*Article*

# Enhancing Retrieval-Augmented Generation with Entity Linking for Educational Platforms

**Francesco Granata** [1], **Francesco Poggi** [2], **Misael Mongiovì** [1,2]

1   Department of Mathematics and Computer Science, University of Catania, Italy;
    grnfnc03t17c351z@studium.unict.it, misael.mongiovi@unict.it
2   Institute of Cognitive Sciences and Technologies (ISTC), National Research Council of Italy (CNR);
    francesco.poggi@cnr.it

**Abstract**

In the era of Large Language Models (LLMs), Retrieval-Augmented Generation (RAG) architectures are gaining significant attention for their ability to ground language generation in reliable knowledge sources. Despite their impressive effectiveness in many areas, RAG systems based solely on semantic similarity often fail to ensure factual accuracy in specialized domains, where terminological ambiguity can affect retrieval relevance. This study proposes an enhanced RAG architecture that integrates a factual signal derived from Entity Linking to improve the accuracy of educational question-answering systems in Italian. The system includes a Wikidata-based Entity Linking module and implements three re-ranking strategies to combine semantic and entity-based information: a hybrid score weighting model, reciprocal rank fusion, and a cross-encoder re-ranker. Experiments were conducted on two benchmarks: a custom academic dataset and the standard SQuAD-it dataset. Results show that, in domain-specific contexts, the hybrid schema based on reciprocal rank fusion significantly outperforms both the baseline and the cross-encoder approach, while the cross-encoder achieves the best results on the general-domain dataset. These findings confirm the presence of an effect of domain mismatch and highlight the importance of domain adaptation and hybrid ranking strategies to enhance factual precision and reliability in retrieval-augmented generation. They also demonstrate the potential of entity-aware RAG systems in educational environments, fostering adaptive and reliable AI-based tutoring tools.

**Keywords:** Retrieval-Augmented Generation; Entity Linking; Domain Adaptation; Wikidata; Educational Question Answering

## 1. Introduction

In the current technological landscape, dominated by *Large Language Models (LLMs)*, the way we approach information and data has profoundly changed. These models are capable of producing fluent and context-aware text across a wide range of domains, making them powerful tools for knowledge access, task automation, and decision support. Their impressive generative capabilities, however, come with inherent limitations. Despite these strengths, LLMs are not always the best way to access knowledge, especially in critical, specialized or high-precision domains, because they can sometimes provide incorrect, inconsistent or unverifiable information. This phenomenon - commonly referred to as *hallucination* [1] - poses significant risks when accuracy and reliability are essential, such as in scientific research [2], medicine [3,4], and education [5,6].

To address this problem and increase factual consistency, *Retrieval-Augmented Generation (RAG)* architectures have proven to be one of the most effective approaches. By combining the generative abilities of LLMs with external knowledge sources, RAG systems can retrieve relevant documents or passages and ground the model's output in real, verifiable information. This allows them to provide precise and up-to-date answers without requiring model retraining or large-scale parameter updates. However, these systems are not perfect: most retrieval components rely heavily on semantic similarity, which can fail to adequately capture subtle distinctions in technical language, polysemous terms, or domain-specific terminology. As a result, retrieval errors propagate into the generation phase, reducing the overall reliability of the system.

These limitations become particularly evident in the educational domain, where materials often contain highly specialized vocabulary, hierarchical conceptual structures, and terminology that differs across subjects or instructional levels. Standard RAG approaches may struggle to identify the correct referents for ambiguous or overlapping terms, leading to suboptimal retrieval and factually incorrect responses. To improve performance in this context we explore the integration between *Entity Linking* and RAG architectures. Entity Linking provides a structured mechanism for identifying and disambiguating mentions of concepts within text, mapping them to unique identifiers in a knowledge base. This additional layer of semantic grounding has the potential to refine retrieval queries, reduce ambiguity, and enhance the factual stability of generated outputs. The aim of this work is to investigate how Entity Linking can enhance retrieval accuracy and factual reliability in RAG systems applied to educational content.

## 1.1. Related Work

To improve the factual grounding of large language models (LLMs), several studies have examined *Retrieval-Augmented Generation* (RAG) systems as a promising framework to address this open research challenge. The original architecture proposed by Lewis et al. (2020) [7] integrates a neural retriever with a generative model, enabling language models to access external knowledge sources dynamically. A standard RAG pipeline can be conceptually divided into two main components: a retrieval phase and a generation phase. The retrieval component is often based on semantic similarity between query and document passages, although alternative approaches have been explored. For example, Mongiovì et al. (2024) [8] introduced *GRAAL*, a graph-based retrieval system that collects related passages across multiple documents by exploiting inter-entity relationships. A wide range of retrieval strategies have been proposed to improve access to relevant information; however, many still face challenges in capturing fine-grained semantic meaning, particularly in specialized domains that contain ambiguous or domain-specific terminology.

*Entity Linking* (EL) has been proposed as a method to improve the disambiguation of named entities by aligning textual mentions with entries in structured knowledge bases such as Wikidata [9] or DBpedia. Shen et al. (2015) [10] provided an extensive overview of EL challenges and techniques. Möller et al. (2022) [11] presented a comprehensive survey of Wikidata-based EL methods and datasets. Several frameworks have been developed in this area, including BLINK [12], ReLiK [13], and OpenTapioca [14], each adopting different architectures and disambiguation strategies. Most current EL systems, however, have been primarily developed for the English language, which limits their direct applicability in multilingual or non-English domains.

At the same time, *Educational Artificial Intelligence (EAI)* has become a rapidly expanding field of research. AI-based systems are increasingly applied to support learning, personalize content delivery, and facilitate knowledge acquisition. Mageira et al. (2022) [15] discussed the pedagogical potential of educational chatbots designed for content

and language-integrated learning. More recently, Swacha et al.(2025) [16] and Li et al.
(2025) [17] provided a systematic survey of RAG applications in education, illustrating how
retrieval-augmented approaches can enhance information access and automated question
answering in academic contexts.

