

UNIVERSITÀ
DEGLI STUDI
DI PADOVA



DEPARTMENT OF INFORMATION ENGINEERING

MASTER THESIS IN COMPUTER ENGINEERING

FactCheck: Knowledge Graph Fact Verification Through Retrieval-Augmented Generation Using a Multi-Model Ensemble Approach

MASTER CANDIDATE

Farzad Shami

Student ID 2090160

SUPERVISOR

Prof. Stefano Marchesin

University of Padova

CO-SUPERVISOR

Prof. Gianmaria Silvello

University of Padova

ACADEMIC YEAR
2024/2025

DATE: 27/03/2025

Give It Up, Sid. You Know Humans Can't Talk
Ice Age (2002)

Abstract

In today’s world of Artificial intelligence (AI) and big data, knowledge graphs (KGs) play an important role in powering many AI systems, search engines, and decision-support systems. Small errors can propagate through connected systems and cause big problems, so ensuring their accuracy is a critical task. This thesis addresses this challenge by introducing FactCheck, a fact-checking system for KGs. Our method uses Retrieval-Augmented Generation (RAG) coupled with multiple language models to verify facts. FactCheck works by generating questions about each KG fact, retrieving relevant documents, splitting them into chunks, and then feeding chunks as input to the large language models (LLMs). Then the majority vote system with dispute resolution decides on the fact’s correctness by considering the generated responses. We tested our approach on three real-life datasets—FactBench, YAGO, and DBpedia—whereby comparing the FactCheck output with gold standard labels, we achieved prediction performance rates of 90, 87, and 70 percent, respectively. On average, verifying a single fact requires processing about 1,550 tokens per LLM, and takes about 7 minutes to reach a final decision. These metrics demonstrate the system’s resource usage and performance. For achieving these results, we tuned different components of RAG pipeline by selecting the best parameters/models for document selection, embedding, and chunking through systematic testing. The system offers a reliable and scalable solution that is compatible with various KG environments and can be adapted to handle different types of facts.

Contents

List of Figures	xi
List of Tables	xiii
List of Algorithms	xvii
List of Code Snippets	xvii
List of Acronyms	xix
1 Introduction	1
1.1 Problem Statement	3
1.2 Proposed System	3
1.3 Contributions	4
1.4 Thesis Structure	5
1.5 Significance and Potential Applications	6
2 Related Works	7
2.1 Knowledge Graph and Fact-Checking	7
2.1.1 DeFacto	8
2.1.2 Validating RDF Triples Using Textual Evidence	8
2.2 LLMs and Retrieval-Augmented Generation	9
2.2.1 Entailment Verification and Language Models	10
2.2.2 Claim Verification in the Age of Large Language Models .	11
2.2.3 RAG-Based Fact Verification by Synthesizing Contrastive Arguments	13
2.2.4 RAGAR: RAG-Augmented Reasoning for Political Fact-Checking using Multimodal LLMs	15
2.2.5 Best Practices	16

CONTENTS

2.3 Cost Estimation	17
3 FactCheck	19
3.1 KG Dataset	21
3.1.1 Definition and Purpose	21
3.1.2 Structure and Components	22
3.1.3 Role in the Overall system	22
3.2 Query Generation and Processing	22
3.2.1 Human-Understandable Text Generation	23
3.2.2 Question Formulation Techniques	23
3.2.3 Query Scoring	25
3.3 Information Retrieval Mechanisms	27
3.3.1 Google Search Integration	27
3.3.2 Process and Extract Links	28
3.3.3 Data Pool Creation	31
3.3.4 Data filtering	31
3.4 Embedding and Retrieval Tasks	32
3.4.1 Embedding Techniques for Smaller Chunks	32
3.4.2 Similarity Cutoff Strategy	33
3.5 LLMs	34
3.5.1 Integration of Multiple LLMs	34
3.5.2 Roles of LLMs in the system	35
3.6 Voting System and Conflict Resolution	36
3.6.1 Majority Voting System	36
3.6.2 Conflict Resolution Strategies	37
3.7 Performance Report	39
3.8 system Flow and Decision Points	39
3.9 Ethical Considerations and Limitations	42
3.9.1 Ethical Considerations	42
3.9.2 Limitations	43
4 Empirical Evaluation	45
4.1 Dataset Analysis	45
4.1.1 FactBench Dataset	45
4.1.2 YAGO Dataset	46
4.1.3 DBpedia Dataset	47

