

UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA



DEPARTMENT OF INFORMATION ENGINEERING

MASTER THESIS IN COMPUTER ENGINEERING

# Knowledge Graph Fact Verification using Retrieval-Augmented Generation

MASTER CANDIDATE

**Farzad Shami**

Student ID 2090160

SUPERVISOR

**Prof. Gianmaria Silvello**

University of Padova

Co-SUPERVISOR

**Prof. Stefano Marchesin**

University of Padova

ACADEMIC YEAR  
2024/2025

DATE: NTH MARCH 2025



*To all those who have believed in me  
and encouraged me to pursue my passions.*



## **Abstract**

Ensuring the truthfulness of knowledge graphs is critical as they serve as foundational tools in powering many AI systems, search engines, and decision-support systems. This thesis proposes a novel approach using retrieval-augmented generation (RAG) to verify facts in knowledge graphs. The system combines large language models, information retrieval techniques, and a multi-stage verification process to evaluate the veracity of knowledge graph facts. Key components include generating queries from knowledge graph triples, integrating web search for external evidence retrieval, employing embedding-based document chunking and retrieval, and utilizing an ensemble of language models for fact verification. The system aggregates outputs from multiple models using adaptive dispute resolution techniques and majority voting. Comprehensive experiments evaluate various text chunking techniques, embedding models, and document selection procedures. Results demonstrate both the effectiveness of the proposed approach in verifying knowledge graph facts and areas for further optimization. This work advances the evolving fields of knowledge base curation and automated fact checking.



# Contents

<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xiii</b>
<b>List of Algorithms</b>	<b>xvii</b>
<b>List of Code Snippets</b>	<b>xvii</b>
<b>List of Acronyms</b>	<b>xix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Background and Motivation . . . . .	1
1.2 Problem Statement . . . . .	2
1.3 Proposed Approach . . . . .	2
1.4 Contributions . . . . .	3
1.5 Thesis Structure . . . . .	4
1.6 Significance and Potential Applications . . . . .	4
<b>2 Related Works</b>	<b>7</b>
2.1 Entailment Verification and Language Models . . . . .	7
2.2 Claim Verification in the Age of Large Language Models . . . . .	9
2.3 Retrieval-Augmented Fact Verification by Synthesizing Contrastive Arguments . . . . .	11
2.4 RAGAR: RAG-Augmented Reasoning for Political Fact-Checking using Multimodal LLMs . . . . .	12
<b>3 Pipeline</b>	<b>15</b>
3.1 Knowledge Graph Dataset . . . . .	17
3.1.1 Definition and Purpose . . . . .	17

## CONTENTS

3.1.2	Structure and Components . . . . .	17
3.1.3	Role in the Overall Pipeline . . . . .	18
3.2	Query Generation and Processing . . . . .	18
3.2.1	Human-Understandable Text Generation . . . . .	18
3.2.2	Question Formulation Techniques . . . . .	19
3.2.3	Cross-Encoder for Query Relevance Scoring . . . . .	20
3.2.4	Relevance Threshold and Sorting . . . . .	21
3.3	Information Retrieval Mechanisms . . . . .	22
3.3.1	Google Search Integration . . . . .	22
3.3.2	Process and Extract Links . . . . .	23
3.3.3	Data Pool Creation . . . . .	25
3.3.4	Data filtering . . . . .	26
3.4	Embedding and Retrieval Tasks . . . . .	26
3.4.1	Embedding Techniques for Smaller Chunks . . . . .	27
3.4.2	Similarity Cutoff Strategy . . . . .	28
3.5	LLMs . . . . .	29
3.5.1	Integration of Multiple LLMs . . . . .	29
3.5.2	Roles of LLMs in the Pipeline . . . . .	30
3.6	Voting System and Conflict Resolution . . . . .	31
3.6.1	Majority Voting System . . . . .	31
3.6.2	Conflict Resolution Strategies . . . . .	32
3.7	Pipeline Flow and Decision Points . . . . .	33
3.8	Performance Report . . . . .	36
3.9	Ethical Considerations and Limitations . . . . .	37
3.9.1	Ethical Considerations . . . . .	37
3.9.2	Limitations . . . . .	38
<b>4</b>	<b>Empirical Evaluation</b> . . . . .	<b>39</b>
4.1	Dataset Analysis . . . . .	39
4.1.1	FactBench Dataset . . . . .	39
4.1.2	YAGO Dataset . . . . .	40
4.1.3	DBpedia Dataset . . . . .	41
4.2	Candidate Models . . . . .	43
4.2.1	Gemma2 . . . . .	43
4.2.2	Qwen2.5 . . . . .	43
4.2.3	Llama3.1 . . . . .	44

4.2.4	Mistral . . . . .	44
4.3	Experimental Setup . . . . .	45
4.3.1	Datasets . . . . .	45
4.3.2	Performance Metrics and Evaluation . . . . .	45
4.3.3	System Configurations . . . . .	47
4.4	Comparative Analysis . . . . .	49
4.4.1	Discussion of Results . . . . .	49
4.4.2	Key Findings . . . . .	50
4.4.3	Error Analysis . . . . .	50
<b>5</b>	<b>Ablation Study</b>	<b>51</b>
5.1	Evaluation Methodology . . . . .	52
5.1.1	Iterative Optimization Process . . . . .	52
5.1.2	Evaluation Metrics . . . . .	52
5.1.3	Significance of the Methodology . . . . .	53
5.2	Document Selection . . . . .	53
5.2.1	Unsupervised Methods . . . . .	54
5.2.2	Supervised Methods . . . . .	55
5.2.3	Evaluation with Large Language Models . . . . .	57
5.3	Embedding Models . . . . .	58
5.3.1	Gte-large-en-v1.5 . . . . .	59
5.3.2	Jina-embeddings-v3 . . . . .	60
5.3.3	Stella_en_1.5B_v5 . . . . .	61
5.3.4	Multilingual-e5-large-instruct . . . . .	61
5.3.5	bge-small-en-v1.5 . . . . .	62
5.3.6	Comparative Analysis . . . . .	63
5.4	Chunking Strategies . . . . .	66
5.4.1	Parsing Documents into Text Chunks . . . . .	66
5.4.2	Smaller Child Chunks Referring to Bigger Parent Chunks (Small2Big) . . . . .	67
5.4.3	Sentence Window Retrieval . . . . .	68
5.4.4	Evaluation . . . . .	69
5.5	Similarity Cut-off . . . . .	70
5.6	Top K . . . . .	71
5.7	Evaluation . . . . .	71
5.8	Failure Analysis . . . . .	72

## CONTENTS

<b>6 Conclusions and Future Works</b>	<b>77</b>
<b>Appendices</b>	<b>79</b>
<b>A Prompt Templates</b>	<b>79</b>
A.1 Human-understandable text generation Prompt . . . . .	79
A.2 Question Generation Prompt . . . . .	80
A.3 RAG Prompt . . . . .	81
<b>B Chunking Strategies</b>	<b>83</b>
<b>References</b>	<b>87</b>
<b>Acknowledgments</b>	<b>93</b>

# List of Figures

2.1	<b>Distinctions between human and LLM Inferences.</b> The entailment prediction performance of humans and LLMs are depicted by a 5-star rating scale [30]. . . . .	8
2.2	Comparison of claim verification systems between NLP-based (traditional) and LLM-based for claim veracity. [7]. . . . .	10
2.3	The proposed RAFTS [40], which performs few-shot fact verification by incorporating informative in-context demonstrations and contrastive arguments with nuanced information derived from the retrieved documents . . . . .	12
2.4	An overview of the fact-checking pipeline [16] contrasting the baseline Sub-Question Generation approach from the RAGAR: Chain of RAG and RAGAR: Tree of RAG approach followed by a veracity explanation generated by a Veracity Prediction module. .	13
3.1	Cross-Encoder Architecture . . . . .	20
3.2	Knowledge Distillation Process . . . . .	21
3.3	Fetching the results from Google Search . . . . .	23
3.4	Extracted Content from the Crawled URLs using newspaper4k .	25
3.5	Node sentence window replacement [20] . . . . .	27
5.1	Document Retrieval Confusion Matrix . . . . .	57
5.2	Positive Labels . . . . .	73
5.3	Negative Labels . . . . .	73
5.4	All Labels . . . . .	73
5.5	Distribution of Fully Incorrect Predictions 4/4 . . . . .	75
5.6	Distribution of Partially Incorrect Predictions 3/4 . . . . .	75



# List of Tables

3.1 Examples of Generation of Human-Understandable Text . . . . .	19
3.2 Question Generation and Scoring Procedure . . . . .	21
3.3 Comprehensive list of decision points in the pipeline flow. . . . .	35
36table.caption.34	
3.5 Performance of our Information Retrieval Mechanisms. . . . .	37
3.6 Limitions of the pipeline. . . . .	38
4.1 Summary of Datasets Used in Empirical Evaluation . . . . .	42
4.2 System Configurations for Empirical Evaluation . . . . .	48
4.3 Empirical Evaluation Results of the Proposed System and Candidate Models over the Datasets . . . . .	49
5.1 Performance of Pre-trained Cross-Encoders . . . . .	56
5.2 Evaluation Results for Different Methods through the Pipeline (just with the Gemma2 model) . . . . .	58
5.3 Comparison of Embedding Models . . . . .	64
5.4 Evaluation Results for Different Embedding Models through the Pipeline (just with the Gemma2 model) . . . . .	65
5.5 Evaluation Results for Different Chunking Strategies through the Pipeline (just with the Gemma2 model) . . . . .	69
5.6 Evaluation Results for Different Similarity Cut through the Pipeline (just with the Gemma2 model) . . . . .	71
5.7 Evaluation Results for Different Top_k through the Pipeline (just with the Gemma2 model) . . . . .	71
5.8 Evaluation Results for Whole Pipeline by Category . . . . .	72
5.9 Evaluation Results for Whole Pipeline . . . . .	72
5.10 Examples of Failure Cases and Error Analysis . . . . .	74

## LIST OF TABLES

6.1	Table example . . . . .	77
B.1	Text Splitter - Chuck Size 512 . . . . .	84
B.2	Small to Big (base retriever chunk size set to 1024) . . . . .	84
B.3	Sliding Window - Window Size 3 . . . . .	85

# List of Algorithms

1	Resolve Ties in Majority Voting System . . . . .	34
2	Stability Across Queries . . . . .	46
3	Similarity Cutoff Postprocessor . . . . .	70



# List of Code Snippets

3.1	Crawling the Extracted URLs . . . . .	23
5.1	Small2Big Chunking Method . . . . .	67



# List of Acronyms

**CSV** Comma Separated Values

**NLI** Natural Language Inference

**NLP** Natural Language Processing

**LLMs** Large Language Models

**LLM** Large Language Model

**QA** Question Answering

**RAG** Retrieval-Augmented Generation

**ML** Machine Learning

**IR** Information Retrieval

**HTML** HyperText Markup Language

**RoPE** rotary position embeddings

**MTEB** Massive Text Embedding Benchmark

**LoRA** Low-Rank Adaptation

**MRL** Matryoshka Representation Learning

**RLHF** Reinforcement Learning with Human Feedback

**SFT** Supervised Fine-Tuning

**GQA** Grouped-Query Attention

**SWA** Sliding Window Attention



# 1

## Introduction

With the rise of big data and AI, it's crucial to have the ability to verify facts and analyze information. Knowledge graphs, which represent knowledge through entities and their relationships, have emerged as an effective tool for arranging and querying massive amounts of structured data. Maintaining the accuracy and dependability of knowledge graphs is a significant challenge. This thesis presents a novel technique for checking facts in knowledge graphs using retrieval-augmented generation, which combines the benefits of Information Retrieval (IR), Natural Language Processing (NLP), and Machine Learning (ML).

### 1.1 BACKGROUND AND MOTIVATION

A lot of different kinds of apps use knowledge graphs now, from search engines and recommendation systems to question-answering sites and virtual helpers. They store information in a structured way that makes it easy to question and draw conclusions. But current knowledge graphs are so big and complicated that it's hard to check every fact they contain by hand. Because of this, automated fact-checking systems are needed to keep these knowledge sources legitimate and reliable.

Rule-based systems or simple statistical methods are often used in traditional ways to check facts. These methods can work for some types of facts, but they don't work well for more complicated or subtle data. Recent progress in ML and NLP has made it possible for fact checking systems to become smarter.

## 1.2. PROBLEM STATEMENT

Large Language Models (LLMs) have shown amazing skills in understanding and producing text that sounds like it was written by a person. This makes them very likely to be successful in tasks that need to verify facts. On the other hand, LLMs have some problems. They can sometimes make up information that sounds reasonable but isn't true, This is called "hallucination." Also, the data they were taught on is all they know, and that data may become out-of-date over time. To get around these problems, this research has come up with retrieval-augmented generation (RAG) methods that mix the best parts of LLMs with information from outside sources.

## 1.2 PROBLEM STATEMENT

This thesis tackles the issue of automated fact verification in knowledge graphs with a RAG-based methodology. Our objective is to create a system capable of:

- Retrieve relevant information from external sources to support or refute claims in a knowledge graph.
- Utilize LLMs to reason about the retrieved information and generate accurate assessments of fact truthfulness.
- Handle a wide range of fact types and domains, from simple statements to more complex relational facts.
- Provide multiple responses for its verification decisions, enhancing transparency and trust in the system.

## 1.3 PROPOSED APPROACH

Our proposed approach combines several key components to create a robust fact verification system:

- **Knowledge Graph Representation:** We start by representing facts from the knowledge graph in a format suitable for processing by language models and IR systems. This involves converting the subject-predicate-object triples of the knowledge graph into natural language statements.
- **Query Generation:** For each fact to be verified, we generate multiple queries designed to retrieve relevant information from external sources. These queries are formulated to capture different aspects of the fact and potential supporting or contradicting evidence.

- **IR:** We use advanced IR techniques to search for relevant documents or passages from a large corpus of trusted sources. This step leverages both traditional search algorithms and dense retrieval methods based on neural networks.
- **Context Processing:** The retrieved information is processed and combined to create a comprehensive context for each fact. This may involve techniques such as text summarization, entity linking, and coreference resolution to create a coherent representation of the relevant information.
- **Large Language Model (LLM) Integration:** We use several LLMs at the same time to look at the context that was retrieved and decide if the original fact is true. By putting together the results of several models, we hope to reduce the flaws in each one and make the whole thing more accurate.
- **Fact Verification Decision:** The system makes a final decision on the truthfulness of the fact based on the consensus of the language models and the strength of the supporting or contradicting evidence. This decision is accompanied by the reasoning process and relevant evidence.

## 1.4 CONTRIBUTIONS

This thesis makes several key contributions to the field of knowledge graph fact verification:

- A novel pipeline for fact verification that integrates state-of-the-art techniques in IR, NLP, ML.
- A comprehensive evaluation of different retrieval methods, embedding techniques, and language models for the task of fact verification.
- New strategies for generating effective queries and processing retrieved information to support fact verification.
- Analysis of the advantages and drawbacks of employing LLMs for reasoning regarding factual knowledge.
- A detailed analysis of the system's performance across different types of facts and knowledge domains.

## 1.5 THESIS STRUCTURE

The remainder of this thesis is organized as follows:

Chapter 2 provides a comprehensive review of the related works in fact verification, IR, and language model applications through LLMs. It situates our work within the broader context of these research areas and highlights the gaps that our approach aims to address.

Chapter 3 presents a detailed description of our proposed pipeline for fact verification. It explains each component of the system, including the rationale behind design choices and implementation details.

Chapter ?? describes the experimental setup used to evaluate our system. This includes details on the datasets used, evaluation metrics, and baseline systems for comparison and offers an in-depth discussion of the results, exploring the implications of our findings and their potential impact on the field of knowledge graph fact verification.

Chapter 5 Presents a study that investigates the impact of various pipeline components on the system's overall performance, while also exploring different methodologies for each component to determine the optimal final pipeline configuration.

Finally, chapter 6 concludes the thesis by summarizing the key contributions, discussing limitations of the current approach, and outlining promising directions for future research.

## 1.6 SIGNIFICANCE AND POTENTIAL APPLICATIONS

The development of effective fact verification systems for knowledge graphs has far-reaching implications across various domains:

- **Information Integrity:** Our solution facilitates the automatic verification of facts, hence enhancing the correctness and dependability of extensive knowledge bases. This is especially crucial in a time when disinformation may disseminate swiftly online.
- **Decision Support:** In sectors such as healthcare, finance, and law, where judgments frequently depend on factual information, our technology could function as an essential instrument for validating crucial data points.
- **Educational Applications:** Fact verification systems can be used in educational settings to help students critically evaluate information and develop digital literacy skills.

- **Content Moderation:** Social media platforms and content aggregators may employ similar techniques to detect and flag potentially inaccurate or misleading information.
- **Scientific Research:** In the scientific community, our approach could assist in fact-checking research claims, cross-referencing findings, and identifying potential inconsistencies in the literature.

By addressing the challenge of knowledge graph fact verification, this thesis aims to contribute to the broader goal of creating more reliable and trustworthy information systems. As the volume and complexity of digital information continue to grow, the need for sophisticated fact verification tools becomes increasingly critical. Our work represents a step towards meeting this need, combining the latest advances in artificial intelligence with rigorous IR techniques. In the following chapters, we will delve into the technical details of our approach, present our findings, and explore the implications of this research for the future of knowledge management and information verification.



# 2

## Related Works

This thesis builds upon prior research in knowledge graph fact verification, LLMs, Retrieval-Augmented Generation (RAG), and entailment verification. In this section, we provide an overview of the relevant literature across these areas.

### 2.1 ENTAILMENT VERIFICATION AND LANGUAGE MODELS

In the paper "Minds versus Machines: Rethinking Entailment Verification with Language Models", Sanyal et al. [30] evaluate and compare the inference capabilities of humans and LLMs through a carefully constructed entailment verification benchmark. Their study spans three categories: Natural Language Inference (NLI), contextual Question Answering (QA), and rationales, using multi-sentence premises and diverse types of knowledge to assess inference across complex reasoning scenarios.

The authors found that LLMs generally excel in multi-hop reasoning tasks, particularly those requiring inference over extended contexts, while humans outperform LLMs in simpler deductive reasoning tasks involving substitutions or negations.

Interestingly, both perform comparably in situations requiring inference of missing knowledge. One of the paper's key contributions is the fine-tuning of the Flan-T5 [3] model, which outperforms GPT-3.5 and performs at a comparable level to GPT-4, thus providing a robust, open-source solution for entailment verification tasks. In contrast, the proposed approach to factulizing the

## 2.1. ENTAILMENT VERIFICATION AND LANGUAGE MODELS

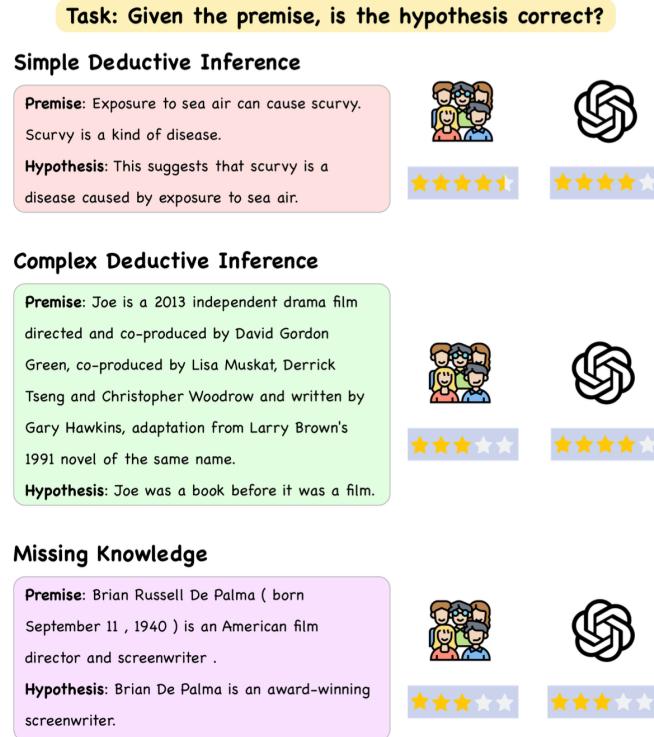


Figure 2.1: **Distinctions between human and LLM Inferences.** The entailment prediction performance of humans and LLMs are depicted by a 5-star rating scale [30].

knowledge graph using RAG emphasizes the integration of external knowledge retrieval to ground factual assertions, which is critical for generating verifiable, accurate knowledge graphs. While Sanyal et al. focus on the entailment between premises and hypotheses in textual inference, my work extends this by incorporating external evidence to ensure not just consistency but also factual correctness.

