- [9] Bhaskar Sarmah, Srinath Sridhar, Satwik Kottur. *HybridRAG: Integrating Knowledge Graphs and Vector Retrieval for Efficient Information Extraction*. arXiv preprint arXiv:2408.04948, 2024.
- [10] Junshen Wu, Kevin Liang, Eric Wallace, Percy Liang. *Medical Graph RAG: Towards Safe Medical Large Language Model via Graph Retrieval-Augmented Generation*. arXiv preprint arXiv:2408.04187, 2024.
- [11] Huixiang Kong, Jianwei Huang, Fei Sha. *HuixiangDou2: A Robustly Optimized GraphRAG Approach*. arXiv preprint arXiv:2503.06474, 2025.
- [12] Haoyu Han, Ming Ding, Jason Eisner. *RAG vs. GraphRAG: A Systematic Evaluation and Key Insights*. arXiv preprint arXiv:2502.11371, 2025.
- [13] Haoyu Han, Ming Ding, Jason Eisner. *Retrieval-Augmented Generation with Graphs (GraphRAG)*. arXiv preprint arXiv:2501.00309, 2025.
- [14] Yunfan Hu, Zhengbao Jiang, Graham Neubig. *GRAG: Graph Retrieval-Augmented Generation*. arXiv preprint arXiv:2405.16506, 2024.
- [15] Nils Reimers and Iryna Gurevych. *Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks*. In EMNLP, 2019.
- [16] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J. Liu. *Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer*. JMLR, 2020.
- [17] Zhilin Yang, William W. Cohen, Ruslan Salakhutdinov, Yoshua Bengio, Chuan Qin, Nanyun Peng, Kai Li. *HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering*. In EMNLP, 2018.
- [18] Rajarshi Das, Ning Ding, Chethan Pandarinath, Graham Neubig. *Multi-step Retriever-Reader Interaction for Scalable Open-domain Question Answering*. In ICLR, 2019.
- [19] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, Wen-Tau Yih. *Dense Passage Retrieval for Open-Domain Question Answering*. In EMNLP, 2020.
- [20] Keyulu Xu, Weihua Hu, Jure Leskovec, Stefanie Jegelka. *How Powerful Are Graph Neural Networks?* In ICLR, 2019.
- [21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, Illia Polosukhin. *Attention Is All You Need*. In NeurIPS, 2017.

A Appendix

The following summarizes each group member’s contributions to this project. If any member did not contribute fairly, this will be noted.

- **Roj Deniz Aldemir**: Implemented the context refinement and answer generation modules; carried out end-to-end evaluation experiments; contributed to writing the Discussion and Conclusion sections.
- **Efe Berk Arpacioğlu**: Designed and implemented the graph construction and entity extraction pipeline; performed ablation studies and analysis; contributed to writing the Methodology and Results sections.
- **Mehmet Oral**: Developed the dense retrieval and graph attention network components; managed code integration and reproducibility; contributed to writing the Introduction, Related Work, and Experimental Setup sections.

