

Relationship Extraction using Retrieval Augmented Generation for biomedical Dataset

Anonymous ACL submission

Abstract

With the increasing number of structured and unstructured data, obtaining reliable information effectively has become crucial. In the biomedical domain, extracting information from the scientific papers is crucial in order to stay up-to-date with accurate information, given the increased pace by which new research studies are published. This work focuses on identifying relationships between entities that are extracted from the abstracts and titles of biomedical research papers. In this work, we developed a Retrieval Augmented Generation (RAG) based system to automatically identify relations between biomedical entities. We evaluate multiple open source Large Language Models (LLMs) and the number of examples (shots) required to improve the LLM's results. We evaluate our methods using precision, recall and F-1 scores and compare our approach to traditional deep learning methods using DeBERTa with a Convolutional Neural Network (CNN). Our results indicate that Qwen models using the RAG approach with 10-shot examples achieved the highest macro F1 score compared to the baseline and other LLMs under the same setting. At 35 shots, Qwen reasoning and Qwen non-reasoning model performed best, exhibiting the fewest hallucinated labels and maintaining high macro F1 scores.

1 Introduction

With the rapid development of technology, the Internet has become part of daily life for everyone. This makes it difficult to stay up to date with useful information all the time. Information Extraction (IE) is an important task in Natural Language Processing. It usually focuses on a specific domain or task, looking for information that is relevant to a user's interests (Hobbs and Riloff, 2010). The need for IE from large-scale scientific resources continues to increase over time (Brokman et al., 2025).

Relationship Extraction (RE) is a subpart of Information Extraction (IE) with many use cases across different domains like automatically extracting relations between entities from biomedical text. RE aims to identify semantic relations between entities in text and build structured knowledge bases (Wang et al., 2022). RE from biomedical corpora is important research to biomedical experts. One common application of biomedical RE is extracting Drug-drug interaction from texts. Drug-drug interaction (DDI) is defined as a change in the effects of one drug by the presence of another drug (Baxter, 2010). It is important to extract knowledge about DDIs from pharmaceutical papers comprehensively to understand the information and have new research around it. It can also be useful for people who need summarization and quick understanding of abstracts when they lack deep domain expertise, helping them save time and avoid being overwhelmed by the technical jargon. Reviewing abstracts from recent and previous biomedical papers is a tedious task, and getting insights from those papers is time-consuming.

Earlier for RE, fine tuned models were fundamentally limited by their pre-trained knowledge scope and often hallucinate, as they are not easily updatable with new biomedical findings. In contrast, large language models (LLMs) such as GPT (Mann et al., 2020) have demonstrated strong performance through prompt engineering. There are a number of prompting techniques such as few-shot, chain-of-thought, and generated knowledge prompting that can enhance a LLM's ability to answer questions. Prompting can serve as a useful tool in a variety of general domains. However, in some cases for specialized domains, such as the biomedical domain, their performance falls short, and finetuning these models is often infeasible due to computational restraints or time. This sets the stage for Retrieval-Augmented

084 Generation (RAG) approaches. These approaches
085 incorporate domain specific knowledge from a
086 vector store into prompts fed into an LLM.
087

088 In this project, we explore efficient approaches
089 for RE using state-of-the-art methods (Lewis et al.,
090 2020). For each relationship extracted from the
091 text between different entities, the system com-
092 putes the embedding similarity between the query
093 entities and those in the training dataset. Then
094 it retrieves the top-*k* most similar candidates us-
095 ing vector search (Olasunkanmi et al., 2025). The
096 results demonstrate that the RAG approach effec-
097 tively mitigates hallucination and incorrect reason-
098 ing in relation extraction tasks (Chen et al., 2024).

099 To take the limitations of LLMs into account,
100 Retrieval-Augmented Generation (RAG) integrates
101 external information, dynamically retrievable re-
102 sources with LLM generation, increasing factual
103 grounding, and enhancing robustness against noisy,
104 vague or counterfactual evidence-key concerns in
105 biomedicine.

106 The diverse applications of Retrieval Augmented
107 Generation (RAG) opens new perspectives, giving
108 dynamic knowledge access and superior robustness
109 against hallucinations and has recent applications
110 in biomedical information extraction challenges
111 like GutBrainIE. Inspired by the potential of RAG
112 for biomedical relationship extraction, the research
113 goals of this work are to address the following
114 questions:

- 115 • How is the performance of Retrieval-
116 Augmented Generation for RE (RAG4RE)
117 pipeline when using different available LLMs
118 compared to traditional neural RE models
119 (such as DeBERTa + CNN) across the
120 Biomedical dataset?
- 121 • How many prompt examples (*shots*) are nec-
122 essary for a given model to achieve optimal
123 results in relationship extraction for biomed-
124 ical dataset? At what point does increasing
125 the number of shots result in saturated perfor-
126 mance for the best performing LLMs?