The present work builds upon these lines of research by combining RAG and Entity
Linking within a hybrid retrieval architecture applied to the educational domain in Italian.
While prior research, such as Shlyk et al. (2024) [18], has explored the integration of RAG
and EL for biomedical concept recognition, our approach differs in both pipeline design
and linguistic scope, focusing on hybrid retrieval in an Italian-language educational setting.

### *1.2. Research Objectives and Contributions*

The objective of this study is to investigate how the integration of an *Entity Linking*
component can enhance the performance of *Retrieval-Augmented Generation* (RAG) systems
in specialized educational contexts. In particular, the research focuses on developing and
evaluating a hybrid retrieval architecture that combines semantic similarity and entity-
based information to improve the relevance and factual reliability of retrieved passages. The
study is conducted entirely on Italian-language educational data, addressing the challenge
of adapting retrieval and generation techniques to a non-English education domain.

The main contributions of this work can be summarized as follows:

1. We design and implement *ELERAG*: a hybrid RAG architecture that integrates a
   Wikidata-based Entity Linking module to incorporate entity-level knowledge during
   retrieval.
2. We evaluate the efficacy of our proposed *RRF-Based Re-ranking* strategy by comparing
   it against a baseline Weighted-Score Re-ranking and a high-complexity RRF + Cross-
   Encoder Re-ranking.
3. We analyze the results across both custom educational data and standard benchmarks,
   highlighting the effects of domain adaptation and linguistic specificity on retrieval
   performance.
4. We provide clear experimental evidence of a *domain mismatch*, demonstrating that a
   domain-adapted hybrid model can outperform a generic State-Of-The-Art re-ranker
   on specialized data, whereas the SOTA model excels on standard benchmarks.

By addressing the intersection of RAG, Entity Linking, and Educational AI, this study
contributes to ongoing efforts to improve factual grounding and domain-sensitive retrieval
in large language model applications. It introduces and analyzes a new class of hybrid
architectures designed to support more trustworthy, transparent and pedagogically aligned
AI-driven educational tools.

This paper is organized as follows: Section 2 details the methodology, describing the
hybrid RAG architecture, the Entity Linking module, and the implemented re-ranking
strategies. Section 3 presents the experimental setup, including the construction of the
custom educational dataset, the benchmarks, and the evaluation metrics. Section 4 reports
the quantitative and qualitative results obtained from the experiments. Finally, Section 5
discusses the findings and draws the main conclusions, highlighting the domain mismatch
phenomenon and outlining limitations and future directions.

## 2. Methodology

This section details the architecture and implementation of *ELERAG* (Entity Linking
Enhanced RAG), the hybrid retrieval system proposed in this study. We first describe the
baseline RAG configuration used as a reference, followed by the Entity Linking module
designed to improve semantic disambiguation. Next, we present the core re-ranking strat-

egy adopted in *ELERAG*, along with alternative strategies implemented for comparative analysis. Finally, we illustrate the complete end-to-end workflow of the proposed method.

### 2.1. RAG baseline

We built a baseline *Retrieval-Augmented Generation* (RAG) system as a starting point for our architecture. This baseline follows the standard structure of embedding-based retrieval combined with generative models, a paradigm introduced by Lewis et al (2020) [7] and widely explored in recent educational applications [16,17]. Our approach combines the vector store *FAISS* [19] with `multilingual-e5-large` [20] as an embedding model and GPT-4o [21] as a generator.

Each chunk in the corpus was encoded into a fixed-length 1024-dimensional vector representation using the `multilingual-e5-large` model by *Sentence Transformers* [1], chosen for its strong multilingual retrieval capabilities and competitive performance on cross-lingual benchmarks. The resulting embeddings were normalized and stored in a *FAISS* index employing an inner-product [2] similarity metric to enable efficient dense retrieval at query time. During inference, a user query is first encoded with the same embedding model, then the top-*K* most similar vectors are probed and the corresponding chunks are retrieved. These retrieved passages are concatenated to form a context window, which is then provided as additional input to the generative model.

We used *GPT-4o* as the generator, prompting it to answer based solely on the retrieved content and to abstain when sufficient evidence is not present in the context. Crucially, the prompt explicitly instructs the model to cite the source chunk IDs for every piece of information used in the answer. Consequently, the generation phase acts as a final refinement stage: only the chunks actually cited in the generated response are considered "retrieved" by the full system, while uncited chunks—even if present in the context window—are effectively discarded as irrelevant. This configuration serves as the reference setup for evaluating the contribution of entity-level enrichment and alternative ranking strategies in the following sections.

### 2.2. Entity Linking Module

While the baseline RAG system relies solely on semantic similarity for retrieval, this approach often struggles with domain-specific ambiguity and the presence of polysemous terms. In educational material, where concepts may appear with slight linguistic variations across disciplines, pure embedding similarity can retrieve semantically close but contextually irrelevant chunks. To address this limitation, we integrated an *Entity Linking* (EL) module designed to ground text spans to canonical entities within *Wikidata* [9].

Every chunk in the dataset was pre-processed to extract named entities using the *SpaCy* pipeline [22]. Specifically, we employed the `it_core_news_lg` [3] model, a large pre-trained pipeline for the Italian language trained on a massive corpus of news and media text, chosen for its superior accuracy in Named Entity Recognition (NER) compared to smaller variants.

Before developing this custom solution, we experimented with standard state-of-the-art Entity Linking systems such as *BLINK* [23]. However, since these models are primarily optimized for English, they yielded unsatisfactory results when applied to our Italian educational corpus. Consequently, we opted for a lightweight, API-based approach tailored to our specific language requirements. For each detected entity mention, a candidate list of Wikidata entities was retrieved through the public Wikidata API.