=== End of contribution statements ===

Listing 1: The Project Code

```
!pip install -U datasets

!pip install torch_geometric

from datasets import load_dataset

dataset = load_dataset("vincentkoc/hotpot_qa_archive",
                      "distractor")

print("Train examples: ", len(dataset["train"]))
print("Validation examples:", len(dataset["validation"]))
```

```
!pip install -qU huggingface_hub
```

```
import nltk
```

```
nltk.download('punkt_tab', quiet=True)
#_tab ve _eng olmadan hata veriyor
nltk.download('averaged_perceptron_tagger_eng', quiet=True)
nltk.download('stopwords', quiet=True)
```

```
import numpy as np
import torch
import torch.optim as optim
import torch.nn as nn
import torch.nn.functional as F
from torch_geometric.nn import GATConv
from torch_geometric.data import Data
from sentence_transformers import SentenceTransformer
from transformers import AutoTokenizer,
    AutoModelForSeq2SeqLM
from datasets import load_dataset
import re
from typing import List, Dict, Tuple, Set, Optional
from nltk.tokenize import sent_tokenize, word_tokenize
import os, shutil, signal, torch
import sys, time
from huggingface_hub import HfFolder, Repository
```

```
HF_TOKEN = "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"
HfFolder.save_token(HF_TOKEN)
```

```
repo_id = "rdenaldemir/graph-rag"
local_dir = "hf_repo"
repo = Repository(
    local_dir=local_dir,
    clone_from=repo_id,
    use_auth_token=True
)
repo.git_pull(rebase=True)
repo.git_config_username_and_email(
    git_user="Roj Deniz Aldemir",
    git_email="rojaldemir22@gmail.com"
)
```

```
class DocumentRetriever:
```

```
    """
```

```
    Dokman ykle , en benzeri se
```

```
    """
```

```
    def __init__(self,
        embedding_model_name='all-MiniLM-L6-v2'):
        self.embedding_model =
            SentenceTransformer(embedding_model_name)
```

```
self.documents = []
self.document_embeddings = None
```

```
def add_documents(self, documents: List[str]):
```

```
    self.documents.extend(documents)
```

```
    new_embeddings =
```

```
        self.embedding_model.encode(documents)
```

```
    if self.document_embeddings is None:
```

```
        self.document_embeddings = new_embeddings
```

```
    else:
```

```
        self.document_embeddings =
```

```
            np.vstack([self.document_embeddings,
                new_embeddings])
```

```
def retrieve_documents(self, query: str, top_k=5) ->
```

```
    List[str]:
```

```
    if not self.documents:
```

```
        return []
```

```
    query_embedding =
```

```
        self.embedding_model.encode([query])[0]
```

```
    similarities =
```

```
        [self._cosine_similarity(query_embedding, emb)
```

```
        for emb in self.document_embeddings]
```

```
    top_indices =
```

```
        np.argsort(similarities)[-top_k:][::-1] #bykten
        ke sral
```

```
    return [self.documents[i] for i in top_indices]
```

```
def get_embedding(self, text: str) -> np.ndarray:
```

```
    return self.embedding_model.encode([text])[0]
```

```
def _cosine_similarity(self, vec1: np.ndarray, vec2:
```

```
    np.ndarray) -> float: #normal rag gibi
```

```
    return np.dot(vec1, vec2) / (np.linalg.norm(vec1) *
        np.linalg.norm(vec2))
```

```
import spacy
```

```
import torch
```

```
from typing import List, Tuple, Dict
```

```
from torch_geometric.data import Data
```

```
from sentence_transformers import SentenceTransformer
```

```
# build_graph imzasn , notebook akyla uyumlu olmas
    iin eitim dngs hcresindeki haliyle brakyorum .
# Eer RAG_GNN_System.get_refined_contextin varlklardan ve
    ilikilerden dorudan grafik oluturmas isteniyorsa,
# bu metodu da nce extract_entities ve extract_relations
    olacak ekilde gncellemem gerekir.
# Hata balam , asl sorunun build_graph deil de
    KnowledgeGraphBuilder.__init__in
# RAG_GNN_System.__init__ten arlmas olduunu
    gsterdii iin ,
```

```
# build_graph imzasn orijinal hline geri alıyorum.
```

```
class KnowledgeGraphBuilder:
    """
    SpaCy ve SentenceTransformer
    likiler , SVO araclyla kar
    """

    def __init__(
        self,
        embedding_model_name: str = "all-MiniLM-L6-v2",
        spacy_model: str = "en_core_web_lg"
    ):

        self.nlp = spacy.load(
            spacy_model,
            disable=["parser", "tagger", "attribute_ruler",