In comparison, the entailment verification tasks handled by Sanyal et al. emphasize reasoning within the constraints of the given context, whereas my RAG-based approach highlights the necessity of retrieval from large external datasets to mitigate hallucinations and improve the factual grounding of generated content. Both approaches deal with inference verification but diverge in their method of contextualizing and validating knowledge, with mine incorporating real-time retrieval for fact-checking.

This distinction is significant in terms of application: while their fine-tuned Flan-T5 model achieves high accuracy in entailment tasks, it remains bound

to the contextual limits of its training data. My work, by integrating retrieval, potentially overcomes this limitation by dynamically accessing external data, thus offering a complementary perspective to entailment verification focused on enhancing factuality.

## 2.2 CLAIM VERIFICATION IN THE AGE OF LARGE LANGUAGE MODELS

Dmonte et al. [7] provide a comprehensive survey of LLM-based approaches to claim verification, highlighting the shift from traditional NLP methods to more sophisticated LLM-driven techniques.

The typical LLM-based claim verification pipeline, as described by Dmonte et al., consists of several key components:

1. Evidence Retrieval: Utilizing techniques like RAG to fetch relevant information from external sources.
2. Prompt Creation: Developing effective prompting strategies to guide LLMs in processing claims and evidence.
3. Transfer Learning: Employing fine-tuning and in-context learning to adapt LLMs to the specific task of claim verification.
4. LLM Generation: Using LLMs to generate veracity labels, supporting evidence, and explanations.

This pipeline represents a departure from traditional fact-checking approaches, leveraging the power of LLMs to improve accuracy and provide more nuanced assessments of claim veracity.

Several studies have demonstrated the effectiveness of LLM-based approaches in claim verification:

- Zhang and Gao [43] introduced HiSS, a hierarchical prompting technique that improved performance on complex news claim verification tasks.
- Lee et al. [18] developed FactualityPrompts, a framework for assessing the factual accuracy of LLM-generated content.

These studies consistently show that LLM-based methods outperform traditional NLP approaches in terms of accuracy, flexibility, and the ability to handle complex claims.

Our approach shares similarities with the LLM-based pipeline described by Dmonte et al., particularly in the use of retrieval-augmented generation and

## 2.2. CLAIM VERIFICATION IN THE AGE OF LARGE LANGUAGE MODELS

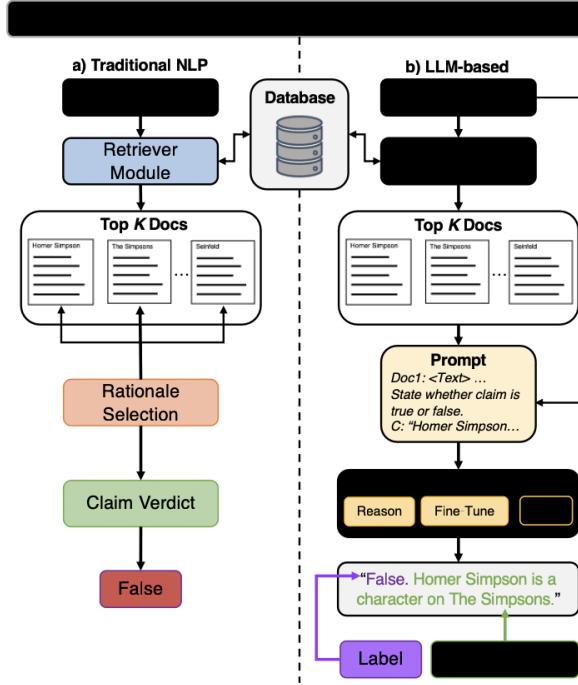


Figure 2.2: Comparison of claim verification systems between NLP-based (traditional) and LLM-based for claim veracity. [7].

the integration of multiple LLMs. However, our method differs in several key aspects:

1. Multi-Query Retrieval: We employ a multi-query strategy for evidence retrieval, potentially improving the coverage and relevance of supporting information.
2. Iterative Refinement: Our system incorporates an iterative process for refining retrieved evidence and generated responses, which is not explicitly mentioned in most LLM-based approaches surveyed.
3. Explainability Focus: While many LLM approaches provide explanations, our method places a stronger emphasis on generating detailed, step-by-step justifications for veracity assessments.
4. Specialized Embedding Models: We utilize domain-specific embedding models for improved semantic understanding in the retrieval phase, which is not commonly discussed in general LLM-based approaches.

These distinctions position our work as a novel contribution to the field, building upon the foundations of LLM-based claim verification while introducing innovative techniques to enhance performance and interpretability.

Despite the promising results of LLM-based claim verification, several challenges remain. Dmonte et al. highlight issues such as handling irrelevant

context, resolving knowledge conflicts, and expanding to multilingual settings. Our approach attempts to address some of these challenges, particularly in the areas of context relevance and explainability. However, there is still significant room for improvement in creating more robust, reliable, and universally applicable claim verification systems.

## 2.3 RETRIEVAL-AUGMENTED FACT VERIFICATION BY SYNTHESIZING CONTRASTIVE ARGUMENTS

The paper Retrieval-Augmented Fact Verification by Synthesizing Contrastive Arguments [40] explores a method for improving fact verification in knowledge graphs using RAG. The proposed framework combines retrieval of external information and the generation of contrastive arguments—claims supported by retrieved evidence, but also those that provide counterpoints. This dual synthesis provides a richer and more nuanced verification process, allowing the system to handle conflicting evidence more effectively. The core contribution of the work lies in the creation of contrastive arguments, a strategy designed to reduce errors in fact verification systems, especially when LLMs may hallucinate or generate incomplete reasoning.

The authors leverage a multi-stage pipeline where external documents are retrieved to support or refute a given claim. Each retrieved piece of evidence is evaluated using a neural network model that ranks the evidence based on its relevance to the claim. By synthesizing contrastive arguments, the system generates explanations for both supporting and refuting the claim, which helps improve the transparency and trustworthiness of the model’s decisions.

In terms of results, the framework shows improvement over traditional fact verification pipelines, particularly in handling ambiguous or conflicting information. The contrastive arguments allow for better handling of cases where facts are not binary but exist in a more complex, nuanced state. The system’s ability to generate arguments for both sides of a claim increases its robustness and provides a more reliable fact verification tool.

The described approach and my work on fact verification in knowledge graphs using RAG share a common goal: improving the factual accuracy of information through the integration of external knowledge retrieval. However, there are key differences in the methodologies used. The contrastive argument

## 2.4. RAGAR: RAG-AUGMENTED REASONING FOR POLITICAL FACT-CHECKING USING MULTIMODAL LLMS

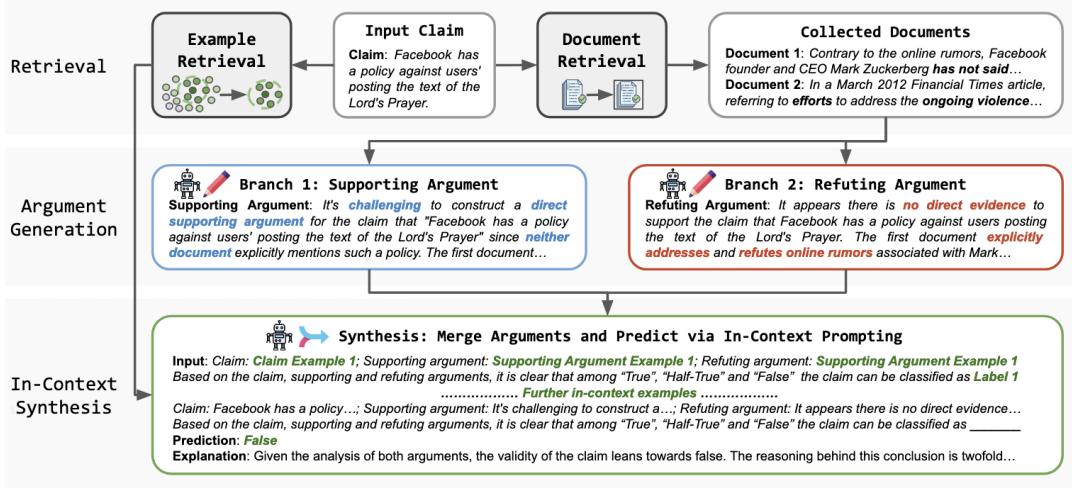


Figure 2.3: The proposed RAFTS [40], which performs few-shot fact verification by incorporating informative in-context demonstrations and contrastive arguments with nuanced information derived from the retrieved documents

synthesis introduced by the authors focuses heavily on generating both supporting and opposing arguments for claims, which provides a more holistic perspective in scenarios where evidence is mixed. In contrast, my approach emphasizes majority voting among multiple models and a multi-query strategy to retrieve a broader range of external evidence, aiming to reduce the incidence of hallucinations in LLM outputs.

While both approaches use retrieval to mitigate the limitations of LLMs, my work incorporates adaptive dispute resolution techniques and focuses on synthesizing outputs from multiple LLMs rather than generating contrastive arguments. This means that my approach leans more towards optimizing model diversity and utilizing the best consensus from several LLMs to ensure factual accuracy, rather than explicitly generating opposing arguments for each claim.

## 2.4 RAGAR: RAG-AUGMENTED REASONING FOR POLITICAL FACT-CHECKING USING MULTIMODAL LLMS

The study titled RAGAR: RAG-Augmented Reasoning for Political Fact-Checking using Multimodal LLMs [16] introduces a novel approach to political fact-checking by leveraging RAG with multimodal LLMs. This work focuses on enhancing fact verification in the politically sensitive domain, where dis-

formation can have far-reaching consequences. The authors integrate various modalities text, images, and other media sources into a unified fact-checking pipeline powered by LLMs, particularly emphasizing RAG’s ability to retrieve and synthesize external evidence to validate or refute claims.

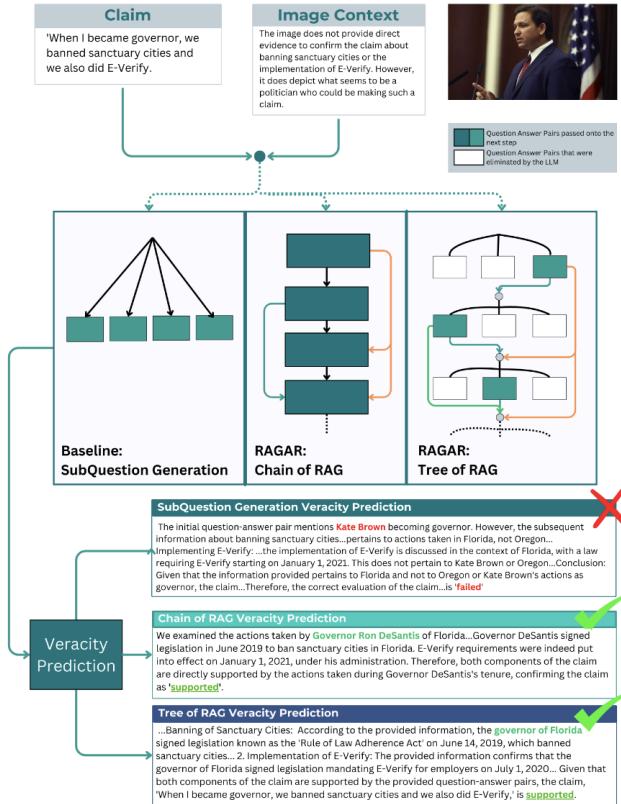


Figure 2.4: An overview of the fact-checking pipeline [16] contrasting the baseline Sub-Question Generation approach from the RAGAR: Chain of RAG and RAGAR: Tree of RAG approach followed by a veracity explanation generated by a Veracity Prediction module.

The central innovation of RAGAR lies in its multimodal reasoning capabilities, which allow the model to handle political claims that involve not only textual content but also visual data, such as images or charts. By extending RAG to this multimodal context, the system improves its ability to assess the veracity of claims in real-time, leveraging external resources such as political databases and live web content. Furthermore, the use of contrastive learning helps the system generate both supporting and opposing arguments for each claim, providing a more balanced and comprehensive fact-checking process.

Results from the paper show significant improvements in fact-checking ac-

## 2.4. RAGAR: RAG-AUGMENTED REASONING FOR POLITICAL FACT-CHECKING USING MULTIMODAL LLMS

curacy, particularly for politically charged claims that are often more nuanced or context-dependent. RAGAR’s ability to synthesize multimodal evidence into a coherent verification report highlights its potential for real-world applications, especially in environments where disinformation spreads quickly, such as social media platforms.

RAGAR focuses on political fact-checking using multimodal data, whereas my work targets the factualization of knowledge graphs with a primary focus on textual information. In contrast to RAGAR’s multimodal pipeline, my system emphasizes multi-query strategies and document chunking techniques to retrieve highly relevant textual evidence for verification.

# 3

## Pipeline

This chapter presents a detailed examination of a multi-stage NLP pipeline designed to enhance information retrieval and question answering capabilities. The pipeline integrates various cutting-edge technologies and methodologies, creating a synergistic system that pushes the boundaries of what is possible in automated information processing and response generation. However, the challenges of accurately interpreting user queries, retrieving relevant information from vast datasets, and generating coherent and contextually appropriate responses remain significant. This pipeline addresses these challenges through a carefully orchestrated series of processes, each designed to refine and enhance the quality of information flow from input to output.

At its core, the pipeline leverages a knowledge graph dataset, serving as the foundational repository of interconnected information. This graph structure allows for the representation of complex relationships between entities, facilitating more nuanced understanding and retrieval of information. The pipeline then employs a series of sophisticated mechanisms, including query generation, cross-encoding for relevance assessment, and multi-tiered information retrieval strategies, to navigate this knowledge landscape effectively.

One of the key innovations in this pipeline is its approach to context processing and information synthesis. By breaking down retrieved information into manageable chunks and employing advanced embedding techniques, the system can perform more granular and accurate analyses of textual data. This granularity, combined with the implementation of multiple LLMs working in concert, allows for a more robust and nuanced interpretation of complex queries

and generation of comprehensive responses.

The integration of external information sources, notably through the incorporation of Google Search capabilities, further enhances the pipeline's ability to access and process up-to-date and diverse information. This hybrid approach, combining structured knowledge graphs with dynamic web-based information retrieval, positions the pipeline at the forefront of adaptive and responsive AI systems.

Critical to the pipeline's effectiveness is its sophisticated decision-making architecture. Employing strategies such as majority voting among multiple models and the implementation of a final judge for conflict resolution, the system strives to achieve a balance between diverse perspectives and the need for coherent, unified outputs. This approach not only enhances the accuracy and reliability of the generated responses but also provides a framework for managing the inherent uncertainties and potential biases in AI-driven decision-making processes.

As we delve deeper into each component of this pipeline, it is crucial to maintain a critical perspective on both its capabilities and limitations. While the system represents a significant advancement in NLP and information retrieval technologies, it also raises important questions about the ethical implications of AI-driven information processing, the potential for bias in knowledge representation and model training, and the broader societal impacts of increasingly sophisticated question-answering systems.

This chapter aims to provide a comprehensive analysis of each stage of the pipeline, examining not only the technical aspects of its implementation but also the theoretical underpinnings and practical implications of its design choices. By understanding the intricacies of this system, we can gain valuable insights into the current state of NLP technologies and the potential future directions for research and development in this rapidly advancing field. As we proceed, we will explore each component in detail, starting with the foundational knowledge graph dataset and progressing through the various stages of query processing, information retrieval, context analysis, and response generation. This exploration will shed light on the complex interplay between different AI technologies and methodologies, offering a holistic view of how modern NLP systems can be architected to tackle some of the most challenging problems in information processing and human-computer interaction.

## 3.1 KNOWLEDGE GRAPH DATASET

The foundation of our multi-stage RAG pipeline is the Knowledge Graph Dataset, which acts as the primary source of factual information for the ensuing processing stages. This section clarifies the attributes and aims of the knowledge graph and its use in our pipeline.

### 3.1.1 DEFINITION AND PURPOSE

In the form of a graph, a knowledge graph is an ordered way to show information that models real-world things and how they relate to each other. The Knowledge Graph Dataset is a huge collection of facts, ideas, and connections that are all linked together in our process. Its main job is to give a lot of background information so that complicated queries can be understood and processed.

The implementation of a knowledge graph fulfills multiple essential functions:

- **Semantic Representation:** Unlike traditional relational databases, knowledge graphs capture semantic relationships between entities, allowing for more nuanced and context-aware information retrieval.
- **Inferential Capabilities:** The interconnected nature of the graph enables the system to make inferences and connections that may not be explicitly stated, enhancing the depth and breadth of responses.
- **Scalability:** Knowledge graphs can efficiently handle large volumes of heterogeneous data, making them ideal for systems that need to process diverse types of information.
- **Flexibility:** The graph structure allows for easy updates and expansions, ensuring that the knowledge base can evolve with new information and changing requirements.

### 3.1.2 STRUCTURE AND COMPONENTS

The Knowledge Graph Dataset in our pipeline is composed of several key components:

- **Nodes:** Representing entities or concepts, nodes are the fundamental units of information in the graph. Each node typically corresponds to a distinct piece of knowledge, such as a person, place, event, or abstract concept.

### 3.2. QUERY GENERATION AND PROCESSING

- **Edges:** These are the connections between nodes, representing relationships or interactions. Edges are often directional and labeled to indicate the nature of the relationship (e.g., "is\_a", "part\_of", "created\_by").
- **Properties:** Nodes and edges can have associated properties or attributes that provide additional details or metadata about the entity or relationship.

#### 3.1.3 ROLE IN THE OVERALL PIPELINE

As illustrated in the pipeline diagram, the Knowledge Graph Dataset is the thing that we want to verify the correctness of it. The pipeline uses the knowledge graph to generate queries, retrieve relevant information, and synthesize responses to find the correctness of the knowledge graph.

### 3.2 QUERY GENERATION AND PROCESSING

The Query Generation and Processing step, following to the basic Knowledge Graph Dataset, is a pivotal point in our pipeline, wherein user inputs are converted into structured queries suitable for efficient processing by later components. This section clarifies the techniques and methodologies utilized in this critical phase for subsequent actions in the information retrieval process.

#### 3.2.1 HUMAN-UNDERSTANDABLE TEXT GENERATION

In the first step of the pipeline, we use a LLM (ie. LLama3 [8], Gemma2 [34]) to easily generate human-readable text by submitting prompts. This approach bridges the gap between representing raw data and conveying it in natural language, making the information more accessible and understandable. For additional information, consult the prompt template in Appendix A.1 to observe the sentence generation process.

Key aspects of this process include:

- **Contextual Awareness:** Integrating relevant context from the Knowledge Graph to guarantee that the output content is relevant and useful.
- **Adaptability:** Customizing the generated text to accommodate various complexity levels, catering to diverse user needs and query types.
- **Semantic Enrichment:** Enhancing the generated text with semantic annotations to facilitate more accurate downstream processing.

Be aware that certain knowledge graph datasets are human-readable, whereas others are not; for instance, the FactBench [10] dataset may be easily comprehended by concatenating the subject, predicate, and object of each triple.