---

[1] https://sbert.net
[2] for normalized vectors it is equivalent to cosine similarity
[3] https://spacy.io/models/it#it_core_news_lg

To select the best candidate, a hybrid scoring function was developed combining two signals:

1. *Popularity* — computed as the inverse of the candidate rank in the list returned by Wikidata:

$$\text{popularity} = \frac{1}{\text{rank} + 1} \tag{1}$$

2. *Semantic similarity* — obtained using `multilingual-e5-large`, applied to the mention context (the sentence where the entity was detected) and the concatenation of the candidate's *label* and its *description*.

The final score is computed as:

$$\text{HybridScore} = \alpha \cdot \text{similarity} + (1 - \alpha) \cdot \text{popularity} \tag{2}$$

where $\alpha$ ($= 0.9$) controls the balance between semantic and popularity signals, a value empirically chosen after a preliminary tuning phase. The candidate with the highest score is selected as the final linked entity and stored, together with all metadata and supporting information, in a JSON structure with the addition of linked entities.

### 2.3. Ranking Strategies and Integration into the Retrieval Pipeline

The enrichment introduced by the *Entity Linking* module was used to enhance the retrieval phase by re-ranking the initially retrieved chunks from the FAISS index. Specifically, after dense retrieval, an additional entity-aware re-ranking stage was applied.

Proposed Strategy: RRF-Based Re-ranking.

The core ranking strategy adopted in our *ELERAG* architecture is based on *Entity-Aware Reciprocal Rank Fusion (RRF)*. In this configuration, chunks are independently ranked according to their dense score (semantic similarity) and their entity score (factual overlap). The two distinct rankings are then fused using the RRF algorithm [24], which assigns a joint score based on the rank position in each list:

$$\text{score}_{\text{RRF}} = \frac{1}{K + \text{rank}_{\text{dense}}} + \frac{1}{K + \text{rank}_{\text{entity}}} \tag{3}$$

where $K = 60$ is the standard smoothing constant. This approach was selected as the primary method for *ELERAG* because it robustly balances semantic relevance with factual entity matching without requiring the manual tuning of weights or the high computational cost of cross-encoders. As demonstrated in Section 4, this strategy yielded the best performance on our specialized educational dataset.

Comparative Strategies.

To rigorously evaluate the effectiveness of our proposed strategy, we implemented two alternative re-ranking methods for comparison:

1. *Weighted-Score Re-ranking.* Each chunk was evaluated using a combined score:

$$\text{final\_score} = \text{dense\_score} + \beta \cdot \text{entity\_score} \tag{4}$$

where dense_score represents the cosine similarity between the query and the chunk embedding computed by `multilingual-e5-large`, and the entity_score represents

the recall-oriented overlap between the set of query linked entities ($Q_E$) and the set of chunk linked entities ($C_E$):

$$\text{entity\_score} = \frac{|Q_E \cap C_E|}{|Q_E|} \text{ if } |Q_E| > 0 \text{ else } 0 \qquad (5)$$

The hyperparameter $\beta$ controls the relative contribution of entity-based evidence compared to semantic similarity.

2. *RRF + Cross-Encoder Re-ranking.* To benchmark our approach against a pure semantic SOTA method, we implemented a three-stage pipeline. After dense retrieval and RRF fusion, the top candidates are re-scored by a transformer-based Cross-Encoder [25] (`mmarco-mMiniLMv2-L12-H384-v1`). Unlike bi-encoders (like E5), the Cross-Encoder processes the query and document simultaneously, capturing finer semantic nuances at a higher computational cost. For this specific configuration, to maximize initial recall before the expensive scoring step, we retrieved a larger pool of candidates (Top-50) from the vector index and selected the Top-20 final chunks after re-ranking.

In all configurations, the progressive re-ranking across stages reduces the candidate set size, allowing the system to focus on the most promising passages for the subsequent generation phase.

Figure 1 illustrates the complete workflow of our proposed *ELERAG* method. The architecture processes the query in parallel streams—extracting entity-based features and computing dense embeddings—before fusing the results via the RRF module to feed the LLM. The other experimental configurations can be understood as variations of this schema: the *Standard RAG* baseline utilizes only the lower branch (Dense Embedding → Vector Index → LLM), bypassing the Entity Linking and Re-ranking stages. The *Weighted-Score Re-ranking* configuration replaces the RRF block with the linear combination strategy described above. Finally, the *RRF + Cross-Encoder Re-ranking* configuration adds a further refinement block between the RRF stage and the LLM generation.
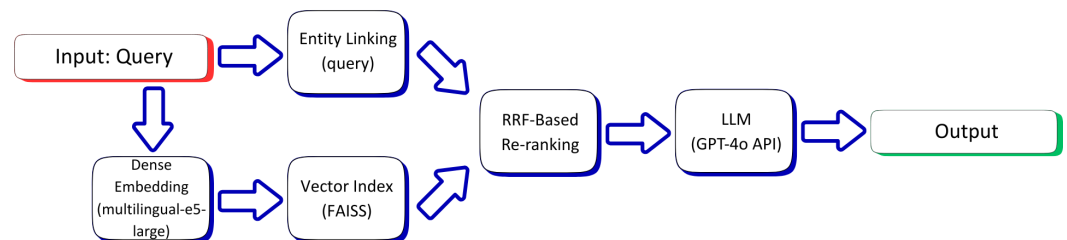


**Figure 1.** Architectural schema of the proposed *ELERAG* method. The system integrates parallel retrieval paths—semantic dense retrieval and entity linking—fusing them via an RRF-based re-ranking module to ground the LLM generation.

## 3. Experimental Setup

This section presents the experimental setup designed to assess the effectiveness of the proposed *ELERAG* architecture. It describes the datasets used for evaluation, both custom and standard, and the criteria adopted to ensure comparability across different domains. Then, it details the metrics employed to quantify retrieval and generation performance, and the evaluation methods applied to each system configuration. Together, these components define a consistent framework for analyzing and comparing the proposed hybrid retrieval models.