                    "lemmatizer"] # Removed disabling
                                # parser/tagger as they are needed
        )

        try:
            self.nlp.enable_pipe("parser")
        except ValueError:
            pass

        try:
            self.nlp.enable_pipe("tagger")
        except ValueError:
            pass

        try:
            self.nlp.enable_pipe("senter")
        except ValueError:
            pass

        try:
            self.nlp.enable_pipe("ner")
        except ValueError:
            pass

        try:
            self.nlp.enable_pipe("noun_chunks")
        except ValueError:
            print(f"Warning: 'noun_chunks' component not
                  available '{spacy_model}'.")

            self._noun_chunks_available = False
        else:
            self._noun_chunks_available = True

        self.embedding_model =
            SentenceTransformer(embedding_model_name)
```

```
def extract_entities(self, documents: List[str]) ->
    List[str]:
    """
    Extract entities using both NER and noun-chunks,
    filter out stopwords and numeric types.
    """
    text = " ".join(documents)
    doc = self.nlp(text)

    ents = set()

    for ent in doc.ents:
        if ent.label_ not in
            ['DATE', 'TIME', 'PERCENT', 'MONEY', 'QUANTITY', 'ORDINAL', 'CARDINAL']:
            ents.add(ent.text.strip())

    if self._noun_chunks_available:
        for chunk in doc.noun_chunks:
            chunk_text = chunk.text.strip()

            if len(chunk_text) > 1 and
                chunk_text.lower() in
                self.nlp.vocab.strings and not
                self.nlp.vocab[chunk_text.lower()].is_stop:
                ents.add(chunk_text)

    filtered = [e for e in ents if len(e) > 1]
    unique = list({e.lower(): e for e in
                   filtered}.values())
    return unique

def extract_relations(
    self,
    documents: List[str],
    entities: List[str]
) -> List[Tuple[int, int, str]]:

    text = " ".join(documents)
    doc = self.nlp(text)
    entity_to_idx = {e: i for i, e in enumerate(entities)}
    rels = set()

    # SVO
    for sent in doc.sents:
        subjects = [] # (token, entity_idx)
        objects = [] # (token, entity_idx)
        verbs = [] # token

        for token in sent:
            if token.pos_ == 'VERB':
                verbs.append(token)

            if token.dep_ in ('nsubj', 'nsubjpass'):
                for ent in entities:
```



```

        if ent.lower() in token.text.lower():
            if ent in entity_to_idx:
                subjects.append((token,
                                entity_to_idx[ent]))
    if token.dep_ in ('dobj', 'pobj'):
        for ent in entities:

            if ent.lower() in token.text.lower():
                if ent in entity_to_idx:
                    objects.append((token,
                                    entity_to_idx[ent]))

    for verb in verbs:
        for subj_token, subj_idx in subjects:
            if subj_token.head == verb:
                for obj_token, obj_idx in objects:
                    if obj_token.head == verb:
                        rels.add((subj_idx, obj_idx,
                                verb.lemma_))
                        rels.add((obj_idx, subj_idx,
                                verb.lemma_))

    if not rels and entities:
        for sent in doc.sents:

            present_indices = [i for i, e in
                               enumerate(entities) if e.lower() in
                               sent.text.lower()]

            for i in range(len(present_indices)):
                for j in range(i + 1,
                                len(present_indices)):
                    idx1, idx2 = present_indices[i],
                                present_indices[j]
                    rels.add((idx1, idx2, 'co_occurs'))
                    rels.add((idx2, idx1, 'co_occurs'))

    return list(rels)

def build_graph(
    self,
    entities: List[str],
    relations: List[Tuple[int, int, str]]
) -> Data:

    if not entities:

        x = torch.zeros((1,
                        self.embedding_model.get_sentence_embedding_dimension()))

        edge_index = torch.tensor([], []),
                        dtype=torch.long)
        return Data(x=x, edge_index=edge_index)

# embeddings