Knowledge Graph			Source	Is Generated ?	Final Text
Subject	Predicate	Object			
<i>Albert Einstein</i>	<i>Birth Place</i>	<i>Ulm, Germany</i>	FactBench [10]	✗	<i>Albert Einstein birth place Ulm, Germany</i>
<i>Chris Benoit</i>	<i>deathPlace</i>	<i>Fayetteville, Georgia</i>	FactBench [10]	✗	<i>Chris Benoit death place Fayetteville, Georgia</i>
<i>Alexander III of Russia</i>	<i>isMarriedTo</i>	<i>Maria_Feodorovna</i> <i>—Dagmar_of_Denmark—</i>	YAGO [33]	✓	<i>Alexander III of Russia is married to Maria Feodorovna, also known as Dagmar of Denmark.</i>
<i>Shock_to_the_System_(Gemma_Hayes_song)</i>	<i>length</i>	221.0	DBpedia	✓	<i>The length of Shock to the System (Gemma Hayes song) is 221.0.</i>
<i>Paora_Winitana</i>	<i>years</i>	2011	DBpedia	✓	<i>Paora Winitana was active in 2011.</i>

Table 3.1: Examples of Generation of Human-Understandable Text

### 3.2.2 QUESTION FORMULATION TECHNIQUES

A cornerstone of our pipeline is its ability to generate 10 questions about the input sentence, as illustrated in the diagram.

This multi-question approach serves several purposes:

- **Comprehensive Coverage:** By generating multiple questions, the system ensures a thorough exploration of the input's various aspects and potential interpretations.
- **Disambiguation:** Multiple questions help in clarifying ambiguities that may be present in the original input.
- **Context Expansion:** Each generated question potentially introduces new contextual elements, broadening the scope of the subsequent information retrieval process.
- **Robustness:** The diversity of questions increases the likelihood of capturing the user's true intent, even if the original input is vague or imprecise.

### 3.2. QUERY GENERATION AND PROCESSING

Implementation of this technique likely involves: Using LLMs to generate 10 questions about the input sentence, leveraging the models' language understanding capabilities to ensure the questions are relevant and contextually appropriate. The prompt template used to guide the question generation process reported in Appendix A.2.

#### 3.2.3 CROSS-ENCODER FOR QUERY RELEVANCE SCORING

The Cross-Encoder component is essential for evaluating the relevance of the generated questions. As indicated in the Figure 3.1, this module takes multiple inputs and produces relevance scores for each question.

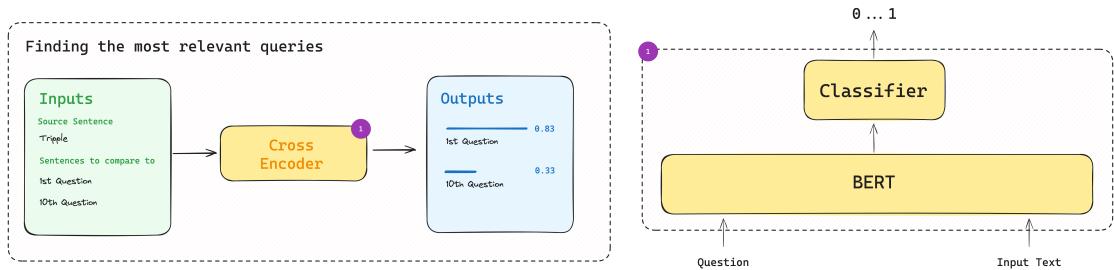


Figure 3.1: Cross-Encoder Architecture

Key features of the Cross-Encoder include:

- Input Processing
  - Source Sentence: The original input text.
  - Sentences to Compare: Likely the 10 generated questions.
- Scoring Mechanism: The Cross-Encoder assigns numerical scores (e.g., 0.83 as shown in the figure 3.1) to each question, indicating its relevance to the source sentence.
- Comparative Analysis: By processing all inputs simultaneously, the Cross-Encoder can perform nuanced comparisons between the original input and each generated question, as well as among the questions themselves.

In this case we use the *jinaai/jina-reranker-v1-turbo-en*<sup>1</sup> model from the Hugging Face Transformers library. This model is designed for blazing-fast re-ranking while maintaining competitive performance. It leverages the power of

<sup>1</sup><https://huggingface.co/jinaai/jina-reranker-v1-turbo-en>

JinaBERT [11] model as its foundation. The model employs a process known as knowledge distillation to attain exceptional speed and efficiency, making it the optimal selection for our pipeline.

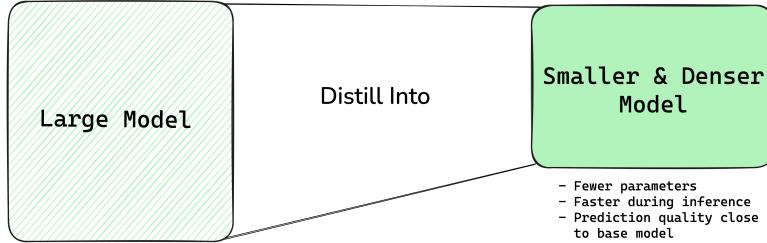


Figure 3.2: Knowledge Distillation Process

Input	Question	score
<i>Frédéric Passy award Nobel Peace Prize</i>	Who was awarded the Nobel Peace Prize?	0.7458
	What award did Frédéric Passy receive?	0.7121
	Is Frédéric Passy a Nobel laureate?	0.8706
	In what category was the Nobel Peace Prize awarded to Frédéric Passy?	<b>0.9491</b>
	Who is known for receiving the Nobel Peace Prize?	0.5823
	What is the name of the award received by Frédéric Passy?	0.6318
	Is Frédéric Passy a recipient of the Nobel Prize in any field?	0.7652
	Who was recognized for his work towards peace?	0.1457
	What is the significance of the award given to Frédéric Passy?	0.6036
	Is there a Nobel laureate with the name Frédéric Passy?	0.8505

Table 3.2: Question Generation and Scoring Procedure

### 3.2.4 RELEVANCE THRESHOLD AND SORTING

Following the Cross-Encoder's scoring, the pipeline implements a crucial decision point:

- **Sorting:** Questions are sorted based on their relevance scores, establishing a priority order for further processing.
- **Threshold Evaluation:** The system checks if the top question's score exceeds an upper threshold. This step ensures that only sufficiently relevant questions proceed further in the pipeline.
- **Feedback Loop:** If the threshold is not met, the process must loop back to generate new questions to adjust the existing ones, maintaining the quality of queries entering subsequent stages.

## 3.3 INFORMATION RETRIEVAL MECHANISMS

The Information Retrieval Mechanisms are an essential element of our pipeline, connecting query processing and content synthesis. This phase is tasked with gathering relevant data from internal and external sources, thereby establishing a comprehensive data repository for further analysis and response formulation. The mechanisms employed in this phase are designed to ensure breadth, depth, and relevance in the retrieved information.

### 3.3.1 GOOGLE SEARCH INTEGRATION

A key feature of our information retrieval process is the integration of Google Search capabilities, as prominently displayed in the pipeline diagram. This integration serves to expand the information horizon beyond the confines of our internal Knowledge Graph Dataset. Key aspects of this integration include:

- **Query Submission:** The system submits the N top questions (where N is a predefined number) along with the main question to Google Search. This approach ensures a multi-faceted search that captures various aspects of the original query.
- **Result Fetching:** As indicated in the diagram, the system retrieves the top 100 search results. This number strikes a balance between comprehensiveness and computational efficiency.
- **Dynamic Information Access:** By leveraging Google Search, the system gains access to up-to-date information, complementing the more static nature of the internal Knowledge Graph.
- **Diverse Source Types:** Google Search results typically include a variety of source types (e.g., websites, news articles, academic papers), enriching the diversity of the retrieved information.

Implementation considerations:

- The system may employ proxies or other mechanisms to manage rate limits and ensure uninterrupted access to search results.
- The search query may be customized based on the specific requirements of the pipeline, such as language restrictions, geolocation preferences, or result quantity. parameters are lr, hl, gl, and num.

Take note that the code base is generic, and it means that you can use any search engine, not only Google Search.

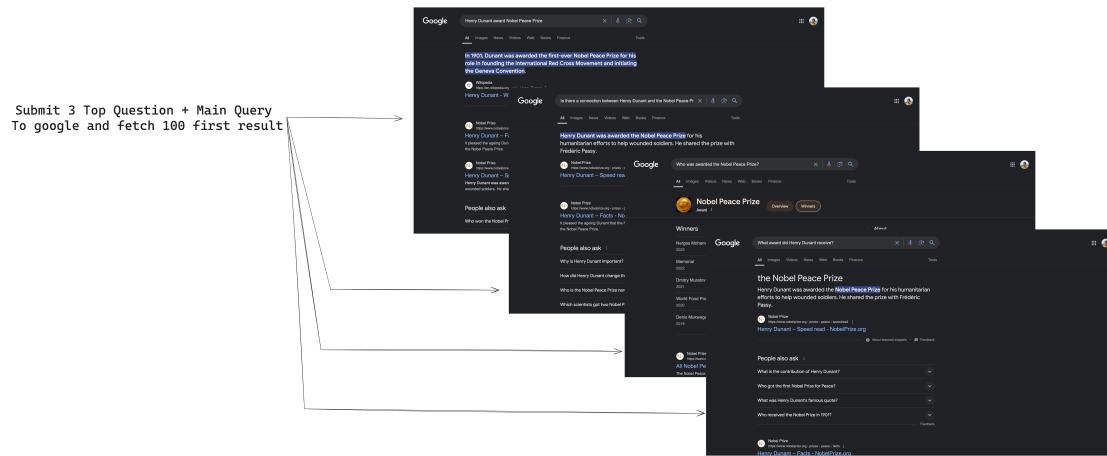


Figure 3.3: Fetching the results from Google Search

### 3.3.2 PROCESS AND EXTRACT LINKS

Following the retrieval of search results, the pipeline incorporates a crucial step of processing and extracting links from the gathered information. This process likely involves:

- **Parsing HyperText Markup Language (HTML) Content:** Extracting relevant textual information from the retrieved web pages.
- **Link Analysis:** Identifying and cataloging hyperlinks within the content, potentially uncovering additional relevant sources.

After Parsing the HTML content of the Google Search results, the system use HTML selectors to extract the links from the search results. In this case we also extract the title, url, description, price, date, duration, missing, rating, availability, and extra details from the search results. Some of the information mentioned above may not be available as it depends on the search results.

Then there is a need to crawl the extracted links to get the content of the page, for doing this we use the Python library called *GRequests*<sup>2</sup>. *GRequests* is a Python library that combines the power of gevent for asynchronous I/O with the simplicity of the Requests library for HTTP operations. It allows developers to perform concurrent HTTP requests easily, significantly speeding up operations that involve multiple API calls or web scraping tasks.

```
1 import os
2 import grequests
```

---

<sup>2</sup><https://pypi.org/project/grequests/>

### 3.3. INFORMATION RETRIEVAL MECHANISMS

```
3 from fake_useragent import UserAgent
4
5 ua = UserAgent(
6     os=['windows'],
7     browsers=["chrome", "edge", "firefox"],
8     platforms=["pc"]
9 )
10
11 urls = [...List of Extracted URLs...]
12
13 rs = [
14     grequests.get(u['url'],
15     timeout=3, headers={"User-Agent": ua.random}) for u in urls
16 ]
17 for index, response in grequests.imap_enumerated(rs, size=50):
18     if response is None or response.status_code != 200:
19         continue
20     # Process the response content ...
```

Code 3.1: Crawling the Extracted URLs

There are several faults in the mentioned approach 3.1 that need to be addressed:

- **Site generated with javascript:** The provided approach does not handle sites that are generated with JavaScript and require dynamic rendering.
- **Protection against scraping:** Sites may have protection mechanisms against scraping, such as CAPTCHAs or IP blocking or behind spam protection services.
- **Login required:** Some sites require login credentials to access the content.

We can use the *Selenium*<sup>3</sup> library to handle the first issue, and for the second and third issues, we can use the *Scrapy*<sup>4</sup> library, but we stick with the provided approach for simplicity and speed.

With the extracted content, the system can now proceed to the next stage of the pipeline, where the information is further processed and analyzed. For extracting the content of the page, as our webpage are from vast sources, we need to use a robust and efficient library to extract the content of the page. We use *newspaper4k*<sup>5</sup> library to extract the content of the page. *newspaper4k* is a Python

---

<sup>3</sup><https://www.selenium.dev/>

<sup>4</sup><https://scrapy.org/>

<sup>5</sup><https://newspaper4k.readthedocs.io/en/latest/>

library designed for extracting and parsing newspaper articles. It's an updated and improved version of the original newspaper3k library, offering enhanced functionality and compatibility with modern Python versions.

Key features of the *newspaper4k* library include:

- **Article Extraction:** Easily extract articles from news websites.
- **Multi-language Support:** Capable of processing articles in various languages.
- **Full-text Extraction:** Extracts the full text of articles, removing ads and extraneous content.
- **Keyword Extraction:** Automatically identifies key topics and keywords from articles.
- **Summary Generation:** Creates concise summaries of article content.
- **Metadata Parsing:** Extracts metadata such as authors, publication dates, and tags.
- **Image Extraction:** Identifies and extracts images associated with articles.

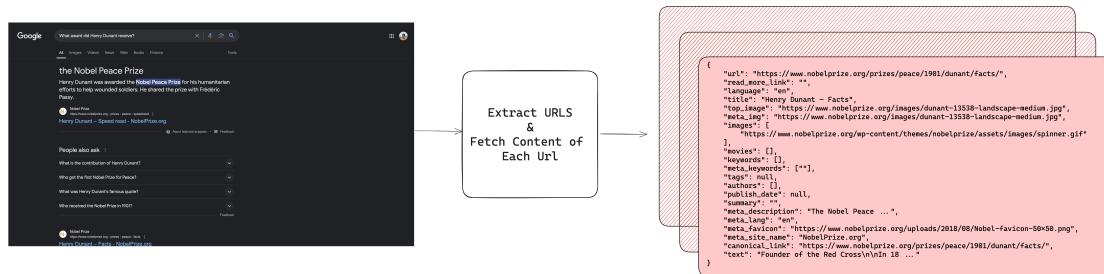


Figure 3.4: Extracted Content from the Crawled URLs using *newspaper4k*

For our purpose, we only use the text content of the page.

### 3.3.3 DATA POOL CREATION

The processed and extracted information culminates in the creation of a data pool, a centralized repository of relevant information that serves as the foundation for subsequent stages of the pipeline.

Key features of the data pool include:

- **Structured Storage:** Organizing the retrieved information in a format that facilitates efficient querying and analysis.
- **Source Diversity:** Maintaining a balance between information from web searches and the internal Knowledge Graph.

### 3.4. EMBEDDING AND RETRIEVAL TASKS

- **Relevance Scoring:** Potentially implementing a scoring system to prioritize more relevant or authoritative pieces of information within the pool.
- **Deduplication:** Employing mechanisms to identify and merge duplicate or highly similar pieces of information to reduce redundancy.

#### 3.3.4 DATA FILTERING

The data filtering process aims to refine our data pool by removing information from sources that are already part of creation of the Knowledge Graph. This ensures the uniqueness and independence of our data, for example, when working with the FactBench dataset, we exclude data from the following sources: wikipedia.org, wikimedia.org, wikidata.org and other similar wiki-based platforms.

On the other hand, we just keep the top 10 relevant information for the pipeline using cross-encoders discussed in Section 3.2.3, this way we ensure about the quality of the data. More details about the data filtering process methodology are provided in the section ??.

Specifically:

- We start with a larger pool of data.
- We identify sources that are already the source of the Knowledge Graph.
- We remove any data points that come from these identified sources.

### 3.4 EMBEDDING AND RETRIEVAL TASKS

An important step in our pipeline is the Embedding and Retrieval Tasks, which connect machine-interpretable vector representations to unprocessed textual input. This component is essential for improving the efficiency and accuracy of information retrieval and subsequent processing stages.

This section emphasizes embedding for retrieval tasks, specifically for small information segments, as depicted in the pipeline diagram.

### 3.4.1 EMBEDDING TECHNIQUES FOR SMALLER CHUNKS

The pipeline employs advanced embedding techniques to transform textual data into dense vector representations, facilitating more efficient and semantically aware retrieval processes.

Key aspects of this embedding process include:

- **Granularity:** The focus on "Smaller Chunks" suggests a fine-grained approach to embedding, where text is broken down into manageable units. This granularity allows for more precise retrieval and relevance assessment.
- **Dimensionality Reduction:** Transforming high-dimensional textual data into lower-dimensional vector spaces while preserving semantic relationships.

we use the *SentenceWindowNodeParser* to parse documents into single sentences per node. Each node also contains a "window" with the sentences on either side of the node sentence. Then, after retrieval, before passing the retrieved sentences to the LLM, the single sentences are replaced with a window containing the surrounding sentences using the *MetadataReplacementNodePostProcessor*. This is most useful for large documents/indexes, as it helps to retrieve more fine-grained details.

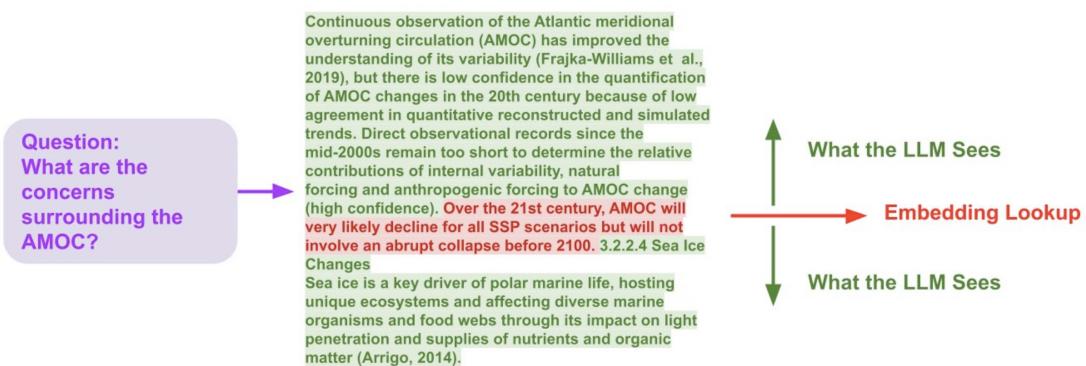


Figure 3.5: Node sentence window replacement [20]

There is example of the *SentenceWindowNodeParser* and *MetadataReplacementNodePostProcessor* in the Appendix ??.

### 3.4. EMBEDDING AND RETRIEVAL TASKS

#### CHALLENGES AND CONSIDERATIONS

Several challenges and considerations are inherent in the Embedding and Retrieval Tasks:

- **Multilingual Support:** Extending embedding capabilities to support multiple languages, potentially through multilingual or language-agnostic models.
- **Embedding Interpretability:** Balancing the trade-off between the performance of dense embeddings and the interpretability often associated with sparse representations.
- **Temporal Dynamics:** Addressing the challenge of embedding temporally sensitive information and ensuring retrieval mechanisms account for time-based relevance.
- **Scalability:** Designing the embedding and retrieval systems to efficiently handle growing volumes of data without compromising on speed or accuracy.
- **Ethical Considerations:** Mitigating potential biases inherent in pre-trained embedding models and ensuring fair representation across diverse topics and perspectives.

#### 3.4.2 SIMILARITY CUTOFF STRATEGY

One of the most important aspects that is highlighted in the pipeline diagram is the "Similarity Cutoff Strategy," which is significant since it is responsible for filtering and prioritizing the information that is included.

This strategy likely involves:

- **Threshold Definition:** Establishing a similarity threshold below which embedded chunks are considered insufficiently relevant or related to the query or context.
- **Similarity Metrics:** Employing appropriate similarity measures (e.g., cosine similarity, Euclidean distance) to quantify the relatedness between embedded representations.
- **Re-run Mechanism:** Potentially re-do the response generation process if the similarity cutoff is not met for all the documents and the response is not generated.

The Similarity Cutoff Strategy serves several critical functions:

- **Noise Reduction:** Filtering out irrelevant or tangentially related information to improve the signal-to-noise ratio in subsequent processing stages.

- **Computational Optimization:** Reducing the volume of data processed in later stages, thereby enhancing overall system efficiency.
- **Relevance Enhancement:** Ensuring that only the most pertinent information is retained for query resolution and response generation.

It can be switched off if the dataset is small or if the similarity cutoff is not needed.

## 3.5 LLMs

LLMs play a pivotal role in our advanced NLP pipeline, serving as the cornerstone for sophisticated information synthesis and response generation. As illustrated in the pipeline diagram, multiple LLMs are employed, contributing to a robust and nuanced approach to question answering and information processing.