4.2	Candidate Models	48
4.2.1	Gemma2	48
4.2.2	Qwen2.5	49
4.2.3	Llama3.1	50
4.2.4	Mistral	51
4.3	Experimental Setup	52
4.3.1	Performance Metrics and Evaluation	52
4.3.2	System Configurations	54
4.4	Comparative Analysis	55
4.4.1	Discussion of Results	55
4.4.2	Qualitative Error Analysis	57
4.4.3	DBpedia Analysis in Depth	62
5	Ablation Study	69
5.1	Evaluation Methodology	70
5.1.1	Iterative Optimization Process	70
5.1.2	Evaluation Metrics	70
5.1.3	Significance of the Methodology	71
5.2	Document Selection	71
5.2.1	Unsupervised Methods	72
5.2.2	Supervised Methods	73
5.2.3	Evaluation with Large Language Models	75
5.3	Embedding Models	77
5.3.1	Gte-large-en-v1.5	77
5.3.2	Jina-embeddings-v3	78
5.3.3	Stella_en_1.5B_v5	79
5.3.4	Multilingual-e5-large-instruct	80
5.3.5	bge-small-en-v1.5	81
5.3.6	Comparative Analysis	81
5.4	Chunking Strategies	84
5.4.1	Parsing Documents into Text Chunks	84
5.4.2	Smaller Child Chunks Referring to Bigger Parent Chunks (Small2Big)	84
5.4.3	Sentence Window Retrieval	85
5.4.4	Advantages and Limitations	86
5.4.5	Evaluation	87

CONTENTS

5.5	Similarity Cut-off	88
5.6	Top K	89
5.7	Evaluation	90
5.8	Failure Analysis	91
6	Conclusions and Future Works	95
	Appendices	99
A	Prompt Templates	99
A.1	Human-understandable text generation Prompt	99
A.2	Question Generation Prompt	100
A.3	RAG Prompt	101
A.4	Reasoning Prompt	102
B	Chunking Strategies	103
	References	107
	Acknowledgments	113

List of Figures

1.1	Overview of the FactCheck system for fact verification. The system processes a knowledge graph triple through multiple stages: (1) Knowledge Graph (KG) humanization, converting structured data into natural language; (2) Generation of aspect-specific questions to probe the fact; (3) Google search retrieval based on generated questions; (4) Analysis of retrieved documents containing relevant information; (5) Deployment of multiple Large Language Models (LLMs) as open-source fact-checkers; (6) Final fact verification through majority voting; and (7) tie-breaker module.	3
2.1	Distinctions between human and LLM Inferences. The entailment prediction performance of humans and LLMs are depicted by a 5-star rating scale [39].	10
2.2	Comparison of claim verification systems between NLP-based (traditional) and LLM-based for claim veracity [7].	12
2.3	The proposed RAFTS [52], which performs few-shot fact verification by incorporating informative in-context demonstrations and contrastive arguments with nuanced information derived from the retrieved documents	14
2.4	An overview of the fact-checking pipeline contrasting the baseline Sub-Question Generation approach from the Chain of RAG and Tree of RAG approach followed by veracity prediction and explanation.	16
3.1	FactCheck: RAG-based KG fact verification system	20

LIST OF FIGURES

3.2	Cross-Encoder component, which assesses the similarity of generated questions by taking multiple input pairs and assigning a relevance score to each, supporting question evaluation and refinement.	25
3.3	Knowledge Distillation Process for Enhanced Re-Ranking Efficiency	26
3.4	Fetching results from Google Search engine for top N questions and the main KG query.	28
3.5	Extracted content from the crawled URLs using newspaper4k library.	31
4.1	Collecting logs and leveraging LLM-generated reasoning, combined with contextual document embeddings (jxm/cde-small-v1) [30], to cluster errors using a hierarchical density-based spatial technique.	58
4.2	Model overlap heatmaps by category and dataset. Each cell shows the percentage overlap in errors between model pairs. Matrices are organized by error category (UnLabeled, Relationship, Role Errors, etc.) and dataset (DBpedia, FactBench, YAGO), revealing patterns in how models agree or disagree when making verification errors.	60
4.3	Normalized distribution of error clusters across datasets.	61
4.4	Distribution of error clusters across selected LLMs.	61
4.5	Distribution of tendency to be wrong across gemma2, qwen2.5, LLama3.1 and mistral models. The right chart illustrates the distribution of fully incorrect predictions (4/4) detailing the instances where all predictions made by the models were incorrect. The left chart depicts the distribution of just one wrong predictions (1/4).	61
4.6	Partition-wise model performance comparison: Accuracy and F1-scores for KG fact verification on DBpedia dataset. Gray bars indicate stratum weights (log scale).	63
4.7	Error distribution analysis across different language models and frequency strata.	64
4.8	distribution of error categories across different language models and frequency strata	65