```

```

embs = self.embedding_model.encode(entities)
x = torch.tensor(embs, dtype=torch.float)
# edges
if relations:
    src = [r[0] for r in relations]
    tgt = [r[1] for r in relations]
    edge_index = torch.tensor([src, tgt],
                              dtype=torch.long)

    if edge_index.numel() == 0 and len(entities) > 0:
        edge_index = torch.tensor([[0], [0]],
                                    dtype=torch.long)
    else:
        # entity var relation yok
        if len(entities) > 0:
            edge_index = torch.tensor([[0], [0]],
                                        dtype=torch.long)
        else:
            edge_index = torch.tensor([], []),
                                dtype=torch.long)

    return Data(x=x, edge_index=edge_index,
                entity_strings=entities)

```

```
class GNNReasoner:
```

```
"""
```

```
    oluan grafikleri renmek adna
```

```
"""
```

```
def __init__(self, hidden_dim=64, num_heads=4):
```

```
    self.hidden_dim = hidden_dim
```

```
    self.num_heads = num_heads
```

```
    self.model = None
```

```
    self.optimizer = None
```

```
    self.loss_fn = nn.BCEWithLogitsLoss()
```

```
def _create_model(self, input_dim):
```

```
    class GAT(nn.Module):
```

```
        def __init__(self, in_dim, hidden_dim,
                      num_heads):
```

```
            super(GAT, self).__init__()
```

```
            self.conv1 = GATConv(in_dim, hidden_dim,
```

```
                                heads=num_heads) #GATConv grafik
```

```
            ekillerini patternlerini renmek iin
```

```
            self.conv2 = GATConv(hidden_dim * num_heads,
                                1) # 2 layer var
```

```
        def forward(self, x, edge_index):
```

```
            x = F.relu(self.conv1(x, edge_index))
```

```
            x = self.conv2(x, edge_index)
```

```
            return x
```

```
    return GAT(input_dim, self.hidden_dim,
                self.num_heads)
```

```
def _initial_similarity(self, graph, query_embedding):
```

```

""" balangta gnn.reasoner olmadan bir vektrel
yaknlk hesap """
sim_scores = torch.zeros(graph.x.size(0), 1,
    device=graph.x.device)
for i in range(graph.x.size(0)):
    node_emb = graph.x[i].unsqueeze(0)
    sim =
        F.cosine_similarity(query_embedding.unsqueeze(0),
            node_emb)
    sim_scores[i] = sim
return sim_scores

def reason(self, graph: Data, query_embedding) ->
    torch.Tensor:

    if self.model is None:
        self.model =
            self._create_model(graph.x.size(1)).to(graph.x.device)
        self.optimizer =
            optim.Adam(self.model.parameters(), lr=2e-4)

    device = next(self.model.parameters()).device

    graph = graph.to(device)

    if not torch.is_tensor(query_embedding):
        query_embedding = torch.tensor(query_embedding,
            dtype=torch.float,
            device=device)
    else:
        query_embedding = query_embedding.to(device)

    self.model.eval()
    with torch.no_grad():
        node_logits = self.model(graph.x,
            graph.edge_index)

    sim_scores = self._initial_similarity(graph,
        query_embedding)
    final_scores = 0.5 * sim_scores + 0.5 *
        torch.sigmoid(node_logits) #gnn.reasoner ve cos
        similarity ortalamas, deerler deiebilir
    return final_scores
def train_on_batch(self, graph: Data, query_embedding:
    torch.Tensor, labels: torch.Tensor) -> float:

    if self.model is None:
        self.model =
            self._create_model(graph.x.size(1)).to(graph.x.device)
        self.optimizer =
            optim.Adam(self.model.parameters(), lr=2e-4)

    device = next(self.model.parameters()).device

```

```

graph = graph.to(device)
query_embedding = query_embedding.to(device)
labels = labels.to(device)

self.model.train()
node_logits = self.model(graph.x, graph.edge_index)
sim_scores = self._initial_similarity(graph,
    query_embedding)
scores = 0.5 * sim_scores + 0.5 *
    torch.sigmoid(node_logits)

loss = self.loss_fn(node_logits, labels)
self.optimizer.zero_grad()
loss.backward() #binary cross-entropy loss backprop
    yaplyor
self.optimizer.step()
return loss.item()

```

```
class ContextRefiner:
```

```
"""
```

```
son filtremiz
```

```
"""
```

```
def __init__(self):
```

```
pass
```

```
def refine_context(self,
```

```
documents: List[str],
```

```
entities: List[str],
```

```
entity_scores: torch.Tensor,
```

```
top_k: int = 10) -> str:
```

```
if not entities or not documents:
```

```
return " ".join(documents)
```

```
# top-k entity indices
```

```
scores = entity_scores.view(-1).tolist()
```

```
sorted_pairs = sorted(zip(entities, scores),
```

```
key=lambda x: x[1], reverse=True)
```

```
top_entities = [entity for entity, _ in
```

```
sorted_pairs[:top_k]]
```

```
all_sentences = []
```

```
for doc in documents:
```

```
sentences = sent_tokenize(doc)
```

```
all_sentences.extend(sentences)
```

```
relevant_sentences = []
```

```
for sentence in all_sentences:
```

```
if any(entity.lower() in sentence.lower() for
```

```
entity in top_entities):
```

```
relevant_sentences.append(sentence)
```

```

        # hi cmle yoksa dokmanu returnle
        if not relevant_sentences:
            return " ".join(documents)

    return " ".join(relevant_sentences)

class LanguageModel:

    def __init__(self, model_name="google/flan-t5-base",
                  max_length=512):

        self.tokenizer =
            AutoTokenizer.from_pretrained(model_name)
        self.model =
            AutoModelForSeq2SeqLM.from_pretrained(model_name)
        self.max_length = max_length

        self.device = torch.device("cuda" if
            torch.cuda.is_available() else "cpu")
        self.model.to(self.device)

    def generate(self, context: str, query: str) -> str:

        prompt = f"Context: {context}\n\nQuestion:
            {query}\n\nAnswer:"

        # Tokenize
        inputs = self.tokenizer(prompt,
            return_tensors="pt").to(self.device)

        # Generate response
        with torch.no_grad():
            output = self.model.generate(
                inputs["input_ids"],
                max_length=self.max_length,
                num_return_sequences=1,
                temperature=0.7,
                do_sample=True,
                top_p=0.95,
                pad_token_id=self.tokenizer.eos_token_id
            )

        # decoder modelinden answer generatele
        response = self.tokenizer.decode(output[0],
            skip_special_tokens=True)

        # Extract just the answer part (after "Answer:")
        if "Answer:" in response:
            response = response.split("Answer:")[1].strip()

    return response