### 3.5.1 INTEGRATION OF MULTIPLE LLMs

The pipeline incorporates multiple LLMs working in concert, a design choice that offers several significant advantages:

- **Diversity of Perspectives:** By utilizing multiple models, the system can capture a broader range of interpretations and approaches to information synthesis.
- **Specialization:** Different LLMs may be fine-tuned or specialized for particular types of queries or domains, allowing for more targeted and accurate responses in specific contexts.
- **Robustness:** The multi-model approach provides redundancy and helps mitigate individual model biases or weaknesses.
- **Scalability:** Parallel processing of information through multiple LLMs can potentially improve the system's throughput and response time.

Implementation considerations:

- **Model Selection:** Choosing a diverse set of LLMs that complement each other in terms of strengths and specializations.
- **Load Balancing:** Implementing efficient mechanisms to distribute work-load across the available models.
- **Version Management:** Maintaining and updating multiple LLMs to ensure they remain current and aligned with the latest advancements in NLP.

### 3.5. LLMS

For running open-source LLMs, we use *Ollama*<sup>6</sup>, Ollama is an open-source project that simplifies the process of setting up, running, and using large language models (LLMs) locally on your machine. It provides a user-friendly area for managing and interacting with various LLMs, making it easier for developers and enthusiasts to experiment with AI without relying on cloud services. Under the hood, *Ollama* is just an API server written in go that serves GGUF models via llama.cpp<sup>7</sup> and a centralized hub of models/settings.

Performance Considerations and Limitations:

- Performance Considerations
  - Performance depends on your hardware, especially CPU and GPU capabilities.
  - Larger models require more RAM and storage space.
  - GPU acceleration can significantly improve inference speed.
- Limitations
  - Resource intensive for larger models.
  - May not match the performance of cloud-based solutions for some use cases.
  - Limited to models that are compatible with Ollama's framework.

#### 3.5.2 ROLES OF LLMs IN THE PIPELINE

Based on the pipeline diagram, the LLMs serve several crucial functions:

- **Information Synthesis:** Integrating and coherently combining information from various sources.
- **Context Processing:** Analyzing and interpreting the broader context of queries and retrieved information to generate more accurate and relevant responses.
- **Reasoning and Inference:** Drawing logical conclusions and making inferences based on the available information, potentially finding the correctness of the knowledge graph.

The prompt template used for RAG approach reported in Appendix A.3, we use the selected chunks from the previous steps to feed the LLMs context and also provide few examples to the LLMs to generate the better response.

---

<sup>6</sup><https://ollama.com/>

<sup>7</sup><https://github.com/ggerganov/llama.cpp>

## 3.6 VOTING SYSTEM AND CONFLICT RESOLUTION

Our state-of-the-art pipeline for natural language processing relies heavily on the integration of many models and the deployment of procedures to resolve conflicts. To ensure robust and trustworthy outcomes, this section covers the ways adopted to resolve disputes and harness model diversity.

### 3.6.1 MAJORITY VOTING SYSTEM

A significant aspect of the LLM integration, is the establishment of a majority voting mechanism for response creation.

This method provides numerous advantages:

- **Consensus Building:** By aggregating outputs from multiple models, the system can identify areas of agreement, potentially leading to more reliable responses.
- **Error Mitigation:** Outlier responses or errors from individual models can be identified and potentially filtered out through the voting process.
- **Confidence Scoring:** The degree of consensus among models can serve as a proxy for the confidence level of the generated response.
- **Handling Ambiguity:** In cases where there's no clear majority, the system can potentially flag the response as uncertain or requiring further clarification.

Implementation challenges:

- **Weighting Mechanism:** Determining whether all LLMs should have equal weight in the voting process or if some models should be prioritized based on their specific strengths or reliability.
- **Threshold Setting:** Establishing the criteria for what constitutes a "majority" and how to handle cases with no clear consensus.
- **Combining Diverse Outputs:** Developing methods to meaningfully aggregate potentially disparate outputs from different models into a coherent final response.

The pipeline incorporates several mechanisms for resolving conflicts that may arise from divergent model outputs:

### 3.6. VOTING SYSTEM AND CONFLICT RESOLUTION

#### 3.6.2 CONFLICT RESOLUTION STRATEGIES

As previously discussed 3.6.1, the system employs a majority voting approach among LLMs to determine the most appropriate response. This serves as a primary conflict resolution mechanism, leveraging the wisdom of the collective to mitigate individual model errors or biases.

##### FINAL JUDGE IMPLEMENTATION:

A crucial component in the conflict resolution process is the "final judge" module, as indicated in the pipeline diagram. This element plays a pivotal role in resolving conflicts and ensuring coherence in the system's outputs.

Key aspects of the final judge implementation:

- **Conflict Identification:** Detecting discrepancies or contradictions in outputs from different models.
- **Resolution Mechanisms:** Using predefined rules or algorithms to determine the final output based on the majority voting results.
- **Tie-Breaking:** Implementing strategies for scenarios where the majority voting system results in a tie or lacks a clear consensus.
- **Consistency Enforcement:** Ensuring that the final output maintains logical consistency and aligns with established knowledge bases.

The final judge module used a system named *Adaptive Conflict Resolution*, which is an adaptive approach to conflict resolution based on the specified settings.

##### ADAPTIVE CONFLICT RESOLUTION

The pipeline demonstrates an adaptive approach to conflict resolution, as evidenced by the following feature: "When two of the language models believe the answer is correct, and two believe it is wrong, we switch to a more advanced model, such as a commercial one or a model with more parameters." There are two ways to apply this adaptive method: 1) Directly selecting paid models (e.g., GPT-4), or 2) Using a more advanced, open-source model with more parameters (e.g., Gemma2:27b). For the second option, we use the algorithm shown in Algorithm 1. This algorithm uses the mapping between the already used and final models to select the final model. It utilizes parameters to identify either

the most stable model or the least stable model to choose based on a majority vote across the entire system.

This adaptive strategy offers several advantages:

- **Escalation Mechanism:** Provides a structured approach for handling ambiguous cases where simpler resolution methods are insufficient.
- **Resource Optimization:** Reserves the use of more advanced (and potentially more computationally expensive) models for cases that truly require their capabilities.
- **Accuracy Enhancement:** Leverages more sophisticated models to resolve complex conflicts, potentially leading to higher-quality outputs in challenging scenarios.
- **Flexibility:** Allows for the integration of specialized or proprietary models in a targeted manner, enhancing the system's overall capabilities without relying on these models for every query.

Implementation challenges:

- **Threshold Definition:** Determining the exact criteria for when to invoke the more advanced models.
- **Model Selection:** Choosing which advanced model to use based on the nature of the conflict and the query context.
- **Integration:** Ensuring smooth handover and result incorporation from the advanced models back into the main pipeline.

## 3.7 PIPELINE FLOW AND DECISION POINTS

The architecture of our pipeline is characterized by a sophisticated flow of information and a series of critical decision points. This structure enables the system to process complex queries, retrieve and synthesize relevant information, and generate accurate responses. This section provides a comprehensive analysis of the pipeline's flow and the key decision points that guide the processing of information.

The pipeline flow can be broadly categorized into several main stages, each with its own set of processes and decision points:

- Input Processing and Query Generation
- Information Retrieval and Enrichment
- Embedding and Relevance Assessment

---

**Algorithm 1** Resolve Ties in Majority Voting System

---

**Require:**

*data* - A list of dictionaries, each representing a model's response

*finalJudger* - A dictionary mapping file indices to final model

*atLeast* - A boolean flag:

True: select model with the highest agreement score

False: select model with the lowest agreement score

```

1: procedure PROCESSFILES(data, finalJudger, atLeast)
2:   majorityValues  $\leftarrow \emptyset$ 
3:   keys  $\leftarrow$  keys from first element of data
4:   for each key in keys do
5:     values  $\leftarrow$  list of values for key from all data
6:     count  $\leftarrow$  count occurrences of 1's and 0's in values
7:     if count of 1's > count of 0's then
8:       majorityValues[key]  $\leftarrow 1$ 
9:     else if count of 1's < count of 0's then
10:      majorityValues[key]  $\leftarrow 0$ 
11:    else
12:      continue to next key
13:    end if
14:   end for
15:   modelScores  $\leftarrow \emptyset$ 
16:   for i  $\leftarrow 0$  to  $|data| - 1$  do
17:     trueCount  $\leftarrow 0$ 
18:     for each (k, v) in data[i] do
19:       if majorityValues[k] = v then
20:         trueCount  $\leftarrow$  trueCount + 1
21:       end if
22:     end for
23:     Add (i, trueCount) to modelScores
24:   end for
25:   if atLeast is True then
26:     maxScore  $\leftarrow$  maximum score in modelScores
27:     candidates  $\leftarrow$  indices with score equal to maxScore
28:   else
29:     minScore  $\leftarrow$  minimum score in modelScores
30:     candidates  $\leftarrow$  indices with score equal to minScore
31:   end if
32:   chosenIndex  $\leftarrow$  random choice from candidates
33:   return finalJudger[chosenIndex]
34: end procedure

```

---

- Multi-Model Processing and Synthesis
- Conflict Resolution and Final Output Generation

Let's examine each of these stages in detail in table 3.3, focusing on the flow of information and the critical decision points within each.

Component	Decision Point
<i>Input Processing and Query Generation</i>	
Knowledge Graph Dataset Integration	Extent and nature of knowledge graph integration based on input complexity.
Human-Understandable Text Generation	Selection of the most appropriate natural language generation technique based on input and context.
Question Generation	Assessment of the quality and relevance of generated questions.
Cross-Encoder for Query Relevance	Determination of whether the top question score exceeds the upper threshold.
<i>Information Retrieval and Enrichment</i>	
Information Source Selection	Determining the most relevant and reliable sources for information retrieval.
Google Search Integration	Balancing between the breadth of search (number of questions submitted) and depth (number of results retrieved).
Process and Extract Links	Determining the relevance and quality of extracted information for inclusion in the data pool.
Data Pool Creation	Structuring the data pool for optimal accessibility in subsequent stages.
<i>Embedding and Relevance Assessment</i>	
Embedding for Retrieval Tasks	Selection of the most appropriate embedding technique based on the nature of the data.
Similarity Cutoff Strategy	Determination of the similarity cutoff threshold.
Context Processing	Determining the optimal chunk size and processing method.
<i>Multi-Model Processing and Synthesis</i>	
Parallel LLM Processing	Allocation of specific tasks or aspects of the query to different models based on their strengths.
Synthesis of Information	Determination of the method of synthesis (e.g., concatenation, abstraction, or hybrid approaches).
Majority Voting System	Assessment of the level of agreement among models.
<i>Conflict Resolution and Final Output Generation</i>	
Conflict Identification	Determination of the threshold for what constitutes a significant conflict requiring resolution.
Adaptive Model Selection	When two models believe the answer is correct and two believe it's wrong, switching to a more advanced model.
Selection of Commercial Models	Choose the commercial model based on the user specified settings.
Final Judge Implementation	Determination of the final response based on aggregated model outputs and conflict resolution results.
Response Generation	Selection of the most appropriate format and level of detail for the response.

Table 3.3: Comprehensive list of decision points in the pipeline flow.

Several challenges and considerations are associated with managing the pipeline flow and decision points:

### 3.8. PERFORMANCE REPORT

- **Computational Efficiency:** Balancing the depth of processing at each stage with the need for timely responses.
- **Error Propagation:** Ensuring that errors or biases introduced at early stages don't disproportionately affect the final output.
- **Adaptability:** Designing decision points that can adapt to different query types and complexity levels.
- **Transparency:** Maintaining traceability of decisions made throughout the pipeline for accountability and debugging purposes.
- **Scalability:** Ensuring that the pipeline can handle increasing query volumes without significant degradation in performance or accuracy.

In conclusion, the Pipeline Flow and Decision Points represent a complex yet well-structured approach to natural language processing and question answering. By implementing a series of carefully designed stages and critical decision points, the system aims to process information in a manner that maximizes accuracy, relevance, and reliability. The adaptive nature of the pipeline, particularly in its approach to conflict resolution and model selection, demonstrates a commitment to handling a wide range of query complexities and scenarios. However, the intricate nature of this flow also underscores the importance of ongoing optimization, monitoring, and refinement to ensure that the system continues to perform effectively and ethically in the face of evolving challenges and requirements in the field of AI and natural language processing.

## 3.8 PERFORMANCE REPORT

In Tables 3.4, 3.5, we provide a detailed performance report, highlighting key metrics that demonstrate the pipeline's efficiency. Other sections of the pipeline, such as RAG and conflict resolution mechanisms, will be evaluated in future chapters.

Task	Avg. Request Duration	Avg. tokens per request
Human understandable text 3.2.1	1.3164 sec	343.16
Question Generation 3.2.2	9.6076 sec	672.58

Table 3.4: Performance of our LLM applications in production, generated by Openlit<sup>8</sup> with *Gemma2* model.

Task	Avg. Duration
Sorting the questions by similarity 3.2.3	0.013 sec
Get documents (Google pages) 3.3.1	3.6 sec
Fetch documents for each triple 3.3.2	350 sec

Table 3.5: Performance of our Information Retrieval Mechanisms.

## 3.9 ETHICAL CONSIDERATIONS AND LIMITATIONS

The development and deployment of advanced pipelines, such as the one described in this thesis, necessitate a thorough examination of ethical considerations and an acknowledgment of system limitations. This section explores the ethical implications of our fact-checking system and discusses its inherent constraints.

### 3.9.1 ETHICAL CONSIDERATIONS

While the pipeline diagram does not explicitly highlight ethical components, several aspects of the system raise important ethical considerations:

#### MISINFORMATION AND HARMFUL CONTENT

The system's ability to synthesize information from various sources poses risks related to misinformation:

- **Propagation of False Information:** Potential for the system to inadvertently spread misinformation present in retrieved data.
- **Generation of Harmful Content:** Risk of producing responses that could be considered harmful, offensive, or inappropriate.

#### ENVIRONMENTAL CONSIDERATIONS

The computational resources required to run multiple LLMs and process large volumes of data raise environmental concerns:

- **Energy Consumption:** High energy usage associated with running complex AI models and large-scale data processing.
- **Carbon Footprint:** Environmental impact of the infrastructure required to support the pipeline.

### 3.9. ETHICAL CONSIDERATIONS AND LIMITATIONS

#### 3.9.2 LIMITATIONS

Understanding and acknowledging the limitations of the system is crucial for ethical deployment and user trust:

Limitation	Description
Temporal Limitations	<i>Scope of Knowledge</i> The system's knowledge base and models have a cutoff date, potentially leading to outdated information.
Domain Specificity	Gaps in specialized or niche areas of knowledge may limit performance.
Language Coverage	<i>Language and Cultural Limitations</i> Biases towards languages well-represented in training data, challenges in handling less common languages.
Cultural Context	Limitations in understanding and responding to culturally specific queries or contexts.
Complex Reasoning	<i>Reasoning and Inference Capabilities</i> Challenges in handling queries requiring advanced logical reasoning or domain-specific expertise.
Causal Understanding	Difficulties in inferring causal relationships.
Contextual Nuances	<i>Handling of Ambiguity and Context</i> Difficulties in capturing subtle contextual cues that humans naturally understand.
Disambiguation	Challenges in resolving ambiguities in queries without additional user input.
Static Knowledge Base	<i>Real-time Adaptation</i> Limitations in adapting to real-time changes without system updates.
Learning from Interactions	Inability to learn and improve from individual user interactions due to privacy and architectural constraints.

Table 3.6: Limitations of the pipeline.

# 4

## Empirical Evaluation

### 4.1 DATASET ANALYSIS

This section presents an analysis of the three datasets used in our empirical evaluation: *FactBench*, *YAGO*, and *DBpedia*. Each dataset offers unique characteristics and challenges, providing a comprehensive basis for assessing our knowledge graph fact verification system.

#### 4.1.1 FACTBENCH DATASET

*FactBench* is a multilingual dataset specifically designed for fact-checking in knowledge graphs<sup>1</sup> [10]. It comprises 2,800 facts, 1,500 true and 1,300 false, across three languages: English, German, and French. The dataset covers various domains, including geography, politics, and entertainment. The data was automatically extracted from Wikipedia<sup>2</sup> (DBpedia respectively) and Freebase<sup>3</sup>.

To obtain positive examples, the repository leverages facts from both DBpedia and Freebase. For each property under consideration, they generated these examples by issuing either a SPARQL (for DBpedia) or MQL (for Freebase) query. They then selected the top 150 results. In Freebase, results are ranked using an internal relevance score, while in DBpedia, the results are sorted by

---

<sup>1</sup><https://github.com/DeFacto/FactBench>

<sup>2</sup><https://www.wikipedia.org/>

<sup>3</sup><https://developers.google.com/freebase>

#### 4.1. DATASET ANALYSIS

the number of inbound links to the resource’s corresponding Wikipedia page. In total, 1500 correct statements were collected, with 750 allocated to both the test and training sets, ensuring that each relation had 150 positive facts equally distributed between the test and training sets.

Generating negative examples is more complex than generating positive ones. To ensure that the negative examples closely resemble true statements (i.e., meaningful triples), the team altered the positive examples while still adhering to domain and range restrictions. Given a triple  $(s, p, o)$  and its timespan (from, to) from the knowledge base, they used different methods to generate sets of negative examples. These methods include modifying the subject, object, both subject and object, or the property. Additionally, they included random modifications, a 20% mix of these methods, and variations in the date.

We don’t consider the time aspect in our evaluation, as our system is not designed to handle time-sensitive issues.

Key characteristics of *FactBench* include:

- Multilingual support (English, German, and French)
- Diverse fact types, including domain-specific and temporal facts
- Manually curated for high-quality ground truth

In our analysis, we found that *FactBench* presents a balanced challenge for our system, with a mix of straightforward and complex fact verification tasks.

##### 4.1.2 YAGO DATASET

YAGO (Yet Another Great Ontology) is a large-scale knowledge base derived from Wikipedia, WordNet<sup>4</sup>, and GeoNames<sup>5</sup>. For our evaluation, we use *YAGO2-sample*<sup>6</sup> [26], a subset of the full YAGO2 knowledge graph derived from AMIE horn clauses [9]. This sample consists of 1,386 beliefs spanning 16 unique predicates.

Key characteristics of the YAGO dataset in our evaluation include:

- High accuracy: The gold standard accuracy of the *YAGO2-sample* is 99.20%, indicating a very high-quality dataset.

---

<sup>4</sup><https://wordnet.princeton.edu/>

<sup>5</sup><https://www.geonames.org/>

<sup>6</sup><https://aclanthology.org/attachments/D17-1183.Attachment.zip>

- Diverse predicates: The sample covers 16 different predicates, allowing for evaluation across a range of relationship types.
- Balanced distribution: Unlike domain-specific datasets, *YAGO2-sample* covers a broad range of topics, reflecting the diverse nature of Wikipedia.

The high accuracy of the *YAGO2-sample* presents a unique challenge for our evaluation system.

### 4.1.3 DBPEDIA DATASET

*DBpedia* serves as a comprehensive, large-scale knowledge base derived from Wikipedia, offering structured information about millions of entities. For our evaluation, we utilize *DBpedia* version 2015-10 as provided by Hasibi et al. [12]. This subset serves as the underlying knowledge graph for our experiments. The dataset was carefully curated to ensure a manageable yet representative evaluation set, derived from the vast scale of *DBpedia*.

The dataset was constructed using the following criteria:

- Entity Selection: The dataset is restricted to entities with essential attributes, including a title (`rdfs:label`) and a short abstract (`rdfs:comment`). This ensures that each entity has sufficient information for meaningful evaluation.
- Query Diversity: The dataset leverages the DBpedia-Entity collection [2], a standard test collection for entity retrieval encompassing a wide range of query types. This includes named entity queries, list queries, natural language questions, and keyword queries, providing a comprehensive evaluation landscape.
- Query-Entity Pairing: For each query in the DBpedia-Entity collection, a single relevant entity was selected. This selection process employed a voting mechanism across multiple entity retrieval approaches, ensuring that the chosen entity is both relevant and readily retrievable.
- Fact Extraction: For each selected entity, the corresponding facts (triples) were extracted from DBpedia. These facts form the basis of our evaluation set for tasks such as fact ranking and entity summarization.