### 3.1. Educational Data

The dataset used in this study was constructed from two Italian university courses, namely *Applied Economics* and *Language and Communication*, offered by an Italian telematic

university. Each course consisted of several video lectures, totalling 50 classes over 32 hours of material, which were first transcribed into text using the *Whisper Turbo* automatic speech recognition model [26]. The transcription process also generated precise word-level timestamps, which were later used to associate each text segment with its corresponding temporal interval in the original lecture.

To prepare the textual material for retrieval, the transcribed lectures were segmented into coherent chunks using the SpaCy Python library [22]. To ensure semantic coherence, segmentation was strictly constrained to sentence boundaries: sentences were grouped to form chunks with a length between 20 and 300 tokens. This strategy prevents cutting sentences in the middle, preserving the syntactic and semantic integrity of the text. For each chunk, start and end timestamps were retained. All chunks were stored in JSON format and enriched with named entities automatically extracted by *SpaCy*, which were later used by the EL module described in Section 2.2.

### 3.2. Evaluation setup

The data produced by the procedure in Section 3.1 does not contain question-answer pairs, necessary to evaluate our approach. Therefore, we adopted a dual evaluation strategy that combines a custom benchmark, specifically generated from the course material using the GPT-4o API, and a standard benchmark (SQuAD-it [27]). This setup allows us to assess the system's performance both within a specialized educational context and in a general-domain question answering scenario.

Custom Benchmark.

The custom benchmark was automatically constructed from the lecture corpus discussed in Section 3.1. Using the GPT-4o API, we prompted the model to generate three types of questions—*factual*, *synthesis*, and *inference*—together with their corresponding gold answers and relevant document references. The resulting benchmark contains 69 questions, each represented as a structured record including query, question type, gold answer, relevant documents, and additional metadata.

After generation, the dataset underwent a two-step validation process. First, a secondary GPT-4o prompt was used to verify the correctness and consistency between questions and their corresponding answers. Subsequently, a manual validation was performed by human annotators to identify and filter out ambiguous, ill-formed, or hallucinated questions that might have bypassed the automated check. This rigorous process ensured the high quality of the final question set used for evaluation.

Standard Benchmark (SQuAD-it).

For comparison with a standard task, we employed the SQuAD-it benchmark, a widely used dataset for Italian question answering derived from Wikipedia. The dataset consists of a collection of passages (contexts), and for each passage, a set of question-answer pairs. To adapt it to our retrieval evaluation setting, we treated each Wikipedia passage as a distinct "chunk" and indexed them exactly as we did for the educational dataset. In this configuration, for each query, the passage referenced in the original dataset (the *context*) is treated as the unique *gold answer* to be retrieved.

### 3.3. Evaluation Metrics

The following metrics were used to compare the retrieval and generation performance of the different systems:

- *Exact Match (EM):* Proportion of queries for which the first retrieved document exactly matches the gold answer.

$$\text{EM} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{1}\{\text{retrieved\_docs}_i[0] = \text{gold\_answer}_i\}$$

- *Recall@k:* Proportion of relevant documents that appear among the top-*k* retrieved results.

$$\text{Recall@k} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\{doc \in \text{retrieved\_docs}_i[:k] : doc \in \text{relevant\_docs}_i\}|}{|\text{relevant\_docs}_i|}$$

- *Precision@k:* Proportion of the top-*k* retrieved documents that are relevant.

$$\text{Precision@k} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\{doc \in \text{retrieved\_docs}_i[:k] : doc \in \text{relevant\_docs}_i\}|}{k}$$

- *Mean Reciprocal Rank (MRR):* Mean reciprocal rank of the first occurrence of a relevant document. Two distinct versions were computed:

  - MRR based on the `gold_answer`:

$$\text{MRR\_gold} = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{\text{rank}_i(\text{gold\_answer}_i)}$$

  - MRR based on all relevant documents:

$$\text{MRR\_rel\_docs} = \frac{1}{N} \sum_{i=1}^{N} \frac{1}{\text{rank}_i(\text{first\_relevant\_doc})}$$

- **General Recall and Precision:** Applied in the full RAG configuration. As described in Section 2.1, the LLM acts as a filter by citing only the sources actually used in the answer. Therefore, for these metrics, the set of `retrieved_docs` consists exclusively of the chunks referenced in the final generation, making the number of retrieved documents variable:

$$\text{Recall} = \frac{|\{doc \in \text{retrieved\_docs} : doc \in \text{relevant\_docs}\}|}{|\text{relevant\_docs}|}$$

$$\text{Precision} = \frac{|\{doc \in \text{retrieved\_docs} : doc \in \text{relevant\_docs}\}|}{|\text{retrieved\_docs}|}$$

- *Subjective metrics:* Completeness, relevance, and clarity, each evaluated by `gpt-4o` on a 1–10 scale for the generated answers obtained by concatenating the top-3 retrieved chunks. Although LLM-based evaluation cannot fully replace human judgment, it provides a scalable and replicable proxy for consistently estimating answer quality across the benchmark.

### 3.4. Evaluation methods

Three complementary evaluation methods were considered to comprehensively compare the different variants of the retrieval pipeline:

1. *Method 1: Classical evaluation on dense retrieval.* The metrics EM, Recall@k, Precision@k, MRR_gold, and MRR_rel_docs were computed directly on the retrieval system, with-

out involving the LLM generation step. This method was applied to both the custom benchmark and the *SQuAD-it* dataset.

2. *Method 2: Subjective evaluation through LLM-based scoring.* To assess the quality of the generated text, final answers were produced by the RAG system using the top-*K* retrieved chunks as context. For this specific evaluation, *K* was set to 3 across all configurations to ensure comparable context lengths. Subsequently, a separate `gpt-4o` instance was prompted to act as an external evaluator, scoring the generated answers on a 1–10 scale according to three criteria: completeness, relevance, and clarity.