```

```

class RAG_GNN_System:
    """
        tm sistem
    """

    def __init__(self,
                  embedding_model_name='all-MiniLM-L6-v2',
                  llm_model_name=None):

        self.retriever =
            DocumentRetriever(embedding_model_name)

        self.graph_builder =
            KnowledgeGraphBuilder(embedding_model_name=embedding_model_name,
                                   # graph olutur
            self.gnn_reasoner = GNNReasoner()
            self.context_refiner = ContextRefiner()

        # decoder tanmla
        self.language_model = None
        if llm_model_name:
            self.language_model =
                LanguageModel(model_name=llm_model_name)

    def add_documents(self, documents: List[str]):
        self.retriever.add_documents(documents)

    def get_refined_context(self, query: str, top_k_docs=5,
                           top_k_entities=10) -> str:

        retrieved_docs =
            self.retriever.retrieve_documents(query,
                top_k=top_k_docs)
        if not retrieved_docs:
            return ""

        entities =
            self.graph_builder.extract_entities(retrieved_docs)
        print("Extracted entities:", entities)
        if not entities:
            return " ".join(retrieved_docs)

        relations =
            self.graph_builder.extract_relations(retrieved_docs,
                entities)

        graph = self.graph_builder.build_graph(entities,
            relations)

```

```

query_embedding = self.retriever.get_embedding(query)

entity_scores = self.gnn_reasoner.reason(graph,
    query_embedding) #GNN kts

refined_context =
    self.context_refiner.refine_context(
        retrieved_docs, entities, entity_scores,
        top_k=top_k_entities
    )

return refined_context

def generate_response(self, query: str,
    external_llm=None) -> str:
    """
    contexti decoder a verip cevap ret
    """
    refined_context = self.get_refined_context(query)

    if external_llm:
        return
        external_llm.generate(context=refined_context,
            query=query)
    elif self.language_model:
        return
        self.language_model.generate(context=refined_context,
            query=query)
    else:
        return f"Based on the refined
            context:\n{refined_context}\n\nThe answer
            to your query is: [LLM generation would go
            here]"

import matplotlib.pyplot as plt
import time
import torch
import shutil
import os
from huggingface_hub import Repository
from sentence_transformers import SentenceTransformer

if __name__ == '__main__':
    device = torch.device("cuda" if
        torch.cuda.is_available() else "cpu")
    print("Total train examples:", len(dataset["train"]))
    train_split = dataset["train"].select(range(15000))

    embedding_model_name = 'all-MiniLM-L6-v2'

retriever =
    DocumentRetriever(embedding_model_name=embedding_model_name)
graph_builder =
    KnowledgeGraphBuilder(embedding_model_name=embedding_model_name)
reasoner = GNNReasoner()

num_epochs = 5
train_losses = []

for epoch in range(num_epochs):
    start_time = time.time()
    total_loss = 0.0
    count = 0

    for ex in train_split:
        ctx = ex['context']
        titles = ctx['title']
        sent_paragraphs = ctx['sentences']
        docs = [" ".join(par) for par in sent_paragraphs]

        entities = graph_builder.extract_entities(docs)
        if not entities:
            continue
        relations =
            graph_builder.extract_relations(docs,
                entities) #yeni data yeni graph olutur
        if not relations: #relation yoksa zaman kayb
            continue
        graph = graph_builder.build_graph(entities,
            relations).to(device)

        labels = torch.zeros((len(entities), 1),
            device=device)
        for title, sid in
            zip(ex['supporting_facts']['title'],
                ex['supporting_facts']['sent_id']):
            if title not in titles:
                continue
            p_idx = titles.index(title)
            sentences = sent_paragraphs[p_idx]
            if 0 <= sid < len(sentences):
                sent = sentences[sid]
                for i, ent in enumerate(entities):
                    if ent.lower() in sent.lower():
                        labels[i] = 1.0

        q_emb_np =
            retriever.get_embedding(ex['question'])
        q_emb = torch.tensor(q_emb_np,
            dtype=torch.float, device=device)

        loss = reasoner.train_on_batch(graph, q_emb,
            labels) #her epochda iteration
        total_loss += loss
        count += 1