127 2 Related Work

128 This section reviews prior work on biomedical Re-
129 lation Extraction and Retrieval Augmented Gen-
130 eration systems that our approach builds upon.
131 RAG utilizes both parametric and non-parametric
132 memory to improve the performance of LLMs

133 over NLP tasks. Parametric memory refers to the
134 knowledge stored internally within the model’s
135 learned weights, while non-parametric memory is
136 the information extracted from the external doc-
137 uments(Lewis et al., 2020). In traditional RAG,
138 retrieval is performed directly from the initial user
139 query, without any refinement or classification
140 along the way (Lewis et al., 2020).

141 Dense Passage Retrieval (DPR) (Karpukhin
142 et al., 2020), replaced the traditional sparse tech-
143 nique (e.g.,BM25) (Robertson et al., 2009) with im-
144 proved dual-encoder embeddings for both queries
145 and passages. In open domain Question Answer-
146 ing, DPR achieved a top-20 passage recall of 85.4%
147 on Natural Question(NQ), which is compared to
148 BM25 with 64.3% and 21.2% absolute gain. When
149 it is integrated with RAG pipeline, DPR systems
150 reached EM (Exact Match) scores exceeding 42%
151 and F1 improvements of 9-11% across major QA
152 datasets, measuring dense retrieval as high-quality
153 selection.

154 The queries and external knowledge base are
155 converted into vector space for efficient retrieval.
156 This is known as the embedding step. Then, these
157 vectors get stored into a vector database, allowing
158 for fast similarity search(Lewis et al., 2020) which
159 is known as indexing. The embedding process has
160 evolved significantly from early approaches. While
161 BERT-based embeddings (Devlin et al., 2019) ini-
162 tially dominated, recent work shows that task-
163 specific fine-tuning of embedders substantially im-
164 proves retrieval quality (Ram et al., 2021). The
165 indexing strategy matters as well: dense passage
166 retrieval (DPR) (Chen, 2024) excels at semantic
167 similarity, but could miss lexically important terms,
168 motivating hybrid approaches that combine dense
169 and sparse representations (Glass et al., 2022a). A
170 critical, but often overlooked aspect is the chunking
171 strategy: how documents are segmented can affect
172 both retrieval precision and context preservation,
173 with optimal chunk sizes varying by domain and
174 query type (Li et al., 2025).

175 For relationship extraction, frameworks like
176 RAG4RE (Efeoglu and Paschke, 2024) used
177 prompt and data augmentation by allowing neigh-
178 boring examples and sentences with context sim-
179 ilarity. This improved the model’s F1 score sig-
180 nificantly in comparison to simple queries, the re-
181 port shows up to 94.6% F1 on TACREV, 93% on
182 Re-TACRED, and robust zero-shot performance
183 while using a fine-tuned LLM i.e Mistral-7B. Aug-
184 mented queries have helped retrieval recall improve

by 30%, demonstrating that quality of prompt and prompt expansion helps RAG systems extract information effectively, even from vague or multiple-hop questions. RAG approach with Multi-Source and Retrieval Augmentation(Lee et al., 2024) Hybrid retrieval pipelines, both HYBGRAG and Re2G(Glass et al., 2022b), combine different retrieval approaches-dense(FAISS), sparse (BM25), and sometimes graph-based or multiple query which can be used as source within ensemble combination and reranking modules. Both of these approaches shows how integrating multiple different context sources, allow for critical selection and lead up to better precision i.e. 88% recall i.e. 65.7% and 34% improvements on multi-hop benchmarks in comparison to traditional RAG. These hybrid models allow to retrieve from more diverse context, which is more relevant and have trustworthy evidence.

3 Dataset

The dataset consists of 2,127 train samples with 18 relation labels and 20,227 test samples with 15 relation labels. It focuses on biomedical titles and abstracts related to the gut microbiota and its effects on mental health, following Martinelli et al. (2025).The dataset focuses on biomedical titles and abstracts related to the gut microbiota and its effects on mental health (Martinelli et al., 2025).

The dataset has 15 distinct relations in Test partition: strike, used by, is linked to, located in, change effect, target, interact, influence, impact, part of, change expression, is a, compared to, administered, and change abundance; and Train partition has 3 more distinct relation i.e NONE, affect, produced by; which are not in the Test partition as shown in the Table 1. The diversity and imbalance in frequency among these relations as represented in the Table 1 provide defined types common and rare but biologically significant which cause a challenge for relationship extraction. These relations show the entire range of annotated connections in the dataset used and forms a base for all metric computation in this work.

4 System Architecture

Here in our system, we have incorporated RAG for RE. It’s a four stage pipeline for RE on biomedical dataset which consists of pre-processing and instance construction, dense embedding and indexing, retrieval and few-shot construction, and

Table 1: Ground truth class frequencies for Train and Test dataset.