3. *Method 3: Classical evaluation on the full RAG pipeline.* In this case, the query was processed by the complete RAG system, where the LLM filters and synthesizes the retrieved documents. The computed metrics included EM, Recall, Precision, MRR_gold, and MRR_rel_docs. Unlike the dense retrieval case, @k-based metrics were not used, since the number of documents returned by the LLM is not fixed.

This evaluation framework provides the basis for the experimental analysis presented in the next section.

## 4. Results

This section presents the results obtained from the evaluation of the different pipelines. We compare the *Standard RAG* baseline against the variants incorporating Entity Linking, specifically focusing on the performance of our proposed *ELERAG* method relative to the alternative re-ranking strategies defined in Section 2:

- *ELERAG (RRF-Based Re-ranking)*: Our proposed method, applied to the Top-30 candidate chunks retrieved by FAISS.
- *Weighted-Score Re-ranking*: Implemented with a weighting factor $\beta = 0.5$, applied to the Top-30 candidate chunks retrieved by FAISS.
- *RRF + Cross-Encoder Re-ranking*: This method is applied to re-score the Top-20 candidates selected via RRF from an extended initial pool of 50 chunks retrieved by FAISS.

The quantitative and qualitative findings are reported according to the three evaluation methods described in the previous section.

### 4.1. Method 1: Classical Evaluation on Retrieval

Table 1 summarizes the performance of the retrieval pipelines on the custom educational benchmark. To facilitate comparison, the best result for each metric is highlighted in bold.

**Table 1.** Retrieval performance on the custom benchmark. The proposed *RRF-Based Re-ranking* strategy achieves the best results in Exact Match (EM) and Precision@1, surpassing both the baseline and the Cross-Encoder.

| Pipeline | EM | R@1 | R@3 | R@5 | R@10 | P@1 | MRR_G | MRR_RD |
|---|---|---|---|---|---|---|---|---|
| Baseline | 0.522 | 0.377 | 0.640 | **0.729** | **0.807** | 0.652 | 0.652 | 0.759 |
| Weighted-Score | 0.536 | 0.391 | **0.647** | 0.717 | 0.795 | 0.681 | 0.654 | 0.772 |
| **ELERAG (RRF)** | **0.565** | 0.399 | **0.647** | 0.725 | 0.795 | **0.696** | **0.668** | **0.779** |
| RRF+Cross-Encoder | 0.536 | **0.408** | 0.645 | 0.698 | 0.802 | 0.652 | 0.647 | 0.760 |

Analysis of Results.

The experimental data confirms that integrating Entity Linking significantly improves retrieval precision. All entity-enhanced pipelines outperformed the baseline in terms of Exact Match (EM), indicating a superior ability to place the precise gold chunk at the very top of the ranking.

*Comparison of Ranking Strategies.* Our proposed *ELERAG* strategy emerged as the most effective method for this domain, achieving the highest scores in EM (0.565), Precision@1 (0.696), and MRR_gold (0.668) and MRR_rel_docs (0.779). This superiority demonstrates that fusing the entity signal with dense retrieval via RRF is more robust than a simple linear combination (Weighted-Score), effectively filtering out semantically similar but factually incorrect chunks.

*Cross-Encoder vs. Entity Signal.* Comparing our method with the *RRF + Cross-Encoder* reveals a critical insight consistent with the Domain Mismatch hypothesis. The Cross-Encoder achieved the highest Recall@1 (0.408), proving its strong capability to identify relevant content in a broad sense. However, it failed to translate this into superior Exact Match or MRR scores, performing worse than *ELERAG* on these precision metrics. This suggests that while the generic, pre-trained Cross-Encoder captures general semantics well, it lacks the specific domain grounding provided by the explicit Entity Linking signal. Consequently, *ELERAG* proves to be not only more accurate for the top-ranking position but also significantly more computationally efficient than the Cross-Encoder pipeline.

*Recall Trade-off.* It is worth noting that the Baseline maintains higher Recall@5 and Recall@10 scores compared to the re-ranking methods. This is a typical behavior of re-ranking architectures, which aggressively optimize the top positions (Top-1 to Top-3) to maximize precision for the limited context window of the LLM. While this may push some marginally relevant documents further down the list (slightly affecting R@10), the substantial gain in MRR and EM is far more valuable for a RAG system, where the generator's performance depends primarily on the quality of the very first retrieved chunks.

### 4.2. Method 2: Subjective Evaluation through LLM-based Scoring

Table 2 reports the average scores assigned by the evaluator LLM to the answers generated by each pipeline. The evaluation focuses on three dimensions: Completeness, Relevance, and Clarity.

**Table 2.** Subjective evaluation scores (scale 1–10). *ELERAG* achieves the highest scores in all categories, confirming that entity-aware retrieval leads to more comprehensive and relevant answers.

| Pipeline | Completeness | Relevance | Clarity |
|---|---|---|---|
| Baseline | 5.99 | 5.45 | 4.51 |
| Weighted-Score | 6.00 | 5.43 | 4.42 |
| **ELERAG (RRF)** | **6.10** | **5.57** | **4.54** |
| RRF+Cross-Encoder | 5.94 | 5.39 | 4.43 |

*Analysis of Response Quality.* Consistent with the retrieval metrics observed in Method 1, the *ELERAG (RRF)* strategy achieves the highest scores across all qualitative dimensions. The improvement is particularly noticeable in *Completeness* (6.10) and *Relevance* (5.57). This indicates that the chunks prioritized by the RRF algorithm—by balancing semantic similarity and entity popularity—provide the generator with a richer and more pertinent context, allowing it to construct answers that are not only factually correct but also more exhaustive.

*Comparison with Cross-Encoder.* Interestingly, the *RRF+Cross-Encoder* pipeline scores slightly lower than both *ELERAG* and the *Baseline* in terms of qualitative scoring. This reinforces the hypothesis derived from the quantitative analysis: although the Cross-Encoder is a powerful semantic filter, in this specific educational domain it may be overly aggressive, discarding chunks that contain necessary context or entity definitions that are valuable for the final generation. Consequently, the simpler and more robust RRF fusion proves to be superior for end-to-end answer quality.