```

```

    avg_loss = total_loss / count if count > 0 else 0
    train_losses.append(avg_loss)

    duration = time.time() - start_time
    if count > 0:
        print(f"Epoch {epoch+1}/{num_epochs} Loss:
              {avg_loss:.4f} Time: {duration:.2f}s")
    else:
        print(f"Epoch {epoch+1}/{num_epochs} No valid
              examples, skipping.")

# hugging face e pushla

try:
    local_dir = os.path.join(os.getcwd(), "hf_repo")

    repo_id = "rdenaldemir/graph-rag"
    repo = Repository(
        local_dir=local_dir,
        clone_from=repo_id,
        use_auth_token=True
    )
except NameError:
    print("Warning")

    pass
except Exception as e:
    print(f"Warning: {e}")
    pass

if 'repo' in locals() and 'local_dir' in locals():
    ckpt_name = f"gnn_reasoner_spacy.pt"
    os.makedirs(local_dir, exist_ok=True)
    torch.save(reasoner.model.state_dict(), ckpt_name)
    shutil.move(ckpt_name, os.path.join(local_dir,
                                          ckpt_name))
    repo.push_to_hub(commit_message=f"Add checkpoint
                              {ckpt_name}")
    print(f"    Saved trained GNNReasoner weights to
          {os.path.join(local_dir, ckpt_name)}")

# Plot training loss
plt.figure(figsize=(10, 5))
plt.plot(range(1, num_epochs+1), train_losses,
         marker='o', label='Training Loss')
plt.xlabel("Epoch")
plt.ylabel("Loss")
plt.title("Training Loss Over Epochs")
plt.grid(True)
plt.legend()
plt.tight_layout()
plt.show()

import torch

```

```

def load_trained_reasoner(model_path: str,
                          retriever: DocumentRetriever,
                          device: torch.device) ->
                          GNNReasoner:
    """
    model tekrar kullan
    """

    reasoner = GNNReasoner()

    emb = retriever.get_embedding("dummy")
    input_dim = emb.shape[0]

    reasoner.model =
        reasoner._create_model(input_dim).to(device)
    reasoner.model.load_state_dict(torch.load(model_path,
                                              map_location=device))
    reasoner.model.eval()
    return reasoner

from huggingface_hub import hf_hub_download
ckpt_path = hf_hub_download(
    repo_id="rdenaldemir/graph-rag",
    filename="gnn_reasoner_spacy.pt",
    repo_type="model"
)
#spacy iyi
#deneme untrained model
#.pt iyi

device = torch.device("cuda" if torch.cuda.is_available()
                      else "cpu")
rag = RAG_GNN_System()
rag.add_documents([])
retriever = rag.retriever

loaded_reasoner = load_trained_reasoner(
    model_path=ckpt_path,
    retriever=retriever,
    device=device
)
rag.gnn_reasoner = loaded_reasoner

correct = 0
total = 0

eval_split = dataset["validation"].select(range(200))
#validation set
TOP_K_DOCS = 1

```

```

for ex in eval_split:
    total += 1
    question = ex["question"]
    ctx = ex["context"]
    # her paragraf bi dokman
    docs = [" ".join(par) for par in ctx["sentences"]]
    doc_titles = ctx["title"]

    # paragraflar birikmesin diye resetle
    rag.retriever.documents = []
    rag.retriever.document_embeddings = None
    rag.add_documents(docs)

    # retrieve top-K
    retrieved = rag.retriever.retrieve_documents(question,
        top_k=TOP_K_DOCS)

    retrieved_titles = {
        doc_titles[docs.index(doc)]
        for doc in retrieved
        if doc in docs
    }

    # gold titles
    gold_titles = set(ex["supporting_facts"]["title"])

    if gold_titles & retrieved_titles:
        correct += 1

recall_at_k = correct / total * 100
print(f"Document Recall@{TOP_K_DOCS}: {correct}/{total} =
    {recall_at_k:.2f}%")

correct = 0
total = 0

eval_split = dataset["validation"].select(range(200))

TOP_K_DOCS = 10
TOP_K_ENTITIES = 5

for ex in eval_split:
    question = ex["question"]
    ctx = ex["context"]

    all_docs = [" ".join(par) for par in ctx["sentences"]]
    rag.retriever.documents = []
    rag.retriever.document_embeddings = None
    rag.retriever.add_documents(all_docs)

    first_pass = rag.retriever.retrieve_documents(question,
        top_k=TOP_K_DOCS)