Key characteristics of this DBpedia-based evaluation set include:

- Scale: The dataset comprises 100 query-entity pairs and 4,069 corresponding facts, with an average of 41 facts per query-entity pair.
- Diversity: The dataset captures a wide range of relationship types and entity attributes across various domains.

#### 4.1. DATASET ANALYSIS

- Quality Assurance: Manual annotation has been performed to assess the relevance and correctness of facts in relation to their corresponding queries and entities. These annotations serve as the gold standard for our evaluation tasks provided by Marchesin et al. [23].
- Task Suitability: The dataset is specifically designed to support various entity-oriented search tasks, including but not limited to entity summarization, fact ranking, and query-dependent fact selection.

By utilizing this curated subset of DBpedia, we benefit from a balance between the richness of a real-world knowledge graph and the practicality required for thorough empirical evaluation.

#### DATASET SUMMARY

To conclude, our empirical evaluation utilizes three distinct datasets: FactBench, YAGO, and DBpedia. Each dataset offers unique characteristics that allow us to assess our knowledge graph veracity framework across diverse scenarios. Table 2 summarizes the key features of these datasets:

Feature	FactBench	YAGO	DBpedia
Number of Facts	1,860	1,386	4,069
Number of Predicates	18	16	Varied
Gold Accuracy	91.34%	99.20%	N/A
Number of Constraints	130	28	N/A
Query-Entity Pairs	N/A	N/A	100
Avg. Facts per Entity	N/A	N/A	41

Table 4.1: Summary of Datasets Used in Empirical Evaluation

*FactBench* provides a good distribution of true and false statements across multiple domains, offering a robust testbed for fact verification. *YAGO*, with its high accuracy, challenges our framework to detect subtle inaccuracies in an otherwise highly reliable knowledge graph. The *DBpedia* subset, curated specifically for entity-oriented search tasks, allows us to evaluate our framework in the context of query-dependent fact ranking and summarization. This diverse selection of datasets enables a comprehensive evaluation of our veracity estimation framework.

## 4.2 CANDIDATE MODELS

### 4.2.1 GEMMA2

Gemma is a family of lightweight, state-of-the-art open models from Google, built from the same research and technology used to create the Gemini models<sup>7</sup>. They are text-to-text, decoder-only large language models, available in English, with open weights for both pre-trained variants and instruction-tuned variants. Gemma 2 implements a novel approach to attention mechanisms:

- Every other layer uses a sliding window attention with a local context of 4096 tokens.
- Alternating layers employ full quadratic global attention across the entire 8192 token context.

This hybrid approach aims to balance efficiency with the ability to capture long-range dependencies in the input.

We selected the *Gemma2-9B* model for our evaluation, which has 9 billion parameters. The 9B model learns from a larger teacher model during initial training in pre-training and use on-policy distillation to refine its performance post-training. This approach allows Gemma2-9B to capture the knowledge and capabilities of the larger model while maintaining a more compact size. As a result, *Gemma2-9B* delivers competitive performance relative to models 2-3 times its size, making it an attractive choice for applications with computational constraints.

### 4.2.2 QWEN2.5

*Qwen2.5* is the latest series of Qwen LLMs [39]. For *Qwen2.5*, Alibaba Cloud<sup>8</sup> release a number of base language models and instruction-tuned language models ranging from 0.5 to 72 billion parameters. All models are pre-trained on our latest large-scale dataset, encompassing up to 18 trillion tokens. Compared to Qwen2, Qwen2.5 has acquired significantly more knowledge and has greatly improved capabilities in coding and mathematics. Additionally, the new models achieve significant improvements in instruction following, generating long

---

<sup>7</sup><https://deepmind.google/technologies/gemini/#introduction>

<sup>8</sup>[https://www.alibabacloud.com/en?\\_p\\_lc=7](https://www.alibabacloud.com/en?_p_lc=7)

## 4.2. CANDIDATE MODELS

texts (over 8K tokens), understanding structured data (e.g, tables), and generating structured outputs especially JSON. *Qwen2.5* models are generally more resilient to the diversity of system prompts, enhancing role-play implementation and condition-setting for chatbots. Like *Qwen2*, the *Qwen2.5* language models support up to 128K tokens and can generate up to 8K tokens. They also maintain multilingual support for over 29 languages [35].

We selected the *Qwen2.5-7b* model for our evaluation, which has 7 billion parameters.

### 4.2.3 LLAMA3.1

The Meta *Llama3.1* collection of multilingual LLMs is a collection of pre-trained and instruction tuned generative models in 8B, 70B and 405B sizes (text in/text out). The *Llama3.1* instruction tuned text only models (8B, 70B, 405B) are optimized for multilingual dialogue use cases and outperform many of the available open source and closed chat models on common industry benchmarks.

*Llama3.1* is an auto-regressive language model that uses an optimized transformer architecture. The tuned versions use Supervised Fine-Tuning (SFT) and Reinforcement Learning with Human Feedback (RLHF) to align with human preferences for helpfulness and safety. All model versions use Grouped-Query Attention (GQA) for improved inference scalability. *Llama3.1* was pre-trained on 15 trillion tokens of data from publicly available sources. The fine-tuning data includes publicly available instruction datasets, as well as over 25M synthetically generated examples [8, 1].

We selected the *Llama3.1-8b* model for our evaluation, which has 8 billion parameters.

### 4.2.4 MISTRAL

The *Mistral* model, released by Mistral AI <sup>9</sup>, is a high-performance LLM, designed to outperform larger models in efficiency and effectiveness. With innovations such as GQA and Sliding Window Attention (SWA), Mistral offers faster inference and better handling of long sequences, reducing computation costs while maintaining high performance [15, 24].

---

<sup>9</sup><https://mistral.ai/>

We selected the *Mistral-7b* model for our evaluation, which has 7.3 billion parameters, its structure allows it to be both cost-effective and memory efficient, making it suitable for a wide variety of real-world applications.

## 4.3 EXPERIMENTAL SETUP

### 4.3.1 DATASETS

### 4.3.2 PERFORMANCE METRICS AND EVALUATION

Performance metrics are essential in assessing the efficacy, efficiency, and reliability of a system or model. The selection of metrics mostly depends on the characteristics of the task, the data, and the objectives. This section emphasizes the principal performance metrics typically employed in systems utilizing LLMs, information retrieval, and various machine learning tasks.

#### CORRECT AND INCORRECT CRITERIA

The system incorporates explicit CORRECT and INCORRECT states, indicating a binary evaluation mechanism for overall performance. This fundamental assessment provides a clear, high-level indication of the system's success in handling queries.

#### RELEVANCE AND ACCURACY METRICS

The evaluation of a fact-checking system typically involves assessing both the correctness and relevance of responses.

Potential metrics include:

- **Recall:** The proportion of relevant responses generated by the system among all possible relevant responses.
- **Precision:** The proportion of correct responses among all generated responses.
- **F1 Score:** The harmonic mean of precision and recall, providing a balanced measure of accuracy.
- **Accuracy:** The proportion of correct responses generated by the system.

### 4.3. EXPERIMENTAL SETUP

#### LATENCY AND EFFICIENCY MEASURES

Given the complexity of the pipeline, evaluating its operational efficiency is crucial:

- **Response Time:** Measuring the end-to-end time from query input to response generation.
- **Component-wise Latency:** Assessing the processing time of individual pipeline components (e.g., embedding generation, LLM processing).
- **Resource Utilization:** Monitoring computational resource usage, particularly important given the use of multiple LLMs.
- **Cost Efficiency:** Evaluating the cost-effectiveness of the pipeline in terms of computational resources and infrastructure.

#### CONSISTENCY EVALUATION

The use of multiple models and a conflict resolution mechanism necessitates specific evaluation of output consistency:

- **Stability Across Models:** Assessing the consistency of responses generated by different LLMs for the same query, refer to algorithm 2.

---

**Algorithm 2** Stability Across Queries

---

```
1: procedure STABILITYACROSSQUERIES(queryResponses)
2:   stabilityScores  $\leftarrow$  empty list
3:   for each responses in queryResponses do
4:     uniqueResponses  $\leftarrow$  set of unique elements in responses
5:     stability  $\leftarrow$  size of uniqueResponses
6:     stabilityScore  $\leftarrow$   $(1 - \frac{\text{stability}-1}{|\text{responses}|}) \times 100$ 
7:     Append stabilityScore to stabilityScores
8:     if stabilityScores is not empty then
9:       return mean of stabilityScores
10:    else
11:      return 0
12:    end if
13:  end for
14: end procedure
```

---

### 4.3.3 SYSTEM CONFIGURATIONS

The system configurations are selected based on the best results obtained from black-box testing the pipeline through a series of experiments, detailed in section 5. Table 4.2 summarizes the key system configurations used in our empirical evaluation.

Section	Method/Model Used	Considerations
Human Understandable Text	Gemma2.9b	Other LLMs can be used, but using instruction-tuned models is recommended. This is skipped for <i>FactBench</i> dataset as discussed on 3.2.1.
Question Generation	Gemma2.9b	Other LLMs can be used, but using instruction-tuned models is recommended.
Question Relevance	Jina-reranker-v1-turbo-en	Cross-encoder models are recommended for this task.
Question Relevance Threshold	0.5	-
Num of Selected Questions	3	-
Google Search	-	Used query params: <code>lr = 'lang_en'</code> , <code>gl = 'us'</code> , <code>hl = 'en'</code> , <code>num = '100'</code> . The lr parameter is set to the language of the query, gl to the country, hl to the language, and num to the number of results.
Num of Selected Documents	10	-

### 4.3. EXPERIMENTAL SETUP

Document Selection	ms-marco-MiniLM-L-6-v2	Filtered out the documents from these origins: "dbpedia.org", "wikipedia.org", "wikimedia.org", "wikidata.org", "quora.com", "britannica.com", "scholarpedia.org", "newworldencyclopedia.org", "everipedia.org", "encyclopedia.com", "wikibooks.org", "wiktionary.org", "wikiversity.org", "wikisource.org", "wikiquote.org", "wikivoyage.org", "academia.edu", "nytimes.com"
Embedding Model	BAAI/bge-small-en-v1.5	-
Chunking Strategy	Sliding Window with window size 3	-
Similarity Cut-off	Simple	Use the threshold to filter out irrelevant documents.
Similarity Cut-off Threshold	0.3	-
Top_k	6	-
Tie-Breaking	-	Based on the models we selected, we use model with higher-param for each model, for Gemma2.9b → 27b, Qwen2.5:7b → 14b, Llama3.1:8b → 70b, and Mistral:7b → Mistral nemo:12b.

Table 4.2: System Configurations for Empirical Evaluation

The tests are run on a server with the following specifications:

- **Model Name:** Mac Studio
- **Model Identifier:** Mac14,14
- **Model Number:** Z180000M3T/A
- **Chip:** Apple M2 Ultra
- **Total Number of Cores:** 24 (16 performance and 8 efficiency)

- **Memory:** 192 GB
- **System Firmware Version:** 11881.1.1
- **OS Loader Version:** 11881.1.1

The timing tests are conducted over 50 samples for each section, using a MacBook Pro equipped with an Apple M2 Max chip and 32 GB of memory.

## 4.4 COMPARATIVE ANALYSIS

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, Sliding Window (ws 3), Similarity Cut-off (Original), Top\_k 6

Dataset	Model	Stability	Avg. Request Duration*	Avg. tokens per request	ACC	F1
FactBench	Gemma2	0.8720	5.6826s	1605.29	<b>0.9032</b>	<b>0.9101</b>
	Qwen2.5	0.8685	6.3094s	1652.66	0.8729	0.8888
	LLama3.1	0.8291	6.5270s	1679.04	0.8154	0.8299
	Mistral	0.8650	4.6692s	1594.76	0.8511	0.8733
	Proposed (At_Most)	<b>0.9176</b>	16.815s	1604.196	0.8964	0.9071
	Proposed (At_Least)	N/A	N/A	N/A	N/A	N/A
YAGO	Gemma2	0.8785	5.1982s	1597.50	0.8499	0.9187
	Qwen2.5	0.8773	6.2770s	1653.56	0.8506	0.9191
	LLama3.1	0.8276	6.4070s	1682.78	0.8059	0.8922
	Mistral	0.8817	4.4824s	1589.58	<b>0.9250</b>	<b>0.9610</b>
	Proposed (At_Most)	<b>0.9226</b>	6.4810s	1586.02	0.8737	0.9325
	Proposed (At_Least)	N/A	N/A	N/A	N/A	N/A
DBpedia	Gemma2	N/A	N/A	N/A	N/A	N/A
	Qwen2.5	N/A	N/A	N/A	N/A	N/A
	LLama3.1	N/A	N/A	N/A	N/A	N/A
	Mistral	N/A	N/A	N/A	N/A	N/A
	Proposed (At_Most)	N/A	N/A	N/A	N/A	N/A
	Proposed (At_Least)	N/A	N/A	N/A	N/A	N/A

\* Each query uses two requests, so the average request duration and average tokens per request are calculated based on the two requests.

Table 4.3: Empirical Evaluation Results of the Proposed System and Candidate Models over the Datasets

### 4.4.1 DISCUSSION OF RESULTS

The empirical evaluation results presented in Table 4.3 offer valuable insights into the performance of our proposed system and the individual candidate models across the three datasets: *FactBench*, *YAGO*, and *DBpedia*.

Among the candidate models, *LLama3.1* generates the reasoning most of the time and don't follow the instructions fully, this way the Avg completion tokens is higher than the other models.

#### 4.4. COMPARATIVE ANALYSIS

In terms of overall accuracy and F1 scores, the *Gemma2* model outperforms the other individual models on the FactBench dataset. This superior performance can be attributed to *Gemma2*'s advanced attention mechanisms, which employ a hybrid approach of sliding window attention and global attention to effectively capture both local and long-range dependencies in the input sequences. The model's ability to learn from larger teacher models during pre-training and its use of on-policy distillation post-training likely contribute to its strong performance relative to its size.

Interestingly, the *Gemma2* model achieves the highest accuracy and F1 scores on the YAGO dataset, surpassing even *Gemma2*. This may be due to *Mistral*'s optimized transformer architecture, which incorporates innovations such as Grouped-Query Attention and Sliding Window Attention to improve inference speed and handle longer sequences more effectively.

The *Qwen2.5* and *LLama3.1* models generally fall in the middle of the pack in terms of accuracy and F1 scores. While both models leverage large-scale pre-training and instruction tuning, their performance suggests room for improvement in the specific task of knowledge graph fact verification.

Notably, our proposed system, which employs an ensemble approach using the combination strategy, consistently achieves competitive accuracy and F1 scores compared to the best-performing individual models. This highlights the effectiveness of combining the strengths of multiple models through majority voting and adaptive dispute resolution techniques. However, the increased computational cost of running multiple models in parallel is reflected in the higher average request duration for the proposed system. The stability scores provide another dimension for comparison, indicating the consistency of model outputs across different runs. The proposed system achieves the highest stability scores on datasets, showing the robustness gained by synthesizing results from multiple models.

Among the individual models, *Gemma2* and *Mistral* exhibit relatively high stability, while *LLama3.1* has the lowest stability scores.

##### 4.4.2 KEY FINDINGS

##### 4.4.3 ERROR ANALYSIS

# 5

## Ablation Study

Our proposed framework for knowledge graph fact verification utilizes a unique combination of web search and language model processing. However, to ensure the robustness and effectiveness of our approach, it is crucial to compare our methods with state-of-the-art RAG techniques, particularly in the critical areas of chunking, embedding, and retrieval.

This section aims to provide a comprehensive comparison between our approach and the RAG-based methods. We will focus on three key components of our framework: 1) the retrieval mechanisms utilized to fetch relevant information, 2) the chunking strategies used to segment information, 3) the embedding models employed for representation, and 4) different hyper parameters and configurations. By analyzing these components in light of RAG recommendations, we aim to identify potential areas for improvement and validate the strengths of our current approach.

Through this comparison, we seek to situate our work within the broader context of retrieval-augmented fact verification systems and provide insights into the trade-offs and benefits of our methodological choices. This analysis will not only contribute to the refinement of our framework but also offer valuable perspectives on the application of RAG principles to knowledge graph fact verification tasks.

## 5.1. EVALUATION METHODOLOGY

### 5.1 EVALUATION METHODOLOGY

This study employs a systematic approach to evaluate and optimize various components of our framework, with the ultimate goal of determining the best methods for each section. Our methodology is designed to isolate and assess the impact of different techniques and parameters on overall system performance, while also considering the efficacy of sampling methods compared to full data runs.

#### 5.1.1 ITERATIVE OPTIMIZATION PROCESS

The evaluation process follows an iterative strategy, focusing on specific sections of the framework in each iteration:

1. **Section Isolation:** In each iteration, we isolate a particular section of the framework for investigation, keeping other components constant. This "enclosed box" approach allows for a controlled examination of individual elements.
2. **Parameter Variation:** Within the isolated section, we systematically vary relevant parameters or methods. This includes, but is not limited to, testing different sampling methods against full data runs.
3. **Performance Evaluation:** For each configuration, we assess the system's performance using predefined metrics (detailed in Sections 5.1.2 and 4.3.2).
4. **Best Method Selection:** Based on the evaluation results, we identify the best-performing method or configuration for the section under investigation.
5. **Incremental Optimization:** The optimal configuration from each iteration is incorporated into the framework for subsequent iterations, gradually refining the entire system.

#### 5.1.2 EVALUATION METRICS

The performance of each configuration is assessed using the following metrics, it's same as the metrics in the empirical evaluation section 4.3.2:

- **Accuracy (Acc):** Measures the overall correctness of predictions:

$$Acc = \frac{TP + TN}{TP + TN + FP + FN} \quad (5.1)$$

where TP, TN, FP, and FN are True Positives, True Negatives, False Positives, and False Negatives, respectively.

- **F1 Score:** Provides a balanced measure of precision and recall:

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (5.2)$$

- **Average Latency:** Measured in seconds per query to assess computational efficiency:

$$\text{Avg Latency} = \frac{\text{Total Processing Time}}{\text{Number of Queries}} \quad (5.3)$$

### 5.1.3 SIGNIFICANCE OF THE METHODOLOGY

This methodical approach serves several key purposes:

1. **Optimization of Individual Components:** By isolating sections, we can fine-tune each part of the framework independently.
2. **Holistic System Improvement:** The iterative process ensures that optimizations in one section complement the overall system performance.
3. **Efficiency-Accuracy Trade-off Analysis:** Comparing sampling methods to full data runs helps balance computational efficiency with result accuracy.
4. **Scalability Assessment:** This approach informs decisions on system scalability as data volumes increase.

By employing this rigorous evaluation methodology, we aim to identify the best methods for each section of our framework, potentially enabling more efficient and accurate data processing. The inclusion of sampling method comparisons adds an extra dimension to our optimization efforts, potentially offering insights into cost-effective alternatives to full data processing where applicable.

## 5.2 DOCUMENT SELECTION

We explore various techniques for retrieving relevant documents from search engine results, with a specific focus on Google search engine. The goal is to identify the most effective methods for finding documents that perfectly match the information need expressed in the query. We consider both unsupervised and supervised approaches.

### 5.2.1 UNSUPERVISED METHODS

#### BM25

BM25 [28] is a widely used unsupervised retrieval method that relies on term frequency and inverse document frequency (TF-IDF) weighting. It estimates the relevance of documents to a query based on the frequency of query terms in each document, offset by the rarity of those terms across the full document collection. BM25 has proven to be a robust baseline for many retrieval tasks. However, it relies on lexical matching between query and document terms, which can limit its effectiveness for queries and documents that use different vocabulary to express similar concepts.

$$\text{BM25}(D, Q) = \sum_{i=1}^n \text{IDF}(q_i) \cdot \frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot (1 - b + b \cdot \frac{|D|}{\text{avgdl}})}$$

where:

$D$  : document

$Q$  : query containing keywords  $q_1, \dots, q_n$

$f(q_i, D)$  : frequency of  $q_i$  in  $D$

$|D|$  : length of document  $D$

$\text{avgdl}$  : average document length in the corpus

$k_1, b$  : free parameters

$$\text{IDF}(q_i) = \log \frac{N - n(q_i) + 0.5}{n(q_i) + 0.5}$$

where:

$N$  : total number of documents in the corpus

$n(q_i)$  : number of documents containing  $q_i$

#### CONTRIEVER

Contriever is a more recently proposed unsupervised method by Izacard et al.[14] that leverages contrastive learning to train dense retrieval models. Rather than relying on term matching, Contriever learns to map semantically similar text pairs to nearby embeddings in a continuous vector space. At query time,

Contriever embeds the query and retrieves the documents whose embeddings are nearest to the query under cosine similarity. By operating in this learned semantic space, Contriever can potentially identify relevant documents that use different surface forms than the query. Contriever has shown promising results, outperforming BM25 on a range of benchmarks when large unsupervised pretraining datasets are available. However, details on its performance in this specific multi-query retrieval setup are needed to fully assess its capabilities here.