LIST OF FIGURES

4.9	Comparative analysis of language model error rates across knowledge domains. (Left) domain-specific error rates across 13 knowledge categories, with overlaid sample distribution bars. (Right) distribution of error rates for each model	66
5.1	Document retrieval confusion matrix based on Jaccard similarity between documents retrieved by each model.	75
5.2	Node sentence window replacement technique as described by Liu [25].	86
5.3	Category-wise performance of different models in identifying Positive Labels (left) and Negative Labels (right) on the FactBench dataset.	90
5.4	Prediction accuracy on the FactBench dataset, focusing on incorrect predictions. The right chart illustrates the Distribution of Fully Incorrect Predictions (4/4), detailing the instances where all predictions made by the models were incorrect. The left chart depicts the Distribution of Partially Incorrect Predictions (3/4). . .	93

List of Tables

3.1	Examples of human-understandable text generation, illustrates entries from multiple KGs, detailing the transformation of raw triples into readable sentences.	24
3.2	Example of generated questions by the Gemma2 model with relevance scores assigned by the Jina Re-ranker Cross-Encoder.	27
3.3	Performance of our LLM-based tasks in production, generated by Openlit with Gemma2 model.	40
3.4	Performance of our information retrieval mechanisms.	40
3.5	Comprehensive list of decision points in the system flow.	41
3.6	Limitations of the system, categorized by scope of knowledge. . .	43
4.1	Statistical summary of FactBench, YAGO, and DBpedia datasets .	48
4.2	Summary of key strengths of selected candidate LLMs for KG fact verification.	51
4.3	System configurations for empirical evaluation	54
4.4	Empirical evaluation results of the proposed system and candidate LLMs over the FactBench, YAGO, and DBpedia.	56
4.5	Statistical analysis of output tokens and request times per query across FactBench, YAGO, and DBpedia datasets for each used model.	56
4.6	Statistical analysis of request time per query across FactBench, YAGO, and DBpedia datasets.	57
4.7	Dataset-wise error clustering based on LLM-generated reasoning, using Contextual Document Embeddings for embeddings, UMAP, and HDBSCAN.	60
4.8	Partition-wise evaluation results of the proposed system and candidate LLMs over the DBpedia dataset.	63

LIST OF TABLES

5.1	Performance comparison of various distilled MS MARCO models based on BERT architecture, measured across NDCG@10 on TREC DL 2019 and MRR@10 on MS MARCO Dev benchmarks.	74
5.2	Performance evaluation of various document retrieval methods on the FactBench dataset, using the Gemma2 model.	76
5.3	Comparison of characteristics of embedding models	82
5.4	Performance evaluation of various embedding models on the FactBench dataset, using the Gemma2 model.	83
5.5	Advantages and Limitations of different chunking strategies for RAG systems.	87
5.6	Performance evaluation of various chunking strategy on the FactBench dataset, using the Gemma2 model.	88
5.7	Performance evaluation of similarity cut-off method on the FactBench dataset, using the Gemma2 model.	89
5.8	Performance evaluation of different Top_k retrieval strategies on the FactBench dataset using the Gemma2 model.	89
5.9	Category-wise performance evaluation results of various models on the FactBench dataset.	91
5.10	Performance evaluation of various models on the FactBench dataset. 91	
5.11	Example of failure cases and error analysis observed in the FactBench dataset using generated results and explanations.	92
B.1	Evaluation of text segmentation using a chunk Size of 512, text chunks derived from the entry "Henry Dunant award Nobel Peace Prize".	104
B.2	Evaluation of text segmentation using a Small to Big technique (base chunk size 1024), text chunks derived from the entry "Henry Dunant award Nobel Peace Prize".	105
B.3	Evaluation of text segmentation using a Sliding Window with window size 3, text chunks derived from the entry "Henry Dunant award Nobel Peace Prize".	106

List of Algorithms

1	Resolve Ties in Majority Voting System	39
2	Calculate Model Consistency Per Model	53
3	Similarity Cutoff Postprocessor (re-rank score)	88

List of Code Snippets

3.1	Crawling the extracted URLs	29
5.1	Small2Big chunking method	85

List of Acronyms

CSV Comma Separated Values

KG Knowledge Graph

AI Artificial Intelligence

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

LIST OF CODE SNIPPETS

GQA Grouped-Query Attention

SWA Sliding Window Attention

RS Rejection Sampling

DPO Direct Preference Optimization

SPO Subject-Predicate-Object

RDF Resource Description Framework

1

Introduction

KGs have become an indispensable component of modern information systems and are widely used in different categories, including search engines [12], recommendation systems [14], and question-answering platforms [33]. KG is a directed, multi relational graph, which represents entities as nodes and their relationships as edges, and can be used as an abstraction of the real world [2]. These KGs have appears as an effective tool for arranging and querying massive amounts of structured data. Maintaining the accuracy and dependability of KGs is a significant challenge, as they have become a standard in representing factual information across different domains. Prominent examples of KGs are DBpedia [21], Wikidata [47], Freebase [3], and Knowledge Vault [8].