*Clarity Stability.* The scores for *Clarity* are relatively stable across all systems (ranging from 4.42 to 4.54). Since the evaluation target in this method is formed by directly concatenating the raw retrieved chunks, this stability is expected: the source material (lecture transcripts) is identical for all pipelines. However, the slight edge held by *ELERAG* (4.54) suggests that the entity-driven retrieval tends to select more structurally complete or self-contained passages, resulting in a slightly more coherent reading flow compared to the other methods.

### 4.3. Method 3: Classical Evaluation on the Full RAG Pipeline

Table 3 presents the performance metrics calculated on the final output of the end-to-end RAG system. In this evaluation, the metrics consider only the chunks explicitly cited by the LLM in the generated answer, effectively treating the generator as a final relevance filter.

**Table 3.** Performance of the full RAG system (post-LLM filtering). *ELERAG* maintains its superiority, achieving the highest Exact Match and MRR, proving that better initial ranking translates to better generation support.

| Pipeline | EM | Recall | Precision | MRR_G | MRR_RD |
|---|---|---|---|---|---|
| Baseline | 0.522 | 0.577 | 0.428 | 0.603 | 0.714 |
| Weighted-Score | 0.522 | 0.563 | 0.448 | 0.599 | 0.729 |
| **ELERAG (RRF)** | **0.551** | **0.589** | **0.458** | **0.622** | **0.742** |
| RRF+Cross-Encoder | 0.507 | 0.582 | 0.441 | 0.589 | 0.708 |

*Impact of Entity-Aware Ranking on Generation.* The results confirm that the retrieval quality improvement observed in the dense stage propagates effectively to the final RAG output. *ELERAG (RRF)* consistently outperforms all other configurations, achieving the highest scores in *Exact Match* (0.551), *Recall* (0.589), and *MRR* (0.622). This indicates that when the LLM is fed with chunks prioritized by our RRF-Based strategy, it is more likely to identify, use, and cite the correct gold answer.

*Performance Drop of Cross-Encoder.* A critical observation is the underperformance of the *RRF+Cross-Encoder* pipeline, which records the lowest Exact Match (0.507) and MRR scores among all entity-enhanced methods, falling even slightly behind the Baseline. Even with the LLM acting as a final filter, the Cross-Encoder fails to recover the gap. This reinforces the conclusion that for specialized educational content, a heavy semantic re-ranker may introduce noise or drift away from the specific factual precision required, whereas the explicit entity signal used in *ELERAG* ensures a more robust alignment with the query's intent.

### 4.4. Method 1: Evaluation on the Standard Benchmark (SQuAD-it)

Finally, Table 4 reports the retrieval performance on the SQuAD-it dataset. This benchmark serves as a control experiment on a general-domain corpus (Wikipedia).

**Table 4.** Retrieval performance on SQuAD-it (General Domain). Unlike the educational dataset, here the *Cross-Encoder* achieves the best results, highlighting its strength on standard web-like data.

| Pipeline | EM | R@1 | R@3 | R@5 | R@10 | P@1 | MRR |
|---|---|---|---|---|---|---|---|
| Standard RAG (Baseline) | 0.693 | 0.693 | 0.843 | 0.884 | 0.922 | 0.693 | 0.776 |
| Weighted-Score | 0.645 | 0.645 | 0.804 | 0.853 | 0.904 | 0.645 | 0.735 |
| ELERAG (RRF) | 0.672 | 0.672 | 0.829 | 0.875 | 0.922 | 0.672 | 0.760 |
| **Cross-Encoder** | **0.777** | **0.777** | **0.885** | **0.912** | **0.936** | **0.777** | **0.836** |

Analysis of Results.

On the SQuAD-it dataset, the trend observed in the educational benchmark is sharply reversed. The *RRF + Cross-Encoder* pipeline achieves the best performance across all metrics, with an Exact Match of 0.777 and an MRR of 0.836, significantly outperforming both the Baseline and the entity-based methods. Conversely, the *ELERAG (RRF)* method performs slightly worse than the Baseline (0.672 vs 0.693 in EM).

*Evidence of Domain Mismatch.* This divergence in performance between the two datasets provides strong experimental evidence for the *Domain Mismatch* hypothesis. The Cross-Encoder model (`mMiniLM`) was pre-trained on massive general-domain datasets (MS MARCO, Wikipedia), which are linguistically very similar to SQuAD-it. Consequently, it can leverage its internal knowledge to rank these documents effectively. However, as shown in Section 4, this advantage disappears when applied to the specialized, high-ambiguity domain of university lectures, where our proposed *ELERAG* method prevails.

This leads to a crucial conclusion: while Cross-Encoders are the State-of-the-Art solution for general web-search tasks, they are not universally optimal. For specialized educational platforms, where training data for fine-tuning is scarce, our entity-aware hybrid approach provides a more robust and accurate retrieval mechanism without the computational cost of a heavy neural re-ranker.

## 5. Discussion

### 5.1. Interpretation of Results

The experimental results present a clear and divergent pattern across the two benchmarks, which forms the central finding of this study.

First, the quantitative analysis on our specialized educational corpus demonstrates that our proposed *ELERAG* architecture achieved the best overall performance. It consistently outperformed not only the *Weighted-Score* baseline but also the SOTA *Cross-Encoder* model in key retrieval metrics, such as Exact Match and Mean Reciprocal Rank (MRR). Conversely, on the general-purpose *SQuAD-it* benchmark, this trend inverted sharply: the *Cross-Encoder* configuration obtained the highest scores across all metrics, proving its superiority in a standard, general-domain QA setting.