Label	Train Freq	Test Freq
NONE	1053	0
TARGET	173	1273
CHANGE EXPRESSION	36	494
IMPACT	84	927
INFLUENCE	199	1955
IS LINKED TO	171	1126
LOCATED IN	141	1865
CHANGE ABUNDANCE	55	1568
USED BY	82	225
AFFECT	105	0
PART OF	44	1214
COMPARED TO	5	500
INTERACT	25	2165
IS A	7	4788
PRODUCED BY	6	0
ADMINISTERED	8	974
STRIKE	15	612
CHANGE EFFECT	8	541

LLM-based relation generation as shown in the Figure 1. This work compares relationship extraction on biomedical texts using both transformer-based neural model (DeBERTa+CNN) and different LLMs in a RAG pipeline. The goal is a systematic, head-to-head comparison between a strong transformer-based neural model and a range of modern retrieval-augmented LLM pipelines for relationship extraction. All experiments are conducted on BioASQ Conference and Labs of the Evaluation Forum (CLEF) 2025 GutbrainIE challenge dataset (Martinelli et al., 2025), containing PubMed abstracts on annotated entities. This comparison allows for better understanding of the strengths and limitations to both the methodologies addressing real-world biomedical relationship extraction tasks.

4.1 Pre-processing and instance construction

In this pre-processing pipeline, annotated abstracts and titles are systematically processed to extract entity pairs, generate context windows, and accurately map character offsets necessary for sentence-level cross-sections. For sentence segmentation spaCy is used. The process begins with reading training and development data from JSON files, containing biomedical text samples with annotated head and tail entities and their labels. The workflow not only works for positive relation extraction but also creates representative negative samples to enhance model robustness. Downsampling NONE relations

264	to balance the training set. After working on the	biomedical data being directly compatible, making	311
265	cross-sentences the data.	it possible to accurately compare their meaning	312
266		within most relevant scientific contexts and retrieve	313
267	4.2 Dense Embedding and Indexing	the most relevant sentences every time.	314
268		In this work for every query, we retrieve the most	315
269	4.2.1 Embedding	similar training nodes based on vector similarity.	316
270	All training samples are converted into Document	These nodes are used as few-shot examples for the	317
271	objects with metadata and passed to LlamaIndex’s	generative model which provides the context drawn	318
272	ingestion pipeline. To improve the retrieval by en-	from similar biomedical content as shown in Figure	319
273	hancing the contextual understanding each chunk is	1. Every retrieved node is fed into the prompt by	320
274	converted into a dense vector embedding using the	formatting its passage, entity, labels and relation	321
275	HuggingFaceEmbedding class within LlamaIndex,	information. These components are used into a	322
276	with the all-MiniLM-L6-v2 model. This process	structured prompt which is fed to the LLM.	323
277	encodes the semantic meaning of sentences (which		
278	are cross-section sentences from abstracts and titles),	4.4 LLM-based relation generation	324
279	it also prevents the relationships central to	This work evaluates several open-source large lan-	325
280	biomedical reasoning. The embeddings are created	guage models, to generate the prediction of rela-	326
281	using the train and test dataset, where each item	tions. Each model is fed the augmented prompt	327
282	from dev dataset is a query as shown in Figure 1.	containing k real examples along with the query-	328
283	These embeddings are the foundation for any re-	and is instructed to extract and explain the biomed-	329
284	trieval and inference steps, which make sure the	ical relationship between the given entities. In	330
285	similarities in scientific meaning are being detected	this work, to measure the effectiveness of RAG	331
286	and utilized correctly.	prompting and measure the capabilities of modern	332
287		open-source models in biomedical relationship ex-	333
288	4.2.2 Vector Database with LlamaIndex &	traction, this study evaluates with multiple state-of-	334
289	Qdrant	the-art LLMs. Here, we have used two reasoning	335
290	During index creation all the preprocessed docu-	and two non-reasoning models, each have a dis-	336
291	ments from the train dataset are passed through the	inct combination of reasoning strength, instruction	337
292	embedding model and which are then mapped to	tuning and parameter.	338
293	the vector space. The embeddings are stored and in-	All the models are integrated into the RAG	339
294	indexed in the vector database, in this work we used	pipeline using LlamaIndex and Ollama infrastruc-	340
295	Qdrant is used for embedding model (Ockerman	ture. Every model received the same structured,	341
296	et al., 2025) via QdrantVectorStore in LlamaIndex,	context-enriched prompts and produced predictions	342
297	which is built for the effective semantic search.	in standardized JSON format. By systematically	343
298		using the different number of retrieved examples,	344
299	4.3 Retrieval and few-shot construction	this work measured each model’s accuracy qual-	345
300		ity and robustness against ground-truth biomedical	346
301	4.3.1 Query	relationships from the GutBrainIE dataset.	347
302	When test sample is processed, its query embed-		
303	ding is computed using the same embedding model	5 Models	348
304	as the training data, but it is not stored in the vec-	The Baseline Model (DeBERTa+CNN) focuses on	349
305	tor index. Instead, this query embedding is used	creating context-aware token embeddings and en-	350
306	only at inference time to retrieve the most simi-	hances them with feature extraction using one-layer	351
307	lar training embeddings. LlamaIndex provides the	of convolution.	352
308	abstraction to seamlessly convert this query to an		
309	embedding and leverage Qdrant’s search capabili-	5.1 Baseline Model	353
310	ties.	In this work, the baseline was conducted using	354
		DeBERTa+CNN pipeline which is designed to	355
		extract and classify relationships within biomed-	356
		ical texts. Every document is embedded using	357
		DeBERTa-v3-base transformer, which captures	358