On the other side, in the past few years, the increasing volume of online content, coupled with the rise of misinformation and disinformation, necessitate the need ways to verify the correctness of the information. So, ensuring the truthfulness of the information is an important task to maintain the reliability in knowledge-driven applications, particularly as Artificial Intelligence (AI) and automated decision-making systems become more prevalent.

Considering that KGs are often incomplete and noisy, performing fact verification on KGs is a challenging task. To evaluate the correctness of a fact, humans need to check the information against several external sources, which can be time-consuming and error-prone. Evaluating the reliability of each source is another challenge. This process can take even professional fact-checkers a long time to complete, and it's not always easy to know if the information is correct. As this process is time-consuming and the amount of information is increasing

day by day, we can notice that manual fact-checking is not a feasible solution.

To address these challenges, there were two types of methods in the past that automate the process of fact-checking: (1) Statistical and rule-based techniques and (2) Sampling techniques that try to reduce the human effort in the fact-checking process by using the samples the data from datasets. The first type of methods can work for some types of facts¹, but they don't work well for more complicated or subtle data. The second type of methods can be more effective, but they do not check all aspects of the data, and they may not be able to verify all the facts.

In response to these limitations, the past few years have witnessed the growing adoption of Machine Learning (ML) and Natural Language Processing (NLP) techniques as new possibilities in this field. These techniques have shown remarkable success and change the traditional methods to be more effective and somehow automatic. With advances in NLP and the emergence of Large Language Models (LLMs), the fact-checking process changed dramatically.

LLMs have shown remarkable capabilities in understanding and generating human-like text, making them well-suited for tasks that require reasoning and understanding of textual information. However, LLMs have some limitations, such as (1) hallucination, where they generate plausible but incorrect information; (2) training data cut-off, where they are limited to the data they were trained on, which may become outdated over time; (3) reliability of training data, where the models may not be able to verify the correctness of the information they generate; and (4) the models contain significant information in their parameters, but their performance is not as good as task-specific architectures. To address these limitations, researchers have proposed combining LLMs with external information retrieval to improve the accuracy and reliability of fact verification and named this technique Retrieval-Augmented Generation (RAG) [23].

By combining the strengths of LLMs with external information, we introduce a novel system for automated fact verification in KGs using RAG. In this thesis, we introduce FactCheck, a system that employs a multi-model ensemble approach by merging the outputs of different LLMs through majority voting to verify facts in KGs.

¹In this work, we use the terms fact, claim and triple interchangeably.

1.1 PROBLEM STATEMENT

This thesis tackles the issue of automated fact verification in KGs with a RAG-based methodology. Our objective is to create a system capable of (1) Retrieve relevant information from external sources to support or refute claims in a KG; (2) Utilize LLMs to reason about the retrieved information and generate accurate assessments of fact truthfulness; (3) Handle a wide range of fact types and domains, from simple statements to more complex relational facts; and finally (4) Using multiple LLMs, provide multiple responses for its verification decisions, enhancing transparency and trust in the system.

1.2 PROPOSED SYSTEM

Our proposed system combines several key components to create a fact verification system as you can see the overview in Figure 1.1.