```

```

# first_pass = rag.retriever.retrieve_documents(question,
#
#     top_k=TOP_K_DOCS)

# rag.retriever.documents = []
# rag.retriever.document_embeddings = None
# rag.retriever.add_documents(first_pass)

entities =
    rag.graph_builder.extract_entities(first_pass)
relations =
    rag.graph_builder.extract_relations(first_pass,
        entities)
graph = rag.graph_builder.build_graph(entities,
    relations).to(device)

q_emb_np = rag.retriever.get_embedding(question)
q_emb = torch.tensor(q_emb_np, device=device)
scores = rag.gnn_reasoner.reason(graph,
    q_emb).view(-1).cpu().tolist()

ranked = sorted(zip(entities, scores),
    key=lambda x: x[1],
    reverse=True)
top_entities = [e for e, _ in ranked[:TOP_K_ENTITIES]]

gold_sets = []
for title, sent_id in
    zip(ex["supporting_facts"]["title"],
        ex["supporting_facts"]["sent_id"]):

    if title not in ctx["title"]:
        gold_sets.append(set())
        continue

    p_idx = ctx["title"].index(title)
    sentence = ctx["sentences"][p_idx][sent_id]

    ents_in_sent = {
        ent for ent in entities
        if ent.lower() in sentence.lower()
    }
    gold_sets.append(ents_in_sent)

hit_both = all(
    any(ent in top_entities for ent in gold_set)
    for gold_set in gold_sets
)

if hit_both:
    correct += 1
total += 1

```

```

print(f" bothsupportfact hits@{TOP_K_ENTITIES}: "
      f"{correct}/{total} = {correct/total:.2%}")

!pip install evaluate

from evaluate import load as load_metric

metric = load_metric("squad")

predictions = []
references = []

rag_gnn =
    RAG_GNN_System(llm_model_name="google/flan-t5-large")
rag_gnn.gnn_reasoner = loaded_reasoner

for ex in eval_split:
    q_id = ex["id"]
    question = ex["question"]

    docs = [" ".join(par) for par in
            ex["context"]["sentences"]]

    # dokman resetle birikmesin
    rag_gnn.retriever.documents = []
    rag_gnn.retriever.document_embeddings = None
    rag_gnn.add_documents(docs)

    refined_ctx =
        rag_gnn.get_refined_context(ex["question"],
                                    top_k_docs=5, top_k_entities=15)
    pred_answer =
        rag_gnn.language_model.generate(refined_ctx,
                                         question)

    predictions.append({
        "id": q_id,
        "prediction_text": pred_answer
    })
    references.append({
        "id": q_id,
        "answers": {
            "text": [ex["answer"]], #
            "answer_start": [0] #
        }
    })

```

```

results = metric.compute(predictions=predictions,
                          references=references)

print(f"Exact Match: {results['exact_match']:.2f}")
print(f"F1 Score : {results['f1']:.2f}")

from evaluate import load as load_metric
import pandas as pd

model_names = [
    "google/flan-t5-base",
    "google/flan-t5-large",
    # "google/flan-t5-xl",
]

k_docs_list = [3, 5, 7, 10]
k_ent_list = [5, 10, 20]

metric = load_metric("squad")
records = []

for model_name in model_names:

    rag = RAG_GNN_System(llm_model_name=model_name)
    rag.gnn_reasoner = loaded_reasoner

    for k_docs in k_docs_list:
        for k_ent in k_ent_list:
            preds, refs = [], []

            for ex in eval_split:
                # reset
                paras = ex["context"]["sentences"]
                docs = [" ".join(p) for p in paras]
                rag.retriever.documents = []
                rag.retriever.document_embeddings = None
                rag.add_documents(docs)

                # retrieve, refine, generate
                refined = rag.get_refined_context(
                    ex["question"],
                    top_k_docs=k_docs,
                    top_k_entities=k_ent
                )
                answer =
                    rag.language_model.generate(refined,
                                                ex["question"])

                preds.append({"id": ex["id"],
                             "prediction_text": answer})
                refs.append({
                    "id": ex["id"],

```

```

        "answers": {"text": [ex["answer"]],
                    "answer_start": [0]}
    })

    res = metric.compute(predictions=preds,
                        references=refs)
    records.append({
        "model":      model_name.split("/")[-1],
        "top_k_docs": k_docs,
        "top_k_ent":  k_ent,
        "EM":         res["exact_match"],
        "F1":         res["f1"]
    })

df = pd.DataFrame(records)
print(df.pivot_table(
    index=["model", "top_k_docs"],
    columns="top_k_ent",
    values="F1"
))

predictions = []
references = []
# VAHYSEL RETRIEVAL
for ex in eval_split:
    q_id = ex["id"]
    question = ex["question"]

    ctx = ex["context"]
    titles = ctx["title"]
    paras = ctx["sentences"]

    support_sents = []
    for title, sent_id in
        zip(ex["supporting_facts"]["title"],
            ex["supporting_facts"]["sent_id"]):
        if title in titles:
            p_idx = titles.index(title)

            support_sents.append(paras[p_idx][sent_id])

    oracle_context = " ".join(support_sents)
    # -----

    pred_answer =
        rag_gnn.language_model.generate(oracle_context,
        question)

    # 4) Append for metric
    predictions.append({
        "id": q_id,
        "prediction_text": pred_answer
    })