While Contriever can be used as an unsupervised retriever, for our thesis project focusing on search-related data, we opt to use the MS-MARCO fine-tuned version.<sup>1</sup> Here's why:

- **Relevance to Search Tasks:** MS-MARCO (Microsoft Machine Reading Comprehension) is a large-scale dataset specifically designed for search and question-answering tasks. It contains real queries from Bing search engine and human-annotated relevant passages. By fine-tuning Contriever on MS-MARCO, the model becomes particularly adept at understanding and representing search-like queries and documents.
- **Improved Performance:** Fine-tuning on MS-MARCO significantly boosts Contriever's performance on various retrieval benchmarks, especially those related to web search and question answering. This improvement is crucial for our project, which deals with search-term related data.
- **Domain Adaptation:** Although Contriever's unsupervised training on Wikipedia and CCNet provides a strong foundation, fine-tuning on MS-MARCO helps adapt the model to the specific nuances and patterns present in search queries and web documents. This domain adaptation is valuable for our search-centric application.

### 5.2.2 SUPERVISED METHODS

#### JINA.AI RERANKER

The Jina.ai Reranker is a supervised neural ranking model. Jina Reranker employs a cross-encoder architecture, which represents a paradigm shift from traditional bi-encoder models used in embedding-based search. While bi-encoder models separately encode queries and documents, cross-encoders jointly process query-document pairs, allowing for more nuanced semantic understanding

---

<sup>1</sup><https://huggingface.co/facebook/contriever-msmarco>

## 5.2. DOCUMENT SELECTION

and relevance assessment. The model generates a relevance score for each query-document pair, enabling a more precise ranking of search results. This approach addresses limitations of vector similarity-based methods by capturing complex token-level interactions between queries and documents.

For our project, we use *jina-reranker-v2-base-multilingual*<sup>2</sup>. This model has demonstrated exceptional performance across various benchmarks and practical applications. In multilingual tasks, it achieved state-of-the-art recall@10 scores on the MKQA dataset [22] spanning 26 languages, while also exhibiting superior NDCG@10 scores on English-language tasks in the BEIR benchmark [36]. Notably, it secured the top position on the AirBench leaderboard upon its release<sup>3</sup>.

These capabilities make the model particularly valuable for multilingual information retrieval, agentic Retrieval-Augmented Generation (RAG) systems, and even in programming and software development support.

## MS MARCO MINILM

The MS MARCO MiniLM is another supervised neural model, based on the popular BERT architecture but distilled to a smaller size for efficiency. For our

Model Name	NDCG@10 (TREC DL 19)	MRR@10 (MS Marco Dev)	Docs / Sec
ms-marco-TinyBERT-L-2-v2	69.84	32.56	9000
ms-marco-MiniLM-L-2-v2	71.01	34.85	4100
ms-marco-MiniLM-L-4-v2	73.04	37.70	2500
ms-marco-MiniLM-L-6-v2	74.30	39.01	1800
ms-marco-MiniLM-L-12-v2	74.31	39.02	960

Table 5.1: Performance of Pre-trained Cross-Encoders

project, we use *ms-marco-MiniLM-L-6-v2*<sup>4</sup> Cross-Encoder model. This model, trained on the extensive MS MARCO dataset comprising approximately 500,000 authentic search queries from the Bing search engine [27], demonstrates superior performance within a two-stage Retrieve & Re-rank framework. In this paradigm, an initial retrieval phase employs either lexical search methods or

<sup>2</sup><https://huggingface.co/jinaai/jina-reranker-v2-base-multilingual>

<sup>3</sup><https://huggingface.co/spaces/AIR-Bench/leaderboard>

<sup>4</sup><https://huggingface.co/cross-encoder/ms-marco-MiniLM-L-6-v2>

dense retrieval techniques utilizing a bi-encoder to identify a broad set of potentially relevant documents. Subsequently, the Cross-Encoder refines this candidate set through a simultaneous processing of the query and each retrieved document, generating a relevance score on a scale of 0 to 1.

### 5.2.3 EVALUATION WITH LARGE LANGUAGE MODELS

With checking the similarity between the retrieved documents and the query in different methods, based on figure 5.1 we can figured out that Bm25-Okapi stands out as the most distinct model, with low similarity scores (0.20-0.24) to the others, suggesting it employs fundamentally different retrieval mechanisms. In contrast, the neural models show higher inter-model similarities (0.43-0.44), indicating shared approaches or architectures. The strong relationship (0.43) between contriever-msmarco and re-ranker-msmarco, likely due to shared training data or similar optimizations. The two re-ranker models gain the highest similarity (0.44) but suggest distinctiveness in their retrieval mechanisms and architectures. In general, we can find out with different methods we have different result over the same query.



Figure 5.1: Document Retrieval Confusion Matrix

To assess the quality of the retrieved documents from each of the above methods, we are passing them through one of our models and evaluating the outputs.

The empirical results indicate that the model *ms-marco-MiniLM-L-6-v2* achieved the highest F1 score, thus demonstrating superior performance among the evaluated models. However, it is noteworthy that the performance metrics across all

### 5.3. EMBEDDING MODELS

Method	Acc	F1	Latency
<i>Unsupervised</i>			
Bm25	0.8882	0.8940	N/A
contriever-msmarco	0.8932	0.8988	N/A
<i>supervised</i>			
jina-reranker-v2-base-multilingual	0.9004	0.9065	N/A
ms-marco-MiniLM-L-6-v2	<b>0.9014</b>	<b>0.9077</b>	N/A

Table 5.2: Evaluation Results for Different Methods through the Pipeline (just with the Gemma2 model)

models were closely clustered, suggesting that even traditional methodologies applied within our pipeline yield satisfactory outcomes.

It is crucial to emphasize the significance of data quality in this context, as it substantially influences the efficacy of the results. To validate the factual accuracy of the knowledge graph, we employed a multi-query information fetching through web search engines for each fact. This approach provides a reasonable degree of verification for the facts contained within the knowledge graph.

For subsequent evaluations and analyzes, we will designate the model *ms-marco-MiniLM-L-6-v2* as our baseline for retrieval tasks. This decision is predicated on its superior F1 score relative to the other models under consideration.

## 5.3 EMBEDDING MODELS

Text embeddings are dense vector representations that capture the semantic meaning and relationships between words, sentences, or documents in a low-dimensional space. By mapping text to a continuous vector space, embeddings enable efficient similarity computations and have become a fundamental building block for many NLP applications, such as information retrieval, text classification, clustering, and semantic search. This section provides an in-depth analysis and comparison of five state-of-the-art text embedding models:

- Alibaba-NLP/gte-large-en-v1.5
- jinaai/jina-embeddings-v3
- dunzhang/stella\_en\_1.5B\_v5
- Nextcloud-AI/multilingual-e5-large-instruct
- BAAI/bge-small-en-v1.5

These models leverage recent advancements in transformer architectures, contrastive learning, and instruction fine-tuning to produce high-quality, general-purpose embeddings that excel across a wide range of downstream tasks. We examine their model architectures, training methodologies, supported features, and empirical performance on standard benchmarks. Through this comparative study, we aim to provide insights and guidance for practitioners to select the most suitable embedding model based on their specific use case and computational constraints.

### 5.3.1 GTE-LARGE-EN-v1.5

The Alibaba-NLP model *gte-large-en-v1.5* is text embedding model designed for general text representation and retrieval tasks. It is built upon a Transformer++ encoder architecture, combining the strengths of BERT [6] with advanced techniques such as rotary position embeddings (RoPE) [32] and Gated Linear Units (GLU). This combination allows for highly efficient text encoding over long sequences, with a maximum context length of 8192 tokens, significantly surpassing previous models restricted to shorter context lengths (up to 512 tokens) [41].

One of the major improvements in the gte-v1.5 series is its ability to process long-context text inputs, making it ideal for complex text retrieval and re-ranking tasks, especially in multilingual settings [19]. This series of models has demonstrated superior performance in multiple benchmarks, including the Massive Text Embedding Benchmark (MTEB) [25] and the LoCo long-context retrieval benchmark [29]. In particular, the model ranked second on the MTEB leaderboard and first in the Chinese version of MTEB (C-MTEB).

The model achieves these results by employing a hybrid architecture, including both a text representation model (TRM) and a cross-encoder reranker. The TRM generates dense text embeddings for retrieval tasks, while the reranker refines results through more precise scoring of candidate texts. This architecture is optimized for efficiency, allowing faster inference while maintaining high accuracy during both pretraining and fine-tuning stages.

The *gte-large-en-v1.5* also includes instruction-tuned variants, such as *gte-Qwen1.5-7B-instruct*, which is particularly effective for multilingual text embeddings, leveraging a wide range of unsupervised and supervised contrastive learning techniques. These instruction-tuned models have outperformed vari-

### 5.3. EMBEDDING MODELS

ous other large embedding models, making them highly suitable for industrial applications that require efficient, accurate text representation across diverse languages.

In summary, the *gte-large-en-v1.5* model stands out in its category due to its ability to handle large context lengths, its efficient encoding techniques, and its strong performance across multilingual and long-context benchmarks. This makes it an invaluable tool for a variety of text retrieval, classification, and representation tasks in both academic research and real-world applications.

#### 5.3.2 JINA-EMBEDDINGS-V3

The *Jina-embeddings-v3* model is a cutting-edge multilingual text embedding solution, developed by Jina AI, aimed at addressing a wide range of NLP tasks. Based on the *Jina-XLM-RoBERTa* architecture, this model supports long-context inputs, handling sequences of up to 8192 tokens thanks to its integration of RoPE [32, 31].

This ability to process extended sequences makes the model well-suited for tasks such as text retrieval, clustering, classification, and text matching across multiple languages. One of the key innovations of *Jina-embeddings-v3* is the introduction of task-specific Low-Rank Adaptation (LoRA) [13] adapters. These adapters are used to tailor the model’s embeddings to specific tasks, such as query-document retrieval, clustering, re-ranking, and classification. This task-specific optimization is achieved without significantly increasing the model’s parameter size.

The model excels in multilingual environments, supporting over 30 languages, and is optimized for performance in long-context retrieval tasks. Compared to LLMs like *e5-mistral-7b-instruct*, *jina-embeddings-v3* offers a more efficient solution with fewer parameters (570 million vs. 7.1 billion), while still achieving competitive or superior performance on several benchmarks. For example, it surpasses proprietary models like OpenAI<sup>5</sup> and Cohere<sup>6</sup> on English tasks and achieves high scores on multilingual benchmarks.

*Jina-embeddings-v3* also features flexible Matryoshka Representation Learning (MRL) [17], allowing users to reduce the embedding size from 1024 to as low

---

<sup>5</sup><https://openai.com/>

<sup>6</sup><https://cohere.com/>

as 32 dimensions, making it adaptable to different resource constraints without significant loss of performance.

### 5.3.3 STELLA\_EN\_1.5B\_v5

The Dunzhang *Stella\_en\_1.5B\_v5*<sup>7</sup> is a powerful multilingual text embedding model, built upon the foundations of *Alibaba-NLP/gte-large-en-v1.5* 5.3.1 and *gte-Qwen2-1.5B-instruct*. This model supports two main prompts for diverse tasks: "s2p" (sentence-to-passage) for information retrieval, and "s2s" (sentence-to-sentence) for semantic textual similarity. These prompts simplify its application in NLP tasks, such as retrieving relevant passages or finding semantically similar text based on a given query.

One of the standout features of *Stella\_en\_1.5B\_v5* is its implementation of MRL [17], allowing the model to output embeddings in multiple dimensions ranging from 512 to 8192, depending on user needs. Typically, a 1024-dimensional output offers an optimal balance between performance and efficiency. In benchmark tests, the model achieves highly competitive results, with only a minor performance difference between 1024-dimensional and 8192-dimensional embeddings. The model can be employed using both SentenceTransformers and transformers libraries, supporting flexible input formats. It is trained on shorter sequences (up to 512 tokens), making it most effective for short-to-medium-length text tasks.

### 5.3.4 MULTILINGUAL-E5-LARGE-INSTRUCT

The Multilingual E5-Large-Instruct model is an advanced multilingual text embedding model introduced as part of the E5 model family, which aims to improve the quality and utility of multilingual text embeddings [37]. It is specifically designed to support a wide range of languages and to deliver robust performance across various tasks such as text retrieval, semantic similarity, and multilingual retrieval.

The E5-Large-Instruct model contains 24 layers and features an embedding size of 1024. It builds on the *XLM-RoBERTa-large* [5] architecture, which supports 100 languages, albeit with varying performance depending on the resource

---

<sup>7</sup>[https://huggingface.co/dunzhang/stella\\_en\\_1.5B\\_v5](https://huggingface.co/dunzhang/stella_en_1.5B_v5)

### 5.3. EMBEDDING MODELS

richness of the language in question. The model was initialized from XLM-RoBERTa-large and underwent two key stages of training:

- **Contrastive Pre-training:** The model was pre-trained on approximately 1 billion weakly supervised multilingual text pairs using a contrastive loss function with a temperature of 0.01. This stage aimed to align similar texts in a shared embedding space.
- **Fine-tuning:** Following pre-training, the model was fine-tuned using high-quality labeled datasets from the E5-mistral paper. This second stage involved a more supervised approach, optimizing performance across specific tasks. During this phase, instruction-tuning was incorporated, where the model learned to generate better embeddings by using natural language task instructions.

The E5-Large-Instruct model was evaluated on BEIR and MTEB benchmarks, and its performance is on par with state-of-the-art English-only models. Evaluation on the MIRACL [42] multilingual retrieval benchmark across 16 languages and on Bitext mining tasks across over 100 languages demonstrated its capability to handle diverse languages effectively. Despite the excellent performance on high-resource languages, the model shows degradation in performance for low-resource languages, a common limitation of multilingual models. The use of contrastive learning and instruction tuning enables the model to generate highly effective embeddings for information retrieval tasks. However, the absolute cosine similarity values generated by the model are less critical than the relative order of these values, which is more relevant for retrieval and similarity ranking tasks.

#### 5.3.5 BGE-SMALL-EN-V1.5

The *bge-small-en-v1.5* model is part of the BGE (BAAI General Embeddings) series developed by the Beijing Academy of Artificial Intelligence [38]. It is a compact English-specific model with just 25M parameters, making it highly efficient for deployment in resource-constrained environments. The model architecture is based on RoBERTa [21] with optimizations like dynamic token pruning and embedding factorization to reduce computational costs.

*bge-small-en-v1.5* follows a two-stage training pipeline similar to other BGE models:

- **Pre-training:** Weakly-supervised contrastive pre-training on large-scale web data

- **Fine-tuning:** Supervised fine-tuning on a curated set of high-quality English NLP datasets

The pre-training stage leverages diverse data sources like Wikipedia, Reddit, and CommonCrawl to learn general-purpose text representations. The fine-tuning stage incorporates datasets spanning retrieval, classification, paraphrase detection, and semantic textual similarity tasks to instill task-specific knowledge.

Despite its small size, *bge-small-en-v1.5* punches above its weight on several English benchmarks. On the SentEval suite [4], it outperforms the base-sized BERT and RoBERTa models on most tasks while being 4x more compact. On retrieval challenges like BEIR, it achieves competitive results, often surpassing larger models like MPNet and the original DPR. The model’s strong performance can be attributed to the efficient architecture design and the use of high-quality fine-tuning data. It presents an attractive option for applications requiring low-latency inference or deployment on edge devices.

### 5.3.6 COMPARATIVE ANALYSIS

We examine their model size and efficiency, language coverage, supported features, and overall performance to provide insights for selecting the most suitable model based on specific requirements.

#### MODEL SIZE AND EFFICIENCY

Table 5.3 compares the model size, memory usage, embedding dimensions, and maximum token length of the five models. The *stella\_en\_1.5B\_v5* model has the largest size with 1,543 million parameters, while *bge-small-en-v1.5* is the smallest with only 33 million parameters. Larger models generally require more memory and computational resources, which may be a consideration for resource-constrained environments. In terms of memory usage, *stella\_en\_1.5B\_v5* requires 5.75 GB in fp32 precision, while *bge-small-en-v1.5* only needs 0.12 GB. This substantial difference in memory footprint can be a decisive factor when deploying models on edge devices or serving them in real-time applications with limited resources. The embedding dimensions also vary among the models, ranging from 384 for *bge-small-en-v1.5* to 8192 for *stella\_en\_1.5B\_v5*. Higher-dimensional embeddings can capture more fine-grained semantic information but may increase storage requirements and similarity computation

### 5.3. EMBEDDING MODELS

costs. Practitioners should consider the trade-off between embedding quality and efficiency based on their specific use case.

Model	Model Size (Million Parameters)	Memory Usage (GB, fp32)	Embedding Dimensions	Max Tokens
stella_en_1.5B_v5	1543	5.75	8192	131072
jina-embeddings-v3	572	2.13	1024	8194
gte-large-en-v1.5	434	1.62	1024	8192
multilingual-e5-large-instruct	560	2.09	1024	514
bge-small-en-v1.5	33	0.12	384	51262

Table 5.3: Comparison of Embedding Models

### LANGUAGE COVERAGE

Language coverage is a crucial aspect when selecting an embedding model for multilingual applications. The *Multilingual-e5-large-instruct* model stands out in this regard, as it supports a wide range of languages. This model leverages instruction fine-tuning on multilingual data, enabling it to generate high-quality embeddings for various languages. The *Jina-embeddings-v3* model also offers multilingual support, although the exact language coverage is not specified in the provided context. On the other hand, the *Bge-small-en-v1.5*, *Stella\_en\_1.5B\_v5*, and *Gte-large-en-v1.5* models primarily focus on English embeddings, making them more suitable for monolingual English applications.

### SUPPORTED FEATURES

The embedding models offer various features that cater to different use cases. The *Stella\_en\_1.5B\_v5* model supports two main prompts for diverse tasks: "s2p" (sentence-to-passage) for information retrieval and "s2s" (sentence-to-sentence) for semantic textual similarity. This flexibility allows users to adapt the model for specific NLP tasks with minimal effort. The *Jina-embeddings-v3* model incorporates task-specific LoRA adapters, enabling it to generate high-quality embeddings for query-document retrieval, clustering, classification, and text matching. This feature eliminates the need for manual prompt engineering and provides a more streamlined approach for task-specific embeddings. Additionally, the *Stella\_en\_1.5B\_v5* and *Jina-embeddings-v3* models implement MRL, allowing users to generate embeddings with different dimensions based on their

needs. This adaptability is particularly useful when balancing between embedding quality and storage or computational efficiency.

## CONCLUSION

The choice of text embedding model depends on various factors, including the specific application, language coverage requirements, available computational resources, and desired features. For monolingual English applications, the *Gte-large-en-v1.5* and *Stella\_en\_1.5B\_v5* models offer high-quality embeddings with support for longer input sequences. The *Stella\_en\_1.5B\_v5* model, in particular, provides prompt-based adaptability for information retrieval and semantic similarity tasks. For multilingual applications, the *Multilingual-e5-large-instruct* and *Jina-embeddings-v3* models are strong contenders. The *Multilingual-e5-large-instruct* model supports a wide range of languages, while *Jina-embeddings-v3* offers task-specific LoRA adapters for enhanced performance across various NLP tasks. When computational resources are limited, the *Bge-small-en-v1.5* model presents a lightweight option with competitive performance. Its small size and low memory footprint make it suitable for deployment on edge devices or real-time applications. Ultimately, practitioners should carefully evaluate their specific requirements and constraints before selecting an embedding model. The comparative analysis provided in this section aims to assist in this decision-making process by highlighting the key differences and strengths of each model.