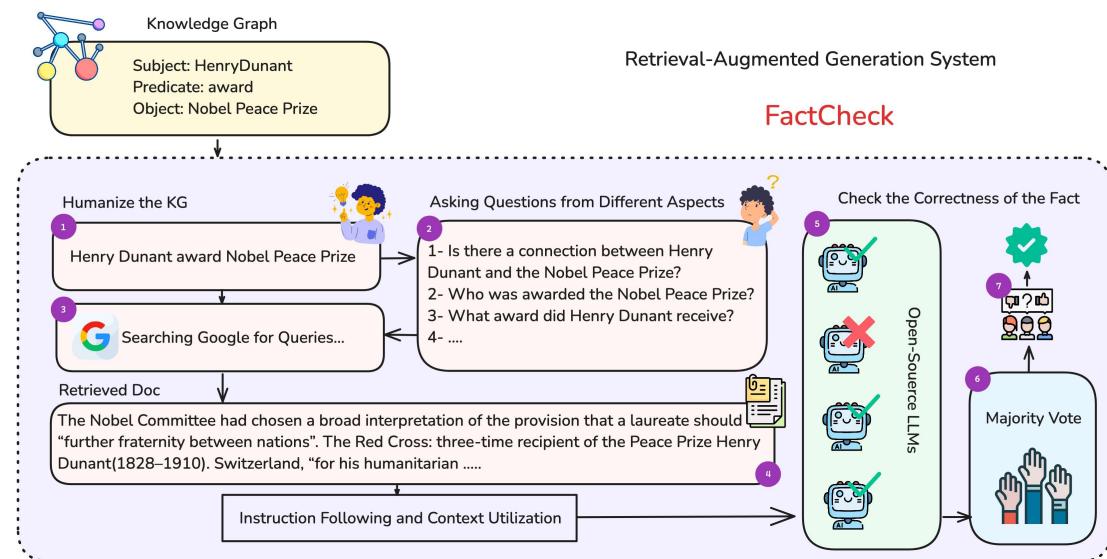


Figure 1.1: Overview of the FactCheck system for fact verification. The system processes a knowledge graph triple through multiple stages: (1) KG humanization, converting structured data into natural language; (2) Generation of aspect-specific questions to probe the fact; (3) Google search retrieval based on generated questions; (4) Analysis of retrieved documents containing relevant information; (5) Deployment of multiple LLMs as open-source fact-checkers; (6) Final fact verification through majority voting; and (7) tie-breaker module.

1.3. CONTRIBUTIONS

The general idea behind the FactCheck system is as follows:

- **KG Representation:** We start by representing facts from the KG in a format suitable for processing by language models and Information Retrieval (IR) systems. This involves converting the subject-predicate-object triples of the KG 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.
- **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.

1.3 CONTRIBUTIONS

This thesis presents several significant contributions to the field of NLP and knowledge verification systems. Our investigation into the capabilities of LLMs within RAG frameworks reveals their considerable potential for fact-verification tasks. Through experimentation, we have proved that these LLMs show remarkable aptitude in adapting and contextualizing external knowledge sources

which is a crucial capability for reliable information validation in contemporary applications. By examining LLMs reasoning patterns and pathways, we provide insights into how these models arrive at verification decisions when confronted with factual claims of varying complexity. This deeper understanding improves our theoretical grasp of these systems and suggests practical improvements for their implementation in real-world scenarios.

A particularly noteworthy aspect of this work is the assessment of computational efficiency. Contrary to prevailing assumptions about LLMs that they require substantial computational resources, our findings indicate that effective fact-verification processes can indeed be executed on modestly provisioned local infrastructures. This finding is important because it makes advanced AI tools more accessible to different research groups, even those with limited resources.

Finally, We conduct error taxonomy that categorizes the failure predication observed during fact-verification task. Through this analysis, we have identified distinct error categories including context mismatches, relationship errors, role attribution discrepancies, geographic inaccuracies, and more. This categorization exposes the current limitations of these systems and suggest ways for targeted improvements in subsequent model post-training or fine-tuning.

To evaluate FactCheck, extensive experiments were conducted using four open-source LLMs — Gemma2, Qwen2.5, Llama3.1, and Mistral. These models ranged in size from 7B to 9B parameters, covering various capabilities for dense retrieval, logical reasoning, and interpretability. Experiments were carried out on 13,530 facts in total. We used accuracy, F1 scores, computational latency, and resource consumption as core metrics on our evaluations. Additional ablation studies further validated the efficiency of our system.

1.4 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 system aims to address. Chapter 3 presents a detailed description of our proposed FactCheck system for fact verification. It explains each component of the system, including the rationale behind design choices and implementation

1.5. SIGNIFICANCE AND POTENTIAL APPLICATIONS

details. Chapter 4 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 KG fact verification. Chapter 5 Presents a study that investigates the impact of various parameters on the system's overall performance, while also exploring different methodologies for each component to determine the optimal final setting configuration for the system. Finally, chapter 6 concludes the thesis by summarizing the contributions, discussing limitations of the current approach, and outlining promising directions for future research.

1.5 SIGNIFICANCE AND POTENTIAL APPLICATIONS

The development of effective fact verification systems for KGs has far-reaching implications across various domains. Our system helps ensure accurate information by automatically checking facts, making large knowledge bases more reliable - especially important in today's world, where false information spreads quickly.

It is useful in fields like healthcare, finance, and law, where decisions rely on correct data. It also helps in education by allowing students to verify information and improve their digital literacy. The same techniques can be used for content moderation, helping social media and news platforms detect false or misleading content. In science, the system can check research claims, compare findings, and spot inconsistencies in studies. By improving fact-checking in KGs, this thesis supports the goal of building more reliable information systems.

2

Related Works

This thesis builds upon prior research in KG fact verification, LLMs, RAG, and entailment verification. In this section, we provide an overview of the relevant literature across these areas. We categorized them into different sections based on their relevance to the proposed work.

2.1 KNOWLEDGE GRAPH AND FACT-CHECKING

KGs can involve in the Fact-Checking task in two ways: (1) Using KG structured data to verify the data or facts and (2) Using the structured or unstructured information to verify the information inside the KG triples.