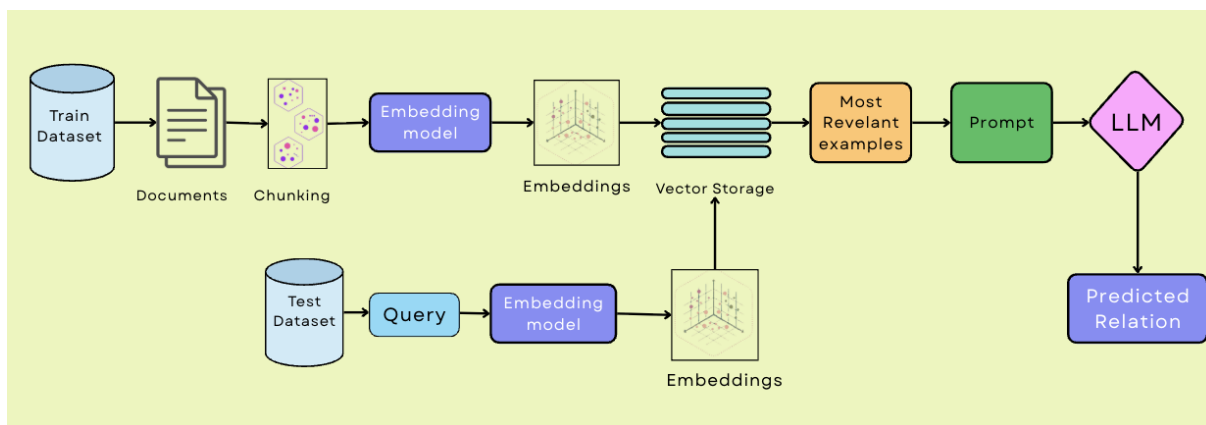


Figure 1: Methodology: RAG for RE on biomedical dataset.

context-dependent meanings of tokens throughout the input text. To further capture semantic information, CNN layer is used directly after embedding. This approach begins with tokenizing input text and feeding it to the DeBERTa-v3-base transformer, which then generates context-aware token embeddings that help the model to have entity recognition capabilities. This transformer-based embeddings are then passed to convolutional layer, consisting of 256 filters and kernel size of three, which allowed for the extraction of local features that are informative for entity interactions. This resulting vector is then passed to a classification layer, giving an 18-dimensional output corresponding to the set of possible relation types defined by the challenge.

5.2 Llama-3.1-8B-Instruct - Non-Reasoning Model

In this work, we used instruction-tuned model from the Llama family is designed for conversational and task-oriented NLP. The purpose of using this model was to allow the analysis being fair how compact models maintain-perform over the relations to larger systems when provided with the real domain-specific, high quality retrieval-augmented prompts. While Llama models (Touvron et al., 2023) are strong generalists, it tends to work slightly behind reasoning-focused architectures in complex biomedical tasks.

5.3 Gemma3:12B - Reasoning Model

This open-source LLM produced by google (Team et al., 2024), was used in this work to add and broaden the comparison scope and assess performance across diverse reasoning-intensive and biomedical scenarios. Gemma delivers consistent results, especially strong in educational reasoning

and standard question-answering tasks.

5.4 Qwen3:14B - Non-Reasoning Model

Qwen is a comprehensive language model series that encompasses distinct models with varying parameter counts (Bai et al., 2023). This version of Qwen model does not explicitly include the reasoning of the prediction it made.

5.5 Qwen3:14B - Reasoning Model

This version of Qwen model have complex reasoning and advanced comprehension, which makes it particularly effective for interpreting biomedical relationships. The reasoning model responded with the relation and a reasoning of the relation provided.

6 Evaluation metrics

To have a fair analysis of relationship extraction performance on biomedical dataset, this work measures Precision and Recall to measure macro and micro F1 score. As in our work there is class imbalance we are measuring both macro F1 score. This capability not only strengthens model interpretability but enhances sensitivity toward minority relations while maintaining overall accuracy.

7 Results

In this work, the results indicates that RAG for RE using LLM have performed better DeBERTa + CNN baseline. Here, we measure the performance of both approaches for relationship extraction on the biomedical dataset for a fair comparison. In Table ??, the macro-averaged metrics is used for evaluation as this method shows how well a model

handles minority classes, which is particularly important in biomedical relationship extraction, when there is class imbalance, which often carry important biomedical meaning.