This quantitative divergence provides strong experimental evidence for our *Domain Mismatch* hypothesis. It indicates that while large, pre-trained re-rankers excel on web-style data (like Wikipedia), a domain-adapted hybrid model like *ELERAG* is significantly more effective on specialized, narrative corpora (like university lectures). In these contexts, the explicit signals provided by Entity Linking align better with the data than the generic semantic patterns learned by the Cross-Encoder.

Having established that *ELERAG* is the most effective solution for the specific domain, qualitative analysis helps to explain why the integration of Entity Linking provides this advantage. Inspection of the results confirms that the entity-based signal is crucial in high-ambiguity scenarios. A clear example is the query "Who is Smith?". The *Baseline* system, relying only on semantic similarity, retrieved scattered and less relevant documents based on broad keyword matching. In contrast, the *ELERAG* pipeline, guided by the unambiguous entity ID from Wikidata, successfully promoted the correct, relevant passages to the top ranks. This allowed the LLM to generate a more coherent, factually dense, and precise answer. This qualitative benefit was less pronounced in broad, low-ambiguity queries, where the baseline's semantic search was already sufficient. This suggests that the primary value of our hybrid approach lies in its ability to resolve factual ambiguity, which is a critical weakness of purely semantic systems in technical domains.

Finally, the analysis confirmed the robustness of the full RAG pipeline in handling out-of-domain questions. Thanks to the LLM acting as a semantic filter, irrelevant queries

(e.g., "What is the capital of France?") correctly resulted in a safe fallback response ("No relevant information found"), effectively mitigating model hallucination and reinforcing the system's reliability as a tutoring tool.

### 5.2. Implications and Contributions

The results offer several important implications for research in retrieval-augmented generation and educational AI. First, the consistent improvements obtained through entity-based re-ranking demonstrate that factual grounding can be enhanced without the need for expensive model retraining, solely through structured post-retrieval refinement. Second, the contrast between *ELERAG* and Cross-Encoder performance emphasizes that high-capacity models do not generalize uniformly across all domains: their success strongly depends on the alignment between pre-training data and the target corpus. This finding suggests that lightweight hybrid methods—such as our entity-enriched retrieval—can be a more efficient and accurate alternative to SOTA re-rankers in domain-specific applications, particularly in low-resource or multilingual settings such as Italian educational materials. Finally, the integration of an *Entity Linking* module provides a clear framework for combining symbolic and neural representations, showing how explicit knowledge bases like Wikidata can complement dense embeddings to improve interpretability and precision.

### 5.3. Limitations and Future Work

Despite the promising results, this study has limitations that open clear avenues for future research. A primary limitation is the size and source of our custom benchmark; while effective for this study, its automatic generation via LLM APIs (even with validation) may introduce inherent biases compared to human-curated datasets. Furthermore, our analysis did not include a systematic evaluation of computational latency, although the architectural simplicity of RRF compared to Cross-Encoders suggests a theoretical advantage.

These limitations directly inform future work. The most critical next step, which follows directly from our "Domain Mismatch" finding, is to fine-tune a Cross-Encoder model specifically on our educational corpus. We hypothesize that this domain-adaptation step would allow the Cross-Encoder to learn the specific linguistic patterns of the lectures, potentially surpassing the RRF system and combining the best of both worlds: domain knowledge and SOTA neural architecture.

Further research could also explore adaptive weighting schemes for the RRF, dynamically adjusting the balance between semantic and entity signals based on the query type (e.g., factual vs. conceptual). Finally, expanding the evaluation to multilingual data and incorporating human-in-the-loop assessments would provide a richer understanding of the system's real-world educational value.

**Author Contributions:** Conceptualization, F.G., F.P. and M.M.; methodology, F.G., F.P. and M.M.; software, F.G.; validation, F.G., F.P. and M.M.; data curation, F.G. and F.P.; writing—original draft preparation, F.G.; writing—review and editing, F.G., F.P.. and M.M.; supervision, F.P. and M.M.; funding acquisition, F.P.. All authors have read and agreed to the published version of the manuscript.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** The source code for the system is publicly available on GitHub at https://github.com/Granataaa/educational-rag-el. The custom dataset generated and analyzed during this study is not publicly available because it has been built from proprietary data from a private company.

**Conflicts of Interest:** The authors declare no conflicts of interest.

## Abbreviations

The following abbreviations are used in this manuscript:

| | |
|---|---|
| RAG | Retrieval-Augmented Generation |
| LLM | Large Language Models |
| EL | Entity Linking |
| RRF | Reciprocal Rank Fusion |
| EM | Exact Match |
| MRR | Mean Reciprocal Rank |
| NER | Named Entity Recognition |
| EAI | Educational Artificial Intelligence |
| FAISS | Facebook AI Similarity Search |
| GPT-40 | Generative Pre-trained Transformer 4.0 |