Now we will test these models through the pipeline and evaluate their performance.

Model	Acc	F1	Latency
stella_en_1.5B_v5	0.8961	0.9028	N/A
jina-embeddings-v3*	0.8852	0.9097	N/A
gte-large-en-v1.5*	0.8971	<b>0.9174</b>	N/A
multilingual-e5-large-instruct	0.8954	0.9018	N/A
bge-small-en-v1.5	<b>0.9014</b>	0.9077	N/A

\* The models were not able to complete the evaluation due to memory constraints, Jina evaluated on the 2238/2800 and Gte-large evaluated on the 2322/2800.

Table 5.4: Evaluation Results for Different Embedding Models through the Pipeline (just with the Gemma2 model)

## 5.4. CHUNKING STRATEGIES

Based on the Table 5.4, we use the *bge-small-en-v1.5* model for the subsequent evaluations and analyses due to its superior performance across Acc metric, the F1 score of *gte-large-en-v1.5* is slightly higher, but the model is not able to complete the evaluation due to memory limitations in the same pipeline.

## 5.4 CHUNKING STRATEGIES

A critical component of RAG systems is the chunking strategy employed to divide documents into smaller, manageable pieces for efficient retrieval and processing. This chapter examines three distinct chunking methods for RAG systems, each with its unique characteristics and potential advantages.

### 5.4.1 PARSING DOCUMENTS INTO TEXT CHUNKS

The first method we will explore involves parsing documents into text chunks, also referred to as nodes, of fixed sizes. This approach is straightforward and widely used in many RAG implementations. We will investigate three different chunk sizes: 256, 512, and 1024 tokens.

#### METHODOLOGY

In this method, documents are sequentially divided into chunks of the specified size. If the final chunk is smaller than the designated size, it is typically padded or left as is, depending on the implementation.

#### CHUNK SIZES

- **256-token chunks:** This size offers fine granularity, potentially allowing for more precise retrieval of relevant information. However, it may result in a loss of context for more complex topics that require broader context.
- **512-token chunks:** This medium-sized chunk strikes a balance between granularity and context preservation. It is often considered a good default choice for many applications.
- **1024-token chunks:** Larger chunks preserve more context but may retrieve more irrelevant information and increase computational overhead during retrieval and processing.

## **ADVANTAGES AND LIMITATIONS**

### Advantages:

1. Simple to implement and understand
  2. Consistent chunk sizes facilitate uniform processing

### Limitations:

1. Fixed chunk sizes may not align with natural breaks in the text
  2. Larger chunks can introduce irrelevant information and increase computational costs

## 5.4.2 SMALLER CHILD CHUNKS REFERRING TO BIGGER PARENT CHUNKS (SMALL2BIG)

The second method, which we will refer to as *Small2Big*, involves creating a hierarchical structure of chunks, where smaller child chunks refer to larger parent chunks. This approach aims to combine the benefits of fine-grained retrieval with the context preservation of larger chunks.

## METHODOLOGY

In this method, we parsed documents into three levels of chunks with appending the original text chunk of size 1024:

- Smallest children: 128-token chunks
  - Intermediate parents: 256-token chunks
  - Largest parents: 512-token chunks

Each smaller chunk maintains a reference to its parent chunks, allowing the system to retrieve additional context when needed.

```
1 # ... previous code
2 sub_chunk_sizes = [128, 256, 512]
3 sub_node_parsers = [SimpleNodeParser.from_defaults(chunk_size=c) for
4     c in sub_chunk_sizes]
5
5 all_nodes = []
6 for base_node in base_nodes:
```

#### 5.4. CHUNKING STRATEGIES

```
7     for n in sub_node_parsers:
8         sub_nodes = n.get_nodes_from_documents([base_node])
9         sub_inodes = [
10             IndexNode.from_text_node(sn, base_node.node_id) for sn in
11             sub_nodes
12         ]
13         all_nodes.extend(sub_inodes)
14
15         original_node = IndexNode.from_text_node(base_node, base_node.
16             node_id) # also add original node to node
17         all_nodes.append(original_node)
18     all_nodes_dict = {n.node_id: n for n in all_nodes}
19 # ... continue processing
```

Code 5.1: Small2Big Chunking Method

#### ADVANTAGES AND LIMITATIONS

Advantages:

1. Allows for fine-grained retrieval with the option to expand context
2. Adapts to different levels of specificity required by queries

Limitations:

1. More complex to implement and manage
2. Increased storage requirements due to redundancy in the hierarchy

#### 5.4.3 SENTENCE WINDOW RETRIEVAL

The third method, Sentence Window Retrieval, focuses on maintaining semantic coherence by chunking based on sentences and incorporating surrounding context through windows.

#### METHODOLOGY

In this approach, documents are first split into individual sentences. For each sentence, a *window* of surrounding sentences is included to provide context. We will examine two window sizes: 3 and 6. For each sentence, the chunk includes the sentence itself, one preceding sentence, and one following sentence.

## ADVANTAGES AND LIMITATIONS

Advantages:

1. Preserves semantic units (sentences) and their immediate context
2. Adapts to the natural structure of the text

Limitations:

1. Variable chunk sizes may complicate processing and indexing
2. Optimal window size may vary depending on the document type and content

### 5.4.4 EVALUATION

Each of the three chunking methods presented in this chapter offers distinct advantages and limitations for RAG systems. The choice of method depends on factors such as the nature of the documents, the specific requirements of the application, and the computational resources available. The fixed-size chunking method provides simplicity and consistency but may sacrifice semantic coherence. The Small2Big hierarchical approach offers flexibility in retrieval granularity but introduces complexity in implementation and storage. Sentence Window Retrieval preserves semantic units and adapts to text structure but may result in variable chunk sizes.

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, *Chunking Strategy*, Similarity Cut-off (Re-rank), Top\_k 3

<b>Method</b>	<b>Parameters</b>	Acc	F1	Latency
Original	Chuck Size: 256	0.8914	0.8914	N/A
	Chuck Size: 512	0.8932	0.8993	N/A
	Chuck Size: 1024	0.8946	0.8993	N/A
small2big	Chuck Size: 1024	0.8889	0.8953	N/A
Sliding Window	Window Size: 6	<b>0.9014</b>	0.9077	N/A
	Window Size: 3	<b>0.9014</b>	<b>0.9080</b>	N/A

Table 5.5: Evaluation Results for Different Chunking Strategies through the Pipeline (just with the Gemma2 model)

Examples of the three chunking methods are available in the Appendix B for further reference.

## 5.5. SIMILARITY CUT-OFF

The high baseline accuracy and F1 scores across chunking strategies indicate that the model's performance is largely unaffected by these variations, suggesting that factors other than chunk size, may have a more significant impact on the retrieval process. The Sliding Window method with a window size of 3 achieves the highest accuracy and F1 score, so we will use this method for subsequent evaluations and analyses due to its F1, accuracy scores and having surrounding context.

## 5.5 SIMILARITY CUT-OFF

In this section we will discuss how similarity cut-off can be used to filter out irrelevant nodes and improve the efficiency of the retrieval process. We use Node postprocessors to apply a similarity cut-off to the retrieved nodes, discarding those with a similarity score below a certain threshold. Node postprocessors are a set of modules that take a set of nodes, and apply some kind of transformation or filtering before returning them. For our experiments, we set the similarity cut-off threshold to 0.3, meaning that nodes with a similarity score below 0.3 are discarded, we use the naive score and the re-ranker score to compare the results.

---

**Algorithm 3** Similarity Cutoff Postprocessor

```
1: procedure POSTPROCESSNODES(nodes, knowledge_graph, similarity_cutoff)
2:   new_nodes  $\leftarrow$  []
3:   node_texts  $\leftarrow$  [node.text for node in nodes]
4:   re_rank_nodes  $\leftarrow$  RE_RANK(knowledge_graph, node_texts)
5:   for each node in nodes do
6:     node.score  $\leftarrow$  get_node_score(node.text, re_rank_nodes)
7:     if node.score > similarity_cutoff then
8:       new_nodes.APPEND(node)
9:     end if
10:   end for
11:   return new_nodes
12: end procedure
```

---

Based on the results in Table 5.6, we apply a similarity cut-off on the original score to provide the model with higher-quality data for further evaluations. This approach yields the highest accuracy and F1 score, and it is faster than using the re-ranked score. Please note that in separate evaluations (not reported in Table 5.6), we tested a similarity cut-off using a normalized version of the re-

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, Sliding Window (ws 3), <i>Similarity Cut-off</i> , Top_k 3			
Method	Acc	F1	Latency
w/o similarity cut-off	0.8971	0.9036	N/A
similarity cut-off (original score)	<b>0.9018</b>	<b>0.9080</b>	N/A
similarity cut-off (re-ranked score)	0.9014	<b>0.9080</b>	N/A

Table 5.6: Evaluation Results for Different Similarity Cut through the Pipeline (just with the Gemma2 model)

ranked score, where re-ranked values were scaled between the minimum and maximum of the original scores. However, this normalized re-ranked score did not perform as well as the original scores.

## 5.6 TOP K

Top\_k mentions how many top embeddings to take into context. Considering a large top\_k might go beyond the max\_tokens of the model, we will evaluate the performance of the pipeline with the top\_k set to 3 and 6, and compare the results to determine the optimal value for this parameter.

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, Sliding Window (ws 3), Similarity Cut-off (Original), <i>Top_k</i>			
Method	Acc	F1	Latency
Top_k 3	0.9018	0.9080	N/A
Top_k 6	<b>0.9032</b>	<b>0.9101</b>	N/A

Table 5.7: Evaluation Results for Different Top\_k through the Pipeline (just with the Gemma2 model)

As we decided to set the similarity cut-off with 0.3 threshold, it's good to use the top\_k 6 to have more high-quality embeddings in the context, and based on results it outperforms the top\_k 3, so we will use top\_k 6 for subsequent evaluations.

## 5.7 EVALUATION

In this ablation study, we evaluate the performance of the merging method, which leverages a novel ensemble approach by combining multiple models (Gemma2, Qwen2.5, Llama3.1, and Mistral) to show the robustness and accuracy

## 5.8. FAILURE ANALYSIS

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, Sliding Window (ws 3), Similarity Cut-off (Original), Top\_k 6

Model	award	birth	death	foundationPlace	leader	nbateam	publicationDate	spouse	starring	subsidiary	Total
<i>Positive Labels</i>											
Gemma2	1.0000	0.9733	0.9200	0.5667	0.9933	1.0000	0.9733	0.8000	0.9733	0.9467	0.9147
Qwen2.5	1.0000	0.9667	0.9600	0.7800	0.9933	1.0000	0.9600	0.9400	0.9800	0.9067	0.9487
LLama3.1	0.9800	0.8933	0.8267	0.5067	0.9333	0.8467	0.9267	0.7867	0.8667	0.8400	0.8407
Mistral	1.0000	0.9867	0.9667	0.7333	0.9933	0.9867	0.9867	0.9867	0.9867	0.9533	0.9580
Proposed (At_Most)	1.0000	0.9867	0.9467	0.7067	0.9933	1.0000	0.9667	0.9200	0.9867	0.9333	0.9440
<i>Negative Labels</i>											
Gemma2	0.9923	0.9538	0.9615	0.8154	0.9308	0.6846	0.9308	0.9462	0.9077	0.7769	0.8900
Qwen2.5	0.8000	0.9000	0.9000	0.6769	0.8769	0.5462	0.7923	0.8538	0.7462	0.7615	0.7854
LLama3.1	0.7462	0.8615	0.8462	0.7615	0.8692	0.6385	0.8538	0.7538	0.8231	0.7077	0.7862
Mistral	0.6538	0.8692	0.8462	0.7000	0.8077	0.6077	0.6769	0.7846	0.6846	0.6462	0.7277
Proposed (At_Most)	0.8846	0.9308	0.9077	0.7692	0.8923	0.6615	0.9000	0.8846	0.8385	0.7462	0.8415

Table 5.8: Evaluation Results for Whole Pipeline by Category

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, Sliding Window (ws 3), Similarity Cut-off (Original), Top\_k 6

Model	Stability	Avg. Request Duration	Avg. tokens per request		ACC	F1
Gemma2	0.8720	N/A	N/A		<b>0.9032</b>	<b>0.9101</b>
Qwen2.5	0.8685	N/A	N/A		0.8729	0.8888
LLama3.1	0.8291	N/A	N/A		0.8154	0.8299
Mistral	0.8650	N/A	N/A		0.8511	0.8733
Proposed (At_Most)	<b>0.9176</b>	N/A	N/A		0.8964	0.9071

Table 5.9: Evaluation Results for Whole Pipeline

of the pipeline. As presented in Tables 5.8 and 5.9, along with Figures 5.2, 5.3, and 5.4, the proposed ensemble method consistently outperforms individual models across both positive and negative labels, achieving a more balanced and comprehensive performance.

Note that the ensemble method is based on the tie-breaking strategy discussed in Section 3.6.2, with *At\_Most* as the merging method. Based on provided configurations, the model selected for *At\_Most* is the *Gemma2:21B* model.

## 5.8 FAILURE ANALYSIS

To gain a deeper understanding of the limitations and challenges faced by our fact-checking system, we conducted a comprehensive failure analysis using the *FactBench* dataset. By examining the instances where our system, based on the majority vote, failed to correctly verify the facts, we aimed to identify the main error types and provide insights into the reasons behind these failures.

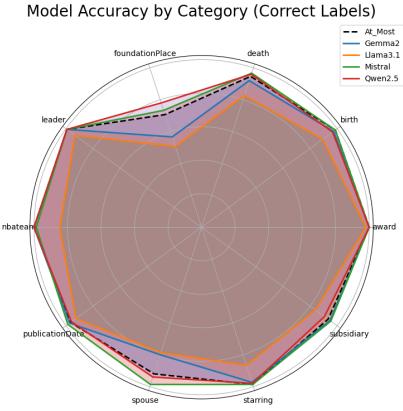


Figure 5.2: Positive Labels

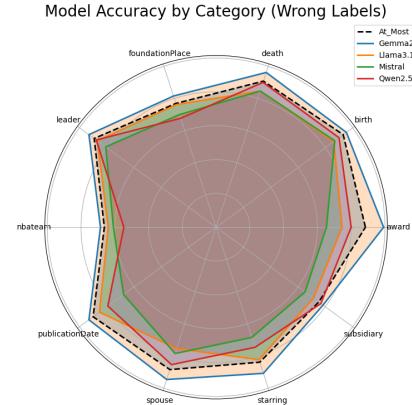


Figure 5.3: Negative Labels

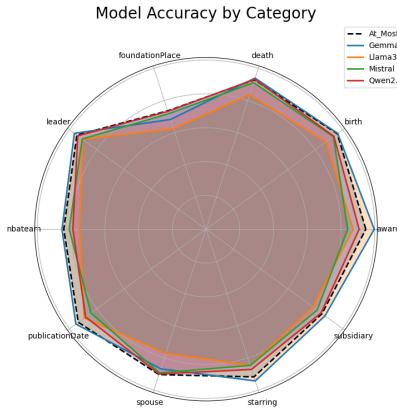


Figure 5.4: All Labels

## ERROR TYPE CATEGORIZATION

After analyzing the failure cases, we categorized the errors into five main types:

- **Insufficient or Irrelevant Context:** In some cases, the provided context information does not directly support or refute the given triple. The LLMs struggle to make accurate judgments when the necessary facts are missing or the available information is tangentially related to the claim.
- **Misinterpretation of Relationships:** The LLMs sometimes misinterpret the relationships between entities mentioned in the context. They may confuse family relations, professional associations, or the nature of events.
- **Over reliance on Keyword Matching:** In some instances, the LLMs rely too heavily on surface-level keyword matching rather than understanding the underlying semantics. The presence of certain words or phrases can lead to incorrect assumptions.

## 5.8. FAILURE ANALYSIS

- **Lack of Common Sense Reasoning:** The LLMs can struggle with applying common sense knowledge or reasoning about the plausibility of claims. They may fail to consider the unlikelihood of certain scenarios or relationships.
- **Difficulty with Negation and Contradiction:** The LLMs sometimes struggle with handling negation or identifying contradictions between the triple and the context information. They may overlook explicit denials or fail to recognize inconsistencies.

To quantify the distribution of error types, we classified each error instance into one or more of five categories, allowing for overlaps when errors fit multiple categories.

Error Type	Triple	Description
Insufficient or Irrelevant Context	Ai Sugiyama birth place Yokohama	Since there's no information in any of the documents about Ai Sugiyama being born in Yokohama, and one document explicitly states her birthplace as Tokyo. So LLMs infer that Ai Sugiyama was born in Tokyo, Japan and not Yokohama, Japan.
Mis Interpretation of Relationships	Mitt Romney office Dallas	LLMs mistakenly infer that Romney has an office in Dallas based on his attendance at a fundraiser there. Attending an event doesn't imply having a permanent office.
Over reliance on Keyword Matching	Robbie Williams office Los Angeles	LLMs wrongly assume Robbie Williams has an office in Los Angeles due to text discussing his purchase or sell of a property there, not an office.
Lack of Common Sense Reasoning	Saul Bellow starring Nobel Prize in Literature	LLMs fail to recognize "starring" is inappropriate for receiving a Nobel Prize. Common sense suggests terms like "awarded" or "received."
Difficulty with Negation and Contradiction	John William Strutt team Nobel Prize in Physics	LLMs confirm the presence of a team, despite no evidence suggesting Strutt was part of one. The context mentions Strutt's Nobel Prize but lacks information about team involvement.

Table 5.10: Examples of Failure Cases and Error Analysis

Based on Figures 5.5 and 5.6, The "foundationPlace" relation shows the highest number of total instances and errors in both charts. This suggests that the model struggles most with verifying facts about the locations where organizations or institutions were founded. The large discrepancy between correct and incorrect predictions for this relation indicates a significant challenge in accurately processing location-based information.

Relations such as "birth", "death", and "spouse" show varying levels of difficulty. While "birth" and "death" have relatively few instances, "spouse" has

a moderate number of cases with a notable error rate. This suggests that verifying personal information, especially relationships, poses challenges for the model, and mostly related to having multiple marriages or relationships during a lifetime.

The "nbateam" and "subsidiary" relations, which involve organizational affiliations, show moderate error rates. This indicates that the model has some difficulty in correctly identifying professional associations and corporate structures.

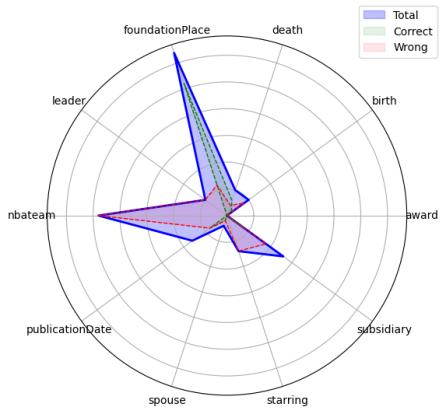


Figure 5.5: Distribution of Fully Incorrect Predictions 4/4

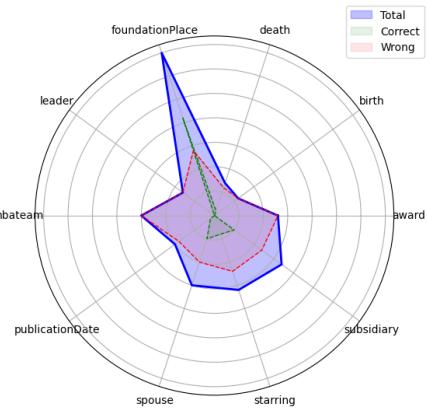


Figure 5.6: Distribution of Partially Incorrect Predictions 3/4

In general, by comparing Figures 5.5 and 5.6, we can observe that the overall shape of the error distribution is similar, indicating consistency in the model's performance across different voting thresholds



# 6

## Conclusions and Future Works

A	B
C	D
E	F
G	H

Table 6.1: Table example





## Prompt Templates

In this section, we present the prompt templates used in the pipeline.