7.1 Comparison between Baseline and other LLMs

Here, in this work DeBERTa + CNN measured Macro F1 score of 0.4404 and Micro F1 of 0.5625, giving a foundation for RE to measure against RAG-based LLM models (Table ??). Since there are few different LLMs, we compared evaluation score across four of them and the baseline model. Empirically, two of the models i.e Qwen3:14B (Reasoning Model) and Qwen3:14B (Non Reasoning Model), had higher evaluation scores with 10 relevant examples (shots) compared to baseline and other LLMs with RAG. The best RAG systems are Qwen3:14B (non-reasoning) with macro F1 0.766 and micro F1 0.753, and Qwen3:14B (reasoning) with macro F1 0.758 and micro F1 0.742, both substantially outperforming the baseline. Using $n = 10$ shots, Llama-3.1-8B-Instruct reaches macro F1 0.211 and micro F1 0.370, while Gemma3:12B reaches macro F1 0.3179 and micro F1 0.658, which remain below the Qwen models but illustrate the gains from RAG over the purely supervised baseline in some metrics.

Llama-3.1-8B-Instruct (Non-Reasoning Model): This model shows the weakest performance with macro f1 of 0.211 and micro F1 of 0.370; which is even less than the baseline performance, macro F1 of 0.4404, micro F1 of 0.5625; showing challenges for the biomedical dataset. As shown in Table 3, this model tend to over-predict some relations with relatively low ground-truth counts, such as *strike* (612 instances vs. 2171 predictions by Llama-3.1-8B-Instruct), and under-predict others, such as *compared to* (500 instances vs. only 64 predictions by Llama-3.1-8B-Instruct). At the same time, prediction behavior is more stable for frequent relations like *is linked to*, *is a*, and *located in*, whose predicted counts stay closer to their ground-truth frequencies across models.

Gemma3:12B (Reasoning Model): This model underperformed the baseline approach than Qwen3 reasoning and non-reasoning models, with macro F1 of 0.317 and micro F1 of 0.658, showing limitation in reasoning or comprehending the biomedical

dataset for this task, and several classes where it predicts far fewer relations than ground truth (e.g. "change abundance" and "strike") as it can be seen in Table 3.

Qwen3:14B (Non Reasoning Model): For this model, explicitly reasoning is disabled, this version performed comparably well; it achieved macro F1 of 0.763 and micro F1 of 0.754. These results shows that even without explicit reasoning outputs, the predicted labels are closer to ground truth(GT) in several categories and performed better than the other models. For example "is a" (3648 vs. GT 4788) and "influence" (3124 vs. GT 1955) as it can be seen in Table 3.

Qwen3:14B (Reasoning Model): This model has consistently outperformed others across evaluation metrics, with macro F1 of 0.758 among all the other models and highest micro F1 of 0.742 as it can be seen in Table 2. The prediction across different classes throughout different number of shots shown in Table 3 have significantly improved the F1 score per class as well, as seen in Table ?? . The alignment between the macro and micro results shows that, this model has strong generalization across both frequent relations and non-frequent ones, likely due to improved context understanding supported by retrieval based reasoning. For relations such as "located in" (2302 vs. GT 1865) and "is a" (3862 vs. GT 4788), the model predicts close to ground truth, reflecting improved context comprehension and less over-prediction compared to non-reasoning. The Qwen3 reasoning model's architecture benefits from retrieval augmentation as the number of relevant retrieved example increases it performs better, but it gets saturated after $n = 60$ as it is shown in Figure 2.

7.2 Qwen Reasoning vs Non-Reasoning behavior

The observed differences between reasoning and non-reasoning variants, between Qwen-Reasoning and Qwen-Non Reasoning models is very minimal. Due to the class imbalance in the dataset, we use macro F1 to give weight to minority classes, which give weight to minority of classes as shown in Table-3, class labels and their count of ground truth.

Here we have conducted ablations across both the Qwen3:14B reasoning and non-reasoning variants in the RAG pipeline over 0–60 shots. Both exhibit dynamic class prediction behavior, with predicted label frequencies evolving. At $n = 10$, both

Table 2: Comparison of baseline and LLM models on Macro and Micro metrics.

Model	Macro P	Macro R	Macro F1	Micro P	Micro R	Micro F1
Baseline	0.5181	0.4330	0.4404	0.6585	0.4909	0.5625
Llama-3.1-8B-Instruct	0.290	0.232	0.211	0.370	0.370	0.370
Gemma3:12B	0.354	0.323	0.3179	0.658	0.658	0.658
Qwen3:14B non-reasoning	0.579	0.554	0.766	0.753	0.753	0.753
Qwen3:14B reasoning	0.599	0.581	0.758	0.742	0.742	0.742

Table 3: Per-class predicted frequencies for all models compared to ground truth.