## References

1.  Huang, L.; Yu, W.; Ma, W.; Zhong, W.; et al. A Survey on Hallucination in Large Language Models: Principles, Taxonomy, Challenges, and Open Questions. *ACM Transactions on Information Systems* **2025**, *43*. https://doi.org/10.1145/3703155.
2.  Zhang, Y.; Li, Y.; Cui, L.; Cai, D.; Liu, L.; et al. Siren's Song in the AI Ocean: A Survey on Hallucination in Large Language Models. *Computational Linguistics* **2025**, pp. 1–46. https://doi.org/10.1162/coli.a.16.
3.  Asgari, E.; Montaña-Brown, N.; Dubois, M.; Khalil, S.; et al. A framework to assess clinical safety and hallucination rates of LLMs for medical text summarisation. *npj Digital Medicine* **2025**, *8*. https://doi.org/10.1038/s41746-025-01670-7.
4.  Pal, A.; Umapathi, L.K.; Sankarasubbu, M. Med-HALT: Medical Domain Hallucination Test for Large Language Models. In Proceedings of the Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL). Association for Computational Linguistics, 2023. https://doi.org/10.18653/v1/2023.conll-1.21.
5.  Qian, K.; Liu, S.; Li, T.; Raković, M.; et al. Towards reliable generative AI-driven scaffolding: Reducing hallucinations and enhancing quality in self-regulated learning support. *Computers and Education* **2026**, *240*. https://doi.org/10.1016/j.compedu.2025.105448.
6.  Vrdoljak, J.; Boban, Z.; Vilović, M.; Kumrić, M.; Božić, J. A Review of Large Language Models in Medical Education, Clinical Decision Support, and Healthcare Administration. *Healthcare* **2025**, *13*. https://doi.org/10.3390/healthcare13060603.
7.  Lewis, P.; Perez, E.; Piktus, A.; Petroni, F.; Karpukhin, V.; Goyal, N.; Küttler, H.; Lewis, M.; Yih, W.t.; Rocktäschel, T.; et al. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In Proceedings of the Advances in Neural Information Processing Systems, 2020, Vol. 33, pp. 9459–9474.
8.  Mongiovì, M.; Gangemi, A. GRAAL: Graph-Based Retrieval for Collecting Related Passages across Multiple Documents. *Information* **2024**, *15*, 318. https://doi.org/10.3390/info15060318.
9.  Wikimedia Foundation. Wikidata, 2025. Accessed: September 2025.
10. Shen, W.; Wang, J.; Han, J. Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions. *IEEE Transactions on Knowledge and Data Engineering* **2015**, *27*, 443–460. https://doi.org/10.1109/TKDE.2014.2327028.
11. Möller, C.; Lehmann, J.; Usbeck, R. Survey on English Entity Linking on Wikidata: Datasets and Approaches. *Semantic Web* **2022**, *13*, 925–966. https://doi.org/10.3233/SW-212986.
12. Wu, L.; Petroni, F.; Josifoski, M.; Riedel, S.; Zettlemoyer, L. Scalable Zero-Shot Entity Linking with Dense Entity Retrieval. *arXiv preprint arXiv:1911.03814* **2019**. https://doi.org/10.48550/arXiv.1911.03814.
13. Orlando, R.; Cabot, P.L.H.; Barba, E.; Navigli, R. ReLiK: Retrieve and LinK, Fast and Accurate Entity Linking and Relation Extraction on an Academic Budget. *arXiv preprint arXiv:2408.00103* **2024**. https://doi.org/10.48550/arXiv.2408.00103.
14. Delpeuch, A. OpenTapioca: Lightweight Entity Linking for Wikidata. *arXiv preprint arXiv:1904.09131* **2019**. https://doi.org/10.48550/arXiv.1904.09131.
15. Mageira, K.; Pittou, D.; Papasalouros, A.; Kotis, K.; Zangogianni, P.; Daradoumis, A. Educational AI Chatbots for Content and Language Integrated Learning. *Applied Sciences* **2022**, *12*, 3239. https://doi.org/10.3390/app12073239.

16. Swacha, J.; Gracel, M. Retrieval-Augmented Generation (RAG) Chatbots for Education: A Survey of Applications. *Applied Sciences* **2025**, *15*, 4234. https://doi.org/10.3390/app15084234.

17. Li, Z.; Wang, Z.; Wang, W.; Hung, K.; Xie, H.; Wang, F.L. Retrieval-Augmented Generation for Educational Application: A Systematic Survey. *Computers and Education: Artificial Intelligence* **2025**, p. 100417. https://doi.org/10.1016/j.caeai.2024.100417.

18. Shlyk, D.; Groza, T.; Montanelli, S.; Cavalleri, E.; Mesiti, M. REAL: A Retrieval-Augmented Entity Linking Approach for Biomedical Concept Recognition. In Proceedings of the Proceedings of the 23rd Workshop on Biomedical Natural Language Processing. Association for Computational Linguistics, 2024, pp. 380–389. https://doi.org/10.18653/v1/2024.bionlp-1.34.

19. Johnson, J.; Douze, M.; Jégou, H. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* **2019**, *7*, 535–547. https://doi.org/10.1109/TBDATA.2019.2921572.

20. Wang, L.; Yang, N.; Huang, X.; Yang, L.; Majumder, R.; Wei, F. Multilingual E5 Text Embeddings: A Technical Report. *arXiv preprint arXiv:2402.05672* **2024**. https://doi.org/10.48550/arXiv.2402.05672.

21. OpenAI. GPT-4o API, 2023. Accessed: September 2025.

22. Explosion AI. spaCy: Industrial-Strength Natural Language Processing in Python. https://spacy.io, 2023. Version 3.7.2.

23. Wu, L.; Petroni, F.; Josifoski, M.; Riedel, S.; Zettlemoyer, L. Scalable Zero-shot Entity Linking with Dense Entity Retrieval. In Proceedings of the Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), 2020, pp. 6397–6405. https://doi.org/10.18653/v1/2020.emnlp-main.519.

24. Cormack, G.V.; Clarke, C.L.A. Reciprocal Rank Fusion Outperforms Condorcet and Individual Rank Learning Methods. In Proceedings of the Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval, 2009, pp. 758–759. https://doi.org/10.1145/1571941.1572114.

25. Reimers, N.; Gurevych, I. Cross-Encoder: Sentence Transformers. https://www.sbert.net/examples/applications/cross-encoder/, 2020. Accessed: September 2025.

26. Radford, A.; Kim, J.W.; Xu, T.; Brockman, G.; McLeavey, C.; Sutskever, I. Robust Speech Recognition via Large-Scale Weak Supervision. *arXiv preprint arXiv:2212.04356* **2022**. https://doi.org/10.48550/arXiv.2212.04356.

27. Croce, D.; Zelenanska, A.; Basili, R. Neural Learning for Question Answering in Italian. In Proceedings of the AI*IA 2018 – Advances in Artificial Intelligence. Springer, 2018, pp. 389–402. https://doi.org/10.1007/978-3-030-03840-3_29.