Fact-Checking using KG Structured Data: By exploiting the structured data in the KG, we can find that whether the claim (*i.e.* text) is correct or not. For example for finding if a person is father of another person, the first step is checking the relations like *fatherOf* in the KG or checking. Afterward we can check the relations like *birthdate* to find the age of the person as the father should be older than the child. As our purpose is to verify the facts in the KG, this approach is out of the scope of this thesis.

Knowledge Graph Fact-Checking using External Data: We focus on the KG fact checking approach, where we use the unstructured information from external sources to verify the information inside the triples.

2.1. KNOWLEDGE GRAPH AND FACT-CHECKING

2.1.1 DeFACTO

DeFacto [11] is designed to check if facts are true by finding supporting evidence on the web. What makes DeFacto special is that it works in multiple languages and can determine when facts were true.

Gerber et al. [11] created a system that takes a fact (in RDF triple format) and searches the web for evidence that confirms it. DeFacto looks for mentions of the fact in English, German, and French websites. It then calculates a confidence score to tell users how likely the fact is to be true. The system also tries to figure out the time period when the fact was valid. DeFacto first converts the fact into natural language patterns that might appear on websites. It then searches the web using these patterns and analyzes the returned pages for confirming text. The system also considers how trustworthy each website is. For the time aspect, DeFacto looks for year mentions near where the fact is discussed. It uses patterns and frequency analysis to determine when a fact was valid. This is important because many facts are only true during specific time periods.

Gerber et al. tested DeFacto on a dataset they created called *FactBench*, which contains 1,500 facts across 10 different types of relations (like births, deaths, marriages, etc.). Their results showed that DeFacto achieves about 85% accuracy in determining if facts are true. For figuring out when facts were true, the accuracy was around 70% for simple time points and 66% for time periods. One important finding was that using multiple languages improved performance significantly. By combining evidence from English, German, and French websites, DeFacto performed better than when using only English.

Our work has some intersection with DeFacto, but there are some key differences. The intersection is as follows: (1) we use web search to verify facts. We feed the web search results as context to the LLMs to verify the facts. and additionally, (2) We use the dataset that they created (*i.e.* FactBench) to evaluate our system. The differences are as follows: (1) We use LLMs to verify the facts. (2) We use the RAG technique to retrieve the evidence and (3) we use the ensemble of LLMs to verify the facts to have more reliable results.

2.1.2 VALIDATING RDF TRIPLES USING TEXTUAL EVIDENCE

Syed et al. [43] introduces a new approach for automatically validating facts in KGs. They present a system that improves upon existing fact validation

methods by combining deep sentence parsing with topic coherence analysis to evaluate the truthfulness of RDF triples based on textual evidence from reference corpora. This represents a significant advancement over previous approaches like DeFacto (§2.1.1), which relied primarily on string matching and word proximity. The system architecture follows a systematic process: First, input RDF triples are verbalized into natural language statements. These verbalizations are used to search reference corpora (either Wikipedia or ClueWeb) for supporting textual evidence. For each piece of evidence, the system extracts features through dependency parsing to identify direct relationships between subject, predicate, and object. Additionally, it calculates topic coherence values to measure the semantic relatedness of terms. These features, along with others shared with DeFacto, serve as input to machine learning classifiers that determine a confidence score for each triple.

The authors also analyzed the impact of corpus size, finding that larger reference corpora improved performance for both systems. Feature analysis confirmed the importance of the newly introduced dependency parsing and topic coherence features in achieving these performance gains.

The system still relies on the presence of explicit textual evidence in the reference corpora, and the computational complexity of dependency parsing could pose scalability challenges for very large knowledge bases. The paper also doesn't extensively discuss how the system handles more complex or composite facts. In our work, we address these limitations by leveraging LLMs to verify facts in knowledge graphs, which can handle more complex reasoning tasks and provide more nuanced assessments of fact veracity. Also, we do the similar analysis on the impact of the corpus size and the features used in the system for the *DBpedia* dataset.

2.2 LLMs AND RETRIEVAL-AUGMENTED GENERATION

In this section we will discuss the related works on LLMs and RAG for general fact-verification task. Eventually, in the § 2.2.5 we will discuss the best practices for the RAG based systems.

2.2.1 ENTAILMENT VERIFICATION AND LANGUAGE MODELS

In the paper “Minds versus Machines: Rethinking Entailment Verification with Language Models”, Sanyal et al. [39] 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.



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 [39].

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 [4] 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

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.

2.2.2 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; and finally, (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. Based on survey, several studies have demonstrated the effectiveness of LLM-based approaches in claim verification:

2.2. LLMS AND RETRIEVAL-AUGMENTED GENERATION

- Zhang and Gao [55] introduced the Hierarchical Step-by-Step (HiSS) prompting method, which directs LLMs to separate a claim into several sub-claims and then verify each via multiple questions-answering steps progressively, improving performance on complex news claim verification tasks.
- Lee et al. [20] 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.