```

```

    references.append({
        "id": q_id,
        "answers": {
            "text": [ex["answer"]],
            "answer_start": [0]
        }
    })

    results = metric.compute(predictions=predictions,
                            references=references)
    print(f"Oracle EM: {results['exact_match']:.2f}")
    print(f"Oracle F1: {results['f1']:.2f}")

    from evaluate import load as load_metric

    rag_plain =
        RAG_GNN_System(llm_model_name="google/flan-t5-large")

    rag_plain.gnn_reasoner = None

    metric = load_metric("squad")

    preds_plain = []
    refs = []
    count = 0

    for ex in eval_split:
        count += 1
        q_id = ex["id"]
        question = ex["question"]
        gold_ans = ex["answer"]

        docs = [" ".join(p) for p in ex["context"]["sentences"]]

        # Reset
        rag_plain.retriever.documents = []
        rag_plain.retriever.document_embeddings = None
        rag_plain.add_documents(docs)

        top_k = 5
        retrieved =
            rag_plain.retriever.retrieve_documents(question,
            top_k)

        context = " ".join(retrieved)

        pred_answer =
            rag_plain.language_model.generate(context,
            question)

```



```

preds_plain.append({
    "id": q_id,
    "prediction_text": pred_answer
})
refs.append({
    "id": q_id,
    "answers": {"text": [gold_ans], "answer_start": [0]}
})
print(count)

```

```

res_plain = metric.compute(predictions=preds_plain,
    references=refs)

```

```

print(f"      Plain RAG      \nExact Match:
      {res_plain['exact_match']:.2f}\nF1 Score :
      {res_plain['f1']:.2f}")

```

```

print(f" GNNRAG      EM: {results['exact_match']:.2f}, F1:
      {results['f1']:.2f}")

```

```

print(f"Plain RAG      EM: {res_plain['exact_match']:.2f},
      F1: {res_plain['f1']:.2f}")

```

```

from torch_geometric.utils import to_networkx
import networkx as nx
import matplotlib.pyplot as plt

```

```

"""
GRAFK      ZDR
"""

```

```

example = dataset["validation"][1020]
question = example["question"]
answer = example["answer"]
all_docs = [" ".join(p) for p in
    example["context"]["sentences"]]

```

```

print("Question:\n", question, "\n Answer:", answer)

```

```

TOP_K = 2
retriever =
    DocumentRetriever(embedding_model_name="all-MiniLM-L6-v2")
retriever.add_documents(all_docs)
docs = retriever.retrieve_documents(question, top_k=TOP_K)

```

```

print(f"Top-{TOP_K} Retrieved Documents:")
for i, doc in enumerate(docs):
    print(f"[Doc {i}]: {doc}\n")

```

```

builder = KnowledgeGraphBuilder(
    embedding_model_name="all-MiniLM-L6-v2",
    spacy_model="en_core_web_lg"
)

```

```

entities = builder.extract_entities(docs)
relations = builder.extract_relations(docs, entities)
graph = builder.build_graph(entities, relations)

```

```

import networkx as nx
import matplotlib.pyplot as plt
from torch_geometric.utils import to_networkx

```

```

G = nx.DiGraph()
for idx, ent in enumerate(entities):
    G.add_node(idx, label=ent)
for src, tgt, rel in relations:
    G.add_edge(src, tgt, rel=rel)

pos = nx.spring_layout(G, k=1.2, seed=42)
plt.figure(figsize=(14,10))
nx.draw_networkx_nodes(G, pos, node_size=800,
    node_color='skyblue', alpha=0.9)
nx.draw_networkx_edges(G, pos, arrowstyle='->',
    arrowsize=12, edge_color='gray', width=1.2)
nx.draw_networkx_labels(G, pos,
    nx.get_node_attributes(G, 'label'),
        font_size=8,
        bbox=dict(facecolor='white',
            edgecolor='none', alpha=0.8))
nx.draw_networkx_edge_labels(G, pos,
    edge_labels=nx.get_edge_attributes(G, 'rel'),
        font_size=7, label_pos=0.5,
        rotate=False)
plt.title(f"Top-{TOP_K} EntityRelation Graph")
plt.axis('off')
plt.tight_layout()
plt.show()

```