Class Label	Ground Truth	Gemma3:12B	Llama-3.1-8B-Instruct	Qwen3:14B Non-Reasoning	Qwen3:14B Reasoning
is linked to	1126	3348	5754	2568	2816
compared to	500	479	64	626	670
located in	1865	1930	1845	2207	2302
change abundance	1568	281	303	746	628
is a	4788	3690	1229	3648	3682
target	1273	958	1386	1091	1052
strike	612	2171	37	356	365
interact	2165	1521	255	1123	1267
impact	927	1469	678	1608	1052
influence	1955	2771	4313	3124	3005
administered	974	1352	1342	777	766
change expression	494	334	414	391	376
used by	225	240	227	117	157
part of	1214	1169	1295	1145	1164
change effect	541	425	748	585	568

Table 4: Macro F1 and hallucinated labels: Qwen3:14B Reasoning vs Non-Reasoning.

#Shots	Reasoning Macro F1	#Halluc.	Non-Reasoning Macro F1	#Halluc.
0	0.084	24	0.106	39
5	0.607	5	0.607	9
10	0.758	4	0.767	5
15	0.816	1	0.799	2
20	0.849	0	0.852	1
25	0.870	0	0.878	1
30	0.876	1	0.885	1
35	0.885	1	0.892	0
40	0.887	4	0.896	4
45	0.893	2	0.905	3
50	0.893	5	0.899	2
55	0.891	12	0.903	10
60	0.891	17	0.902	16

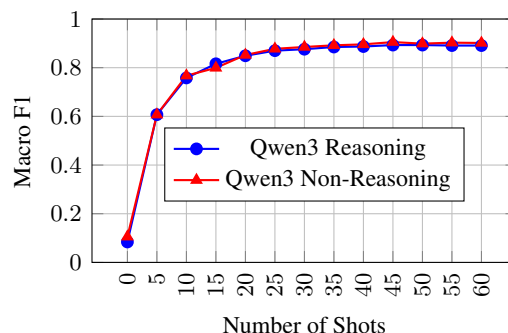


Figure 2: Macro F1 for Qwen3:14B reasoning and non-reasoning across shot counts.

surpass the baseline by ~ 0.3 Macro F1, steadily reaching to saturation beyond $n = 60$ (Figure 2).

At low shots (0 to 10), predictions align closely with ground truth but introduce hallucinations (24 vs 39 in zero-shot). Qwen reasoning variants achieve zero hallucinations at 20 and 25 shots (Table 4), yet hallucinated labels got predicted again beyond 30. At 35 shots under Qwen non reasoning there was zero hallucinated label predicted

At high shots (35 to 60), expanded context drives predictions beyond annotated relations by both the Qwen reasoning and non-reasoning model, capturing subtle biomedical interactions absent from the test set. For example, at $n = 40$, both predict low-frequency spurious classes ("administered to" (2), "is caused by" (1), "is" (1), "located_in" (2)). These

minimally impact Macro F1 due to their rarity (1 to 4 instances) against 20,227 samples, allowing F1 gains despite increased hallucinations (Figures 3, 2).

Non-reasoning slightly leads post-30 shots while matching hallucination patterns, suggesting few-shot RAG regularizes both similarly.

8 Conclusion

This work suggests that retrieval-augmented LLMs can go beyond a strong DeBERTa+CNN baseline for biomedical relationship extraction, especially on rare but important relations, when they are supported with good domain examples. The Qwen3:14B models—most notably the reasoning

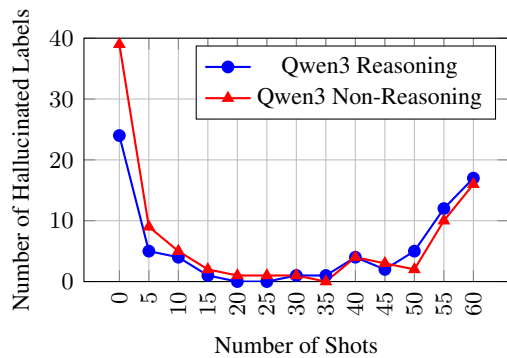


Figure 3: Number of hallucinated labels for Qwen3:14B reasoning and non-reasoning across shot counts.

variant—struck the best balance between macro and micro F1, while also offering clear rationales that make their decisions easier to check and trust. Experiments over different numbers of retrieved shots showed quick improvements up to about 10 examples and a plateau around 60, with a small trade-off: higher performance came alongside a few extra, plausible but unannotated relations.

9 Limitations

There are several limitations that motivate future directions. First, the evaluation was conducted on a dataset where every test instance contained at least one annotated relation, so the models were never required to explicitly decide that “no valid relation” exists between entity pairs. This setting does not reflect many real-world biomedical scenarios, where unrelated entities frequently co-occur and models must reliably abstain from predicting a relation. They can hallucinate relation labels, especially under distribution shift or with poorly chosen shots. Generated outputs should be treated as research signals, not as verified biomedical facts.