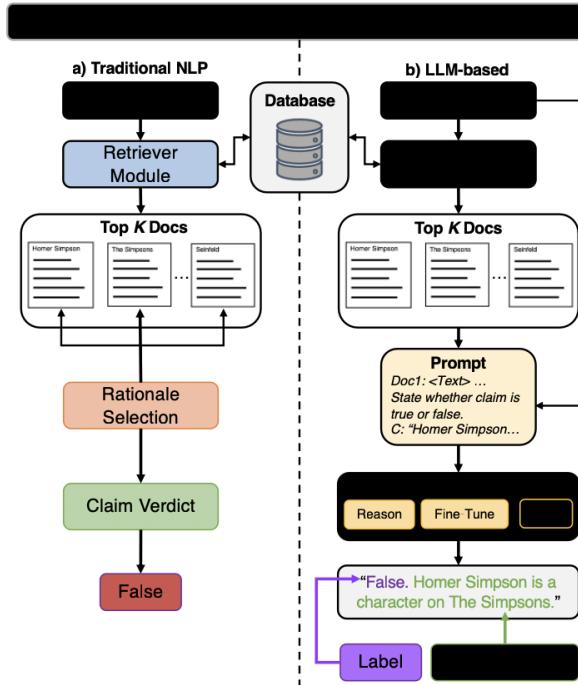


Figure 2.2: Comparison of claim verification systems between NLP-based (traditional) and LLM-based for claim veracity [7].

Our approach shares similarities with the LLM-based pipeline described by Dmonte et al., particularly in the use of retrieval-augmented generation and 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. **Limited Explanation:** While many LLM approaches provide explanations, our method places a stronger emphasis on generating binary (pass/fail) labels for claims to reduce the costs of using LLMs.
4. **Diverse LLM Model:** We use multiple LLMs with diverse architectures to provide more reliable verification result.

These distinctions position our work as a new 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 fact 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.

2.2.3 RAG-BASED FACT VERIFICATION BY SYNTHESIZING CONTRASTIVE ARGUMENTS

The paper Retrieval-Augmented Fact Verification by Synthesizing Contrastive Arguments [52] 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

2.2. LLMS AND RETRIEVAL-AUGMENTED GENERATION

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.

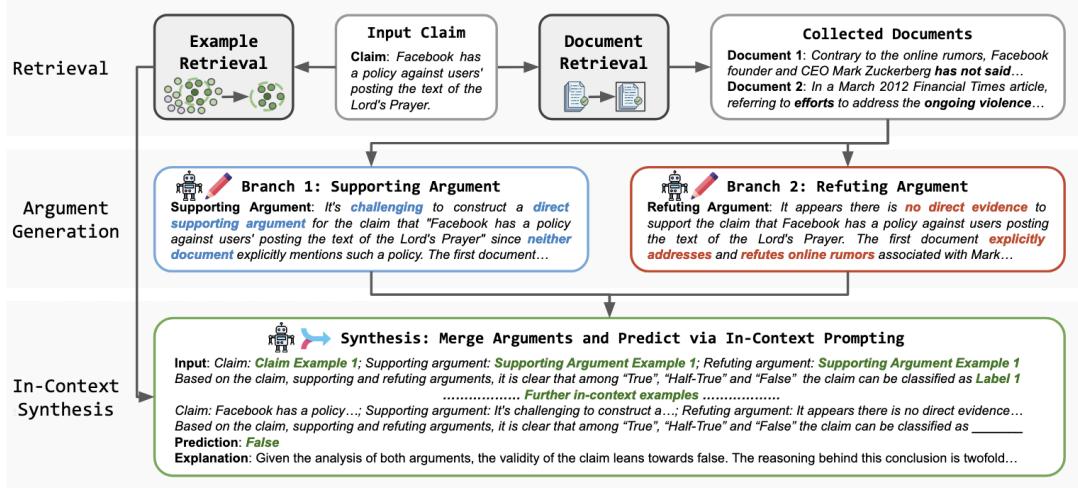


Figure 2.3: The proposed RAFTS [52], which performs few-shot fact verification by incorporating informative in-context demonstrations and contrastive arguments with nuanced information derived from the retrieved documents

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 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.

2.2.4 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” [18] introduces a new methodology to political fact-checking by leveraging RAG with multimodal LLMs. This work focuses on enhancing fact verification in the politically sensitive domain, where disinformation 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.

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 accuracy, 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.