### A.1 HUMAN-UNDERSTANDABLE TEXT GENERATION PROMPT

— Prompt template for generating human-readable text —

**Task Description:**

Convert a kg triple into a meaningful human readable sentence.

**Instructions:**

Given a subject, predicate, and object from a kg, form a grammatically correct and meaningful sentence that conveys the relationship between them.

**Examples:**

**Input:**

Subject: Alexander\_III\_of\_Russia

Predicate: isMarriedTo

Object: Maria\_Feodorovna\_Dagmar\_of\_Denmark

Output: {"output" : "Alexander III of Russia is married to Maria Feodorovna, also known as Dagmar of Denmark."}

**Input:**

Subject: Quentin\_Tarantino

Predicate: produced

## A.2. QUESTION GENERATION PROMPT

```
Object: From_Dusk_till_Dawn  
Output: {"output": "Quentin Tarantino produced the film  
From Dusk till Dawn."}
```

Input:

```
Subject: Joseph_Heller  
Predicate: created  
Object: Catch-22  
Output: {"output": "Joseph Heller created the novel Catch-22."}
```

Do the following:

Input:

```
Subject: {knowledge_graph.subject}  
Predicate: {knowledge_graph.predicate}  
Object: {knowledge_graph.object}
```

The output should be a JSON object with the key "output" and the value as the sentence. The sentence should be human-readable and grammatically correct. The subject, predicate, and object can be any valid string without having extra information.

## A.2 QUESTION GENERATION PROMPT

Prompt template for generating 10 questions for each triple  
You are an intelligent system with access to a vast amount of information. I will provide you with a knowledge graph in the form of triples (subject, predicate, object).

Your task is to generate ten questions based on the kg. The questions should assess understanding and insight into the information presented in the graph.

Provide the output in JSON format, with each question having a unique identifier. Instructions:

1. Analyze the provided knowledge graph.
2. Generate ten questions that are relevant to the information in kg.
3. Provide the questions in JSON format, each with a unique identifier.

Input Knowledge Graph: Albert Einstein bornIn Ulm, Germany

```

Expected Response: {
  "questions": [
    {"id": 1,
     "question": "Where was Albert Einstein born?"},
    {"id": 2,
     "question": "What is Albert Einstein known for?"},
    {"id": 3,
     "question": "In what year was the Theory of Relativity published?"},
    {"id": 4,
     "question": "Where did Albert Einstein work?"},
    {"id": 5,
     "question": "What prestigious award did Albert Einstein win?"},
    {"id": 6,
     "question": "Which theory is associated with Albert Einstein?"},
    {"id": 7,
     "question": "Which university did Albert Einstein work at?"},
    {"id": 8,
     "question": "What did Albert Einstein receive the Nobel Prize in?"},
    {"id": 9,
     "question": "In what field did Albert Einstein win a Nobel Prize?"},
    {"id": 10,
     "question": "Name the city where Albert Einstein was born."}
  ]
}

Considering the above information, please respond to this kg: {query}
The output should be in JSON format with each question having a unique
identifier and question doesn't contain term knowledge graph, without
any additional information

```

### A.3 RAG PROMPT

Prompt template for RAG
<p>Context information is below.</p> <hr/> <p>{context_str}</p> <hr/> <p>Given the context information and without prior knowledge, Evaluate whether the information in the documents supports the triple. Please provide your answer in the form of a structured JSON</p>

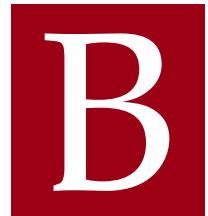
### A.3. RAG PROMPT

format containing a key "output" with the value as "yes" or "no".  
If the triple is correct according to the documents, the value  
should be "yes". If the triple is incorrect, the value should be "no".

{few\_shot\_examples}

Query: {query\_str}

Answer:



## Chunking Strategies

As discussed in Section 5.4, the method used to chunk input text is a critical decision in the design of a RAG system. Here, we present concrete examples of how different chunking strategies affect the segmentation of text, using the *correct\_award\_00000* entry from the FactBench dataset. The text is as follows:

*Henry Dunant award Nobel Peace Prize*

The model used for these examples is *Gemma2* with *similarity\_top\_k* set to 3, and *BAAI/bge-small-en-v1.5* as embedding model. The documents are selected using *ms-marco-MiniLM-L-6-v2* discussed in 5.2.2. We report the best node found by through our pipeline for each chunking strategy.

Chunk	Score
Abstract When Jean Henry Dunant received the first Nobel Peace Prize in 1900, he was praised for "the supreme humanitarian achievement of the nineteenth century." This praise was merited, for Dunant had led the creation of both the international Red Cross and the First Geneva Convention. The Red Cross has since saved countless lives and relieved human suffering around the world. The Geneva Convention established that those treating war wounded, wearing a red cross, would not be attacked. With this Convention, Dunant began the creation of international humanitarian law to reduce the suffering caused by war. Despite Dunant's vital contributions, he has been largely forgotten. This article briefly tells the story of this dedicated humanitarian leader and of his great achievements. Recommended Citation McFarland, Sam (2017) "A Brief History of An Unsung Hero and Leader – Jean Henry Dunant and the Founding of the Red Cross at the Geneva Convention," International Journal of Leadership and Change: Vol. 5: Iss. 1, Article 5. Available at: <a href="https://digitalcommons.wku.edu/ijlc/vol5/iss1/5">https://digitalcommons.wku.edu/ijlc/vol5/iss1/5</a>	90.5791
In 1901, Henry Dunant was co-awarded the first Nobel Peace Prize in recognition of his devotion to the humanitarian cause. (Español): Henry Dunant, Premio Nobel de la Paz - En 1901, Henry Dunant fue co-galardonado con el primer Premio Nobel de la Paz en reconocimiento a su devoción por la causa humanitaria. Credit: ICRC / Vincent Varin / www.icrc.org	90.5830
To celebrate the memory and work of Henry Dunant, on the centenary of the presentation of the first Nobel Peace Prize, rightly awarded to Dunant for his having founded the institution of the International Red Cross, this paper presents the reader with some insights into his activities and sufferings, his trials and tribulations, and the hope and strength of his character. The ceaseless efforts made by Dunant to bring about the Institution which today represents Hope for so many suffering people who are silent victims of wars and atrocities, are fleetingly presented. The authors' intention is to give due recognition to Dunant for his work, and to highlight the humanity and the moral and social worth of the face behind the International Red Cross.	90.6017

Table B.1: Text Splitter - Chuck Size 512

Chunk	Score
Henry Dunant The Nobel Peace Prize 1901 Nobel co-recipient: Frédéric Passy Role: Founder of the International Committee of the Red Cross, Geneva, Originator Geneva Convention (Convention de Genève) Nobel Prize Cash and Philanthropy Jean Henry Dunant, though poor, donated his Nobel Prize money to charity. Hans Daae, a military physician, managed to get the money deposited in a bank in Norway. Thus Dunant's creditors could not claim the money. When Dunant was alive the money remained untouched in the bank. He lived frugally in a Swiss nursing home. Dunant's will bequeathed one half of the money to the Norwegian Red Cross and the Norwegian Women's Public Health Association. The will bequeathed the other half of the money to charities in Switzerland. Daae was also responsible for Dunant being awarded the Noble Prize.	90.6243

Table B.2: Small to Big (base retriver chunk size set to 1024)

---

**Window - Highlighted Text is Original Text**

---

You can read more about that here: From the first Nobel Prize award ceremony, 1901 The announcement that the founder of the Red Cross had been chosen as Peace Prize laureate met with mixed reactions. Dunant had been awarded the prize for ameliorating the suffering of wounded soldiers, not for organising peace congresses or reducing standing forces, as stipulated in Alfred Nobel's will. The Nobel Committee had chosen a broad interpretation of the provision that a laureate should "further fraternity between nations".

The Red Cross: three-time recipient of the Peace Prize Henry Dunant (1828–1910). Switzerland, "for his humanitarian efforts to help wounded soldiers and create international understanding" Frédéric Passy (1822–1912). France, "for his lifelong work for international peace conferences, diplomacy and arbitration."

On 10th of December 1901 the first Nobel Peace Prize was awarded. It went to Henry Dunant, founder of the International Committee of the Red Cross, who shared the first Nobel Peace Prize with Frédéric Passy, a leading international pacifist of the time. Since then, the Red Cross has been awarded the Peace Prize three times. The Red Cross: Three-time recipient of the Peace Prize Four of them given out in Stockholm and one, the Peace Prize, in Christiania, as Oslo was then called. You can read more about that here: From the first Nobel Prize award ceremony, 1901 The announcement that the founder of the Red Cross had been chosen as Peace Prize laureate met with mixed reactions. Dunant had been awarded the prize for ameliorating the suffering of wounded soldiers, not for organising peace congresses or reducing standing forces, as stipulated in Alfred Nobel's will.

Henry Dunant The Nobel Peace Prize 1901 Nobel co-recipient: Frédéric Passy Role: Founder of the International Committee of the Red Cross, Geneva, Originator Geneva Convention (Convention de Genève) Nobel Prize Cash and Philanthropy Jean Henry Dunant, though poor, donated his Nobel Prize money to charity. Hans Daae, a military physician, managed to get the money deposited in a bank in Norway. Thus Dunant's creditors could not claim the money.

When Dunant was alive the money remained untouched in the bank. He lived frugally in a Swiss nursing home. Dunant's will bequeathed one half of the money to the Norwegian Red Cross and the Norwegian Women's Public Health Association.

---

Table B.3: Sliding Window - Window Size 3



# References

- [1] Meta AI. *Introducing LLaMA 3: Advancing Open Foundation Models*. <https://ai.meta.com/blog/meta-llama-3-1/>. Accessed: 2024-10-17. 2023.
- [2] Krisztian Balog and Robert Neumayer. “A test collection for entity search in DBpedia”. In: *Proceedings of the 36th International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR ’13. Dublin, Ireland: Association for Computing Machinery, 2013, pp. 737–740. ISBN: 9781450320344. doi: 10.1145/2484028.2484165. URL: <https://doi.org/10.1145/2484028.2484165>.
- [3] Hyung Won Chung et al. *Scaling Instruction-Finetuned Language Models*. 2022. doi: 10.48550/ARXIV.2210.11416. URL: <https://arxiv.org/abs/2210.11416>.
- [4] Alexis Conneau and Douwe Kiela. “SentEval: An Evaluation Toolkit for Universal Sentence Representations”. In: *arXiv preprint arXiv:1803.05449* (2018).
- [5] Alexis Conneau et al. “Unsupervised Cross-lingual Representation Learning at Scale”. In: *CoRR* abs/1911.02116 (2019). arXiv: 1911.02116. URL: <http://arxiv.org/abs/1911.02116>.
- [6] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: 1810.04805 [cs.CL]. URL: <https://arxiv.org/abs/1810.04805>.
- [7] Alphaeus Dmonte et al. *Claim Verification in the Age of Large Language Models: A Survey*. 2024. arXiv: 2408.14317 [cs.CL]. URL: <https://arxiv.org/abs/2408.14317>.
- [8] Abhimanyu Dubey et al. *The Llama 3 Herd of Models*. 2024. arXiv: 2407.21783 [cs.AI]. URL: <https://arxiv.org/abs/2407.21783>.

## REFERENCES

- [9] Luis Galárraga et al. "AMIE: Association rule mining under incomplete evidence in ontological knowledge bases". In: May 2013, pp. 413–422. doi: [10.1145/2488388.2488425](https://doi.org/10.1145/2488388.2488425).
- [10] Daniel Gerber et al. "DeFacto—Temporal and multilingual Deep Fact Validation". In: *Journal of Web Semantics* 35 (2015). Machine Learning and Data Mining for the Semantic Web (MLDMSW), pp. 85–101. issn: 1570-8268. doi: <https://doi.org/10.1016/j.websem.2015.08.001>. url: <https://www.sciencedirect.com/science/article/pii/S1570826815000645>.
- [11] Michael Günther et al. *Jina Embeddings 2: 8192-Token General-Purpose Text Embeddings for Long Documents*. 2024. arXiv: 2310.19923 [cs.CL]. url: <https://arxiv.org/abs/2310.19923>.
- [12] Faegheh Hasibi, Krisztian Balog, and Svein Erik Bratsberg. "Dynamic Factual Summaries for Entity Cards". In: *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. SIGIR '17. Shinjuku, Tokyo, Japan: Association for Computing Machinery, 2017, pp. 773–782. isbn: 9781450350228. doi: [10.1145/3077136.3080810](https://doi.org/10.1145/3077136.3080810). url: <https://doi.org/10.1145/3077136.3080810>.
- [13] Edward J Hu et al. "LoRA: Low-Rank Adaptation of Large Language Models". In: *International Conference on Learning Representations*. 2022. url: <https://openreview.net/forum?id=nZeVKeeFYf9>.
- [14] Gautier Izacard et al. *Unsupervised Dense Information Retrieval with Contrastive Learning*. 2022. arXiv: 2112.09118 [cs.IR]. url: <https://arxiv.org/abs/2112.09118>.
- [15] Albert Q. Jiang et al. *Mistral 7B*. 2023. arXiv: 2310.06825 [cs.CL]. url: <https://arxiv.org/abs/2310.06825>.
- [16] M. Abdul Khalil et al. *RAGAR, Your Falsehood Radar: RAG-Augmented Reasoning for Political Fact-Checking using Multimodal Large Language Models*. 2024. arXiv: 2404.12065 [cs.CL]. url: <https://arxiv.org/abs/2404.12065>.
- [17] Aditya Kusupati et al. *Matryoshka Representation Learning*. 2024. arXiv: 2205.13147 [cs.LG]. url: <https://arxiv.org/abs/2205.13147>.
- [18] Nayeon Lee et al. *Factuality Enhanced Language Models for Open-Ended Text Generation*. 2023. arXiv: 2206.04624 [cs.CL]. url: <https://arxiv.org/abs/2206.04624>.

- [19] Zehan Li et al. *Towards General Text Embeddings with Multi-stage Contrastive Learning*. 2023. arXiv: 2308.03281 [cs.CL]. URL: <https://arxiv.org/abs/2308.03281>.
- [20] Jerry Liu. *Tweet on Information Retrieval and LLMs*. Accessed: 2024-10-10. 2023. URL: <https://twitter.com/jerryjliu0/status/1708147687084986504>.
- [21] Yinhua Liu et al. *RoBERTa: A Robustly Optimized BERT Pretraining Approach*. 2019. arXiv: 1907.11692 [cs.CL]. URL: <https://arxiv.org/abs/1907.11692>.
- [22] Shayne Longpre, Yi Lu, and Joachim Daiber. *MKQA: A Linguistically Diverse Benchmark for Multilingual Open Domain Question Answering*. 2020. URL: <https://arxiv.org/pdf/2007.15207.pdf>.
- [23] S. Marchesin, G. Silvello, and O. Alonso. “Veracity Estimation for Entity-Oriented Search with Knowledge Graphs”. In: *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management (CIKM '24), October 21–25, 2024, Boise, ID, USA*. ACM, 2024. doi: [10.1145/3627673.3679561](https://doi.org/10.1145/3627673.3679561).
- [24] Mistral AI Team. *Announcing Mistral 7B*. Accessed: 2024-10-17. 2023. URL: <https://mistral.ai/news/announcing-mistral-7b/>.
- [25] Niklas Muennighoff et al. “MTEB: Massive Text Embedding Benchmark”. In: *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*. Ed. by Andreas Vlachos and Isabelle Augenstein. Dubrovnik, Croatia: Association for Computational Linguistics, May 2023, pp. 2014–2037. doi: [10.18653/v1/2023.eacl-main.148](https://doi.org/10.18653/v1/2023.eacl-main.148). URL: <https://aclanthology.org/2023.eacl-main.148>.
- [26] Prakhar Ojha and Partha Talukdar. “KGEval: Accuracy Estimation of Automatically Constructed Knowledge Graphs”. In: *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Ed. by Martha Palmer, Rebecca Hwa, and Sebastian Riedel. Copenhagen, Denmark: Association for Computational Linguistics, Sept. 2017, pp. 1741–1750. doi: [10.18653/v1/D17-1183](https://doi.org/10.18653/v1/D17-1183). URL: <https://aclanthology.org/D17-1183>.
- [27] Nils Reimers and Iryna Gurevych. “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks”. In: *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Nov. 2019. URL: <https://arxiv.org/abs/1908.10084>.

## REFERENCES

- [28] Stephen Robertson et al. “Okapi at TREC-3.” In: Jan. 1994, pp. 0-.
- [29] Jon Saad-Falcon et al. *Benchmarking and Building Long-Context Retrieval Models with LoCo and M2-BERT*. 2024. arXiv: 2402.07440 [cs.IR]. URL: <https://arxiv.org/abs/2402.07440>.
- [30] Soumya Sanyal et al. *Are Machines Better at Complex Reasoning? Unveiling Human-Machine Inference Gaps in Entailment Verification*. 2024. arXiv: 2402.03686 [cs.CL]. URL: <https://arxiv.org/abs/2402.03686>.
- [31] Saba Sturua et al. *jina-embeddings-v3: Multilingual Embeddings With Task LoRA*. 2024. arXiv: 2409.10173 [cs.CL]. URL: <https://arxiv.org/abs/2409.10173>.
- [32] Jianlin Su et al. *RoFormer: Enhanced Transformer with Rotary Position Embedding*. 2023. arXiv: 2104.09864 [cs.CL]. URL: <https://arxiv.org/abs/2104.09864>.
- [33] Fabian Suchanek et al. *YAGO 4.5: A Large and Clean Knowledge Base with a Rich Taxonomy*. 2024. arXiv: 2308.11884 [cs.AI]. URL: <https://arxiv.org/abs/2308.11884>.
- [34] Gemma Team et al. *Gemma 2: Improving Open Language Models at a Practical Size*. 2024. arXiv: 2408.00118 [cs.CL]. URL: <https://arxiv.org/abs/2408.00118>.
- [35] Qwen Team. *Qwen2.5: A Party of Foundation Models*. Sept. 2024. URL: <https://qwenlm.github.io/blog/qwen2.5/>.
- [36] Nandan Thakur et al. *BEIR: A Heterogenous Benchmark for Zero-shot Evaluation of Information Retrieval Models*. 2021. arXiv: 2104.08663 [cs.IR]. URL: <https://arxiv.org/abs/2104.08663>.
- [37] Liang Wang et al. *Multilingual E5 Text Embeddings: A Technical Report*. 2024. arXiv: 2402.05672 [cs.CL]. URL: <https://arxiv.org/abs/2402.05672>.
- [38] Shitao Xiao et al. *C-Pack: Packaged Resources To Advance General Chinese Embedding*. 2023. arXiv: 2309.07597 [cs.CL].
- [39] An Yang et al. “Qwen2 Technical Report”. In: *arXiv preprint arXiv:2407.10671* (2024).
- [40] Zhenrui Yue et al. *Retrieval Augmented Fact Verification by Synthesizing Contrastive Arguments*. 2024. arXiv: 2406.09815 [cs.CL]. URL: <https://arxiv.org/abs/2406.09815>.

## REFERENCES

- [41] Xin Zhang et al. *mGTE: Generalized Long-Context Text Representation and Reranking Models for Multilingual Text Retrieval*. 2024. arXiv: 2407.19669 [cs.CL]. URL: <https://arxiv.org/abs/2407.19669>.
- [42] Xinyu Zhang et al. “MIRACL: A Multilingual Retrieval Dataset Covering 18 Diverse Languages”. In: *Transactions of the Association for Computational Linguistics* 11 (2023), pp. 1114–1131. doi: 10.1162/tacl\_a\_00595. URL: <https://aclanthology.org/2023.tacl-1.63>.
- [43] Xuan Zhang and Wei Gao. *Reinforcement Retrieval Leveraging Fine-grained Feedback for Fact Checking News Claims with Black-Box LLM*. 2024. arXiv: 2404.17283 [cs.CL]. URL: <https://arxiv.org/abs/2404.17283>.



# Acknowledgments