Finally, the current work focused on relation extraction and did not extend the approach to upstream tasks such as Named Entity Recognition or to richer retrieval paradigms like GraphRAG-based, entity-centric document retrieval. Without these extensions, the pipeline does not yet capture graph-structured biomedical context, which may be necessary to fully support complex, multi-hop relations in noisy real-world datasets.

Data documentation. We provide a data sheet in the supplementary material that documents the motivation, composition, collection and preprocessing, intended uses and known limitations of the biomedical relation extraction dataset used in this work.

Model documentation. We also provide model cards for our RAG4RE pipeline and the DeBERTa+CNN baseline in the supplementary material, describing their configurations, intended use, training data, evaluation, and limitations.

A Data Sheet for Biomedical RE Dataset

A.1 Motivation

This dataset is used to study biomedical relation extraction on the gut–brain axis, focusing on relations between biomedical entities in PubMed abstracts and titles.

A.2 Composition

The dataset consists of annotated abstracts and titles with head and tail entities and relation labels (including a NONE label). We downsample NONE relations to reduce class imbalance and improve training stability.

A.3 Collection and Preprocessing

We use the training and development splits provided by the GutBrainIE challenge organizers. Text is segmented into sentences with spaCy, entities and relations are mapped to sentence-level contexts, and cross-sentence are constructed when head and tail entities span multiple sentences.

A.4 Intended Uses and Risks

This dataset is intended for research on biomedical relation extraction and RAG systems. Potential risks include label noise from distant annotations and domain bias toward gut–brain literature.

A.5 Distribution

Researchers can obtain the original data from the GutBrainIE challenge or associated CLEF resources. Our processed splits and scripts can be released under the same conditions as original dataset.

B Model Card for RAG for RE Pipeline

B.1 Model Details

In our work, main system is Retrieval-Augmented Generation (RAG) pipeline that uses all-MiniLM-L6-v2 embeddings with LlamaIndex, a Qdrant vector database, and Qwen-based large language models as generators. We also train DeBERTa-based encoder with a CNN classifier as a supervised baseline.