2.2. LLMS AND RETRIEVAL-AUGMENTED GENERATION

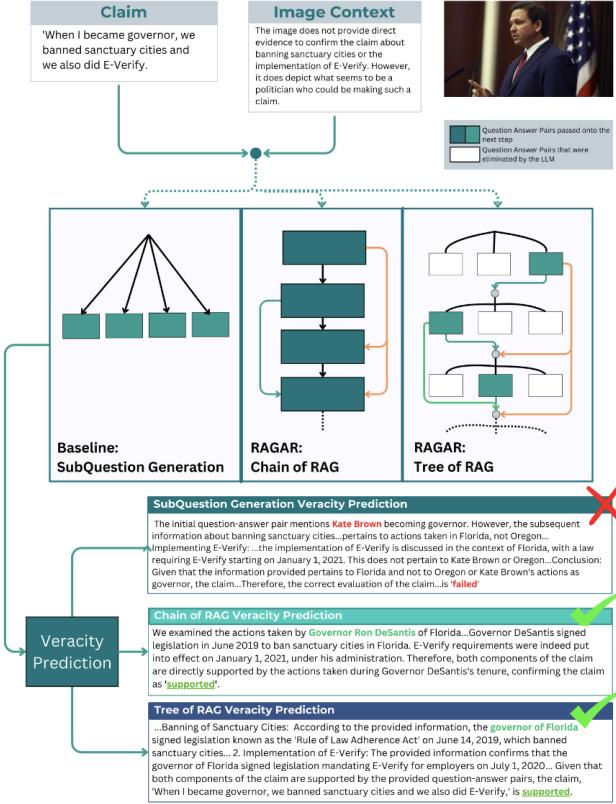


Figure 2.4: An overview of the fact-checking pipeline contrasting the baseline Sub-Question Generation approach from the Chain of RAG and Tree of RAG approach followed by veracity prediction and explanation.

2.2.5 BEST PRACTICES

Different practices for RAG-based systems have been explored in several recent studies for finding the best and optimal parameters/models for different components of the RAG systems. Studies [50, 26] try to find the best practices for different RAG components by examining the impact of different strategies on the performance of LLMs. These studies are evaluating different aspects of the RAG based system and provide detailed comparison between the components from the retrieval to the generation part of the system.

In the same way we try to conduct cold ablation study on fact verification task and provide detailed comparison between the components of the system. Insights from these studies can help us choose initial parameters. Our work also go beyond these studies by providing the timing analysis that help users to know the trade-offs.

2.3 COST ESTIMATION

Heiko Paulheim [34] addresses an important yet underexplored aspect of KG research: the economic cost of creating KGs. While KGs have been extensively analyzed in terms of size, overlap, and quality, Paulheim provides a perspective by quantifying the financial investment required for both manual and automated KG creation. The author offers cost estimates for several well-known KGs. For manually curated resources, Cyc [22] costs approximately \$5.71 per statement (totaling \$120M for 21M assertions), while Freebase costs about \$2.25 per statement (\$6.75B total for 3B facts). In contrast, automatically created KG are substantially more cost-effective [36]: DBpedia costs only \$0.0125 per statement (\$5.1M development cost for 400M statements), YAGO costs \$0.0083 per statement (combining Wikipedia extraction with WordNet), and NELL costs \$0.1425 per statement. The paper demonstrates that automatic KG creation is approximately 15-250 times more cost-effective than manual curation, providing a clear economic argument for automation. In general, We noticed that the manual curation of KGs is expensive and time-consuming, while automatic creation is more cost-effective and scalable.

Here, in this thesis we also provide the cost of using LLMs by providing the amount of tokens used for each model. Through this we can estimate the cost of using LLMs for fact verification in KGs. Additionally, we provide a time analysis of the models to estimate the time required to verify a fact in a KG in the local infrastructure setup.

3

FactCheck

This chapter outlines our a multi-stage NLP system (FactCheck) developed specifically for fact verification in KGs. The system integrates advanced techniques and methodologies to address core challenges in verifying information accuracy, interpreting user queries, retrieving relevant data from external sources (*e.g.* large-scale web searches), and generating factually consistent responses.

The FactCheck build on top of LLMs and focus on graph data, the system is able to provide the correctness of different aspects of each fact inside the KG datasets. The system is designed to handle and effectively process different categories of facts, such as people, places, events, and abstract concepts. This system is specifically designed for fact verification in KG datasets (*e.g.* FactBench, DBpedia, YAGO, NELL), but this approach can be easily adopted and extended to other categories of fact-verification systems by changing some parts of the system.

There are some key features in the system: (1) To access diverse and up-to-date sources, the system incorporates Google Search, blending structured knowledge with real-time information retrieval to improve verification accuracy. (2) To ensure the correctness of facts the system performs precise analyses by dividing retrieved data into manageable chunks and utilizing embedding techniques, and (3) To make the system more reliable, reduce errors and manage potential biases, the decision-making framework in the system leverages majority voting across LLMs to decide and employs a final judge to resolve inconsistencies, creating a balanced and coherent output.

The system's general architecture is illustrated in Figure 3.1, highlighting the key components and their interactions.