638	B.2 Intended Use and Risk		
639	The models are intended for research on biomedical		
640	relation extraction and evaluation of RAG architec-		
641	tures. They are not intended for any high-stakes		
642	medical decision-making.		
643	B.3 Training Data and Setup		
644	The baseline model is trained on the annotated train-		
645	ing split with downsampled NONE relations, using		
646	standard supervised learning. The RAG pipeline		
647	indexes training corpus in Qdrant via LlamaIndex		
648	and retrieves top- k similar examples for each query.		
649	We study different few-shot configurations (e.g., 0,		
650	5, 10, 25 shots) in prompts for Qwen models.		
651	B.4 Evaluation		
652	We evaluate using macro F1, precision, and recall		
653	on the held-out development set. Qwen models		
654	with 10-shot RAG prompting achieve the highest		
655	macro F1 compared to the DeBERTa+CNN base-		
656	line under the same setting; at 25 shots, the non-		
657	reasoning Qwen model performs best with fewer		
658	hallucinated labels.		
659	B.5 Limitations and Risks		
660	The models are trained on a specific biomedical		
661	subdomain and may not generalize to other areas.		
662	References		
663	Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang,		
664	Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei		
665	Huang, and 1 others. 2023. Qwen technical report.		
666	arXiv preprint arXiv:2309.16609 .		
667	Karen Baxter. 2010. Stockley’s Drug interactions		
668	pocket companion 2010 . Pharmaceutical Press 2010.		
669	Aviv Brokman, Xuguang Ai, Yuhang Jiang, Shashank		
670	Gupta, and Ramakanth Kavuluru. 2025. A bench-		
671	mark for end-to-end zero-shot biomedical relation		
672	extraction with llms: Experiments with openai mod-		
673	els . Preprint , arXiv:2504.04083.		
674	Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun.		
675	2024. Benchmarking large language models in		
676	retrieval-augmented generation . In Proceedings of		
677	the AAI Conference on Artificial Intelligence , vol-		
678	ume 38 , pages 17754–17762.		
679	Zisong Chen. 2024. The question-and-answer system		
680	based on dpr system and llava .		
681	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and		
682	Kristina Toutanova. 2019. Bert: Pre-training of deep		
683	bidirectional transformers for language understand-		
684	ing . Preprint , arXiv:1810.04805.		
	Sefika Efeoglu and Adrian Paschke. 2024. Retrieval-	685	
	augmented generation-based relation extraction .	686	
	arXiv preprint arXiv:2404.13397 .	687	
	Michael Glass, Gaetano Rossiello, Md Faisal Mah-	688	
	bub Chowdhury, Ankita Naik, Pengshan Cai, and	689	
	Alfio Gliozzo. 2022a. Re2G: Retrieve, rerank,	690	
	generate . In Proceedings of the 2022 Conference	691	
	of the North American Chapter of the Association	692	
	for Computational Linguistics: Human Language	693	
	Technologies , pages 2701–2715, Seattle, United	694	
	States. Association for Computational Linguistics.	695	
	Michael Glass, Gaetano Rossiello, Md Faisal Mah-	696	
	bub Chowdhury, Ankita Rajaram Naik, Pengshan Cai,	697	
	and Alfio Gliozzo. 2022b. Re2g: Retrieve, rerank,	698	
	generate . arXiv preprint arXiv:2207.06300 .	699	
	Jerry R Hobbs and Ellen Riloff. 2010. Information ex-	700	
	traction . Handbook of natural language processing ,	701	
	15:16.	702	
	Vladimir Karpukhin, Barlas Oguz, Sewon Min,	703	
	Patrick SH Lewis, Ledell Wu, Sergey Edunov, Danqi	704	
	Chen, and Wen-tau Yih. 2020. Dense passage re-	705	
	trieval for open-domain question answering . In	706	
	EMNLP (1) , pages 6769–6781.	707	
	Meng-Chieh Lee, Qi Zhu, Costas Mavromatis, Zhen	708	
	Han, Soji Adeshina, Vassilis N Ioannidis, Huzefa	709	
	Rangwala, and Christos Faloutsos. 2024. Hyb-	710	
	rag: Hybrid retrieval-augmented generation on tex-	711	
	tual and relational knowledge bases . arXiv preprint	712	
	arXiv:2412.16311 .	713	
	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	714	
	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	715	
	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	716	
	täschel, and 1 others. 2020. Retrieval-augmented	717	
	generation for knowledge-intensive nlp tasks . Advances	718	
	in neural information processing systems , 33:9459–	719	
	9474.	720	
	Siran Li, Linus Stenzel, Carsten Eickhoff, and Seyed Ali	721	
	Bahrainian. 2025. Enhancing retrieval-augmented	722	
	generation: a study of best practices . arXiv preprint	723	
	arXiv:2501.07391 .	724	
	Ben Mann, Nick Ryder, Melanie Subbiah, J Kaplan,	725	
	P Dhariwal, A Neelakantan, P Shyam, G Sastry,	726	
	A Askell, S Agarwal, and 1 others. 2020. Lang-	727	
	uage models are few-shot learners . arXiv preprint	728	
	arXiv:2005.14165 , 1(3):3.	729	
	Marco Martinelli, Gianmaria Silvello, Vanessa Bon-	730	
	ato, Giorgio Maria Di Nunzio, Nicola Ferro, Or-	731	
	nella Irrera, Stefano Marchesin, Laura Menotti, Fed-	732	
	ERICA Vezzani, and 1 others. 2025. Overview of gut-	733	
	brainie@ clef 2025: gut-brain interplay information	734	
	extraction . In CLEF .	735	
	Seth Ockerman, Amal Gueroudji, Song Young Oh,	736	
	Robert Underwood, Nicholas Chia, Kyle Chard,	737	
	Robert Ross, and Shivaram Venkataraman. 2025. Ex-	738	
	ploring distributed vector databases performance on	739	
	hpc platforms: A study with qdrant . arXiv preprint	740	
	arXiv:2509.12384 .	741	

742 Olawumi Olasunkanmi, Mathew Saturdays, Hong Yi,
743 Chris Bizon, Harlin Lee, and Stanley Ahalt. 2025.
744 Relate: Relation extraction in biomedical abstracts
745 with llms and ontology constraints. [arXiv preprint](#)
746 [arXiv:2509.19057](#).

747 Ori Ram, Gal Shachaf, Omer Levy, Jonathan Be-
748 rant, and Amir Globerson. 2021. Learning to re-
749 trieve passages without supervision. [arXiv preprint](#)
750 [arXiv:2112.07708](#).

751 Stephen Robertson, Hugo Zaragoza, and 1 others. 2009.
752 The probabilistic relevance framework: Bm25 and
753 beyond. [Foundations and Trends® in Information](#)
754 [Retrieval](#), 3(4):333–389.

755 Gemma Team, Thomas Mesnard, Cassidy Hardin,
756 Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,
757 Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale,
758 Juliette Love, and 1 others. 2024. Gemma: Open
759 models based on gemini research and technology.
760 [arXiv preprint arXiv:2403.08295](#).

761 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier
762 Martinet, Marie-Anne Lachaux, Timothée Lacroix,
763 Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal
764 Azhar, and 1 others. 2023. Llama: Open and effi-
765 cient foundation language models. [arXiv preprint](#)
766 [arXiv:2302.13971](#).

767 Hailin Wang, Ke Qin, Rufai Yusuf Zakari, Guoming
768 Lu, and Jin Yin. 2022. Deep neural network-based
769 relation extraction: an overview. [Neural Computing](#)
770 [and Applications](#), 34(6):4781–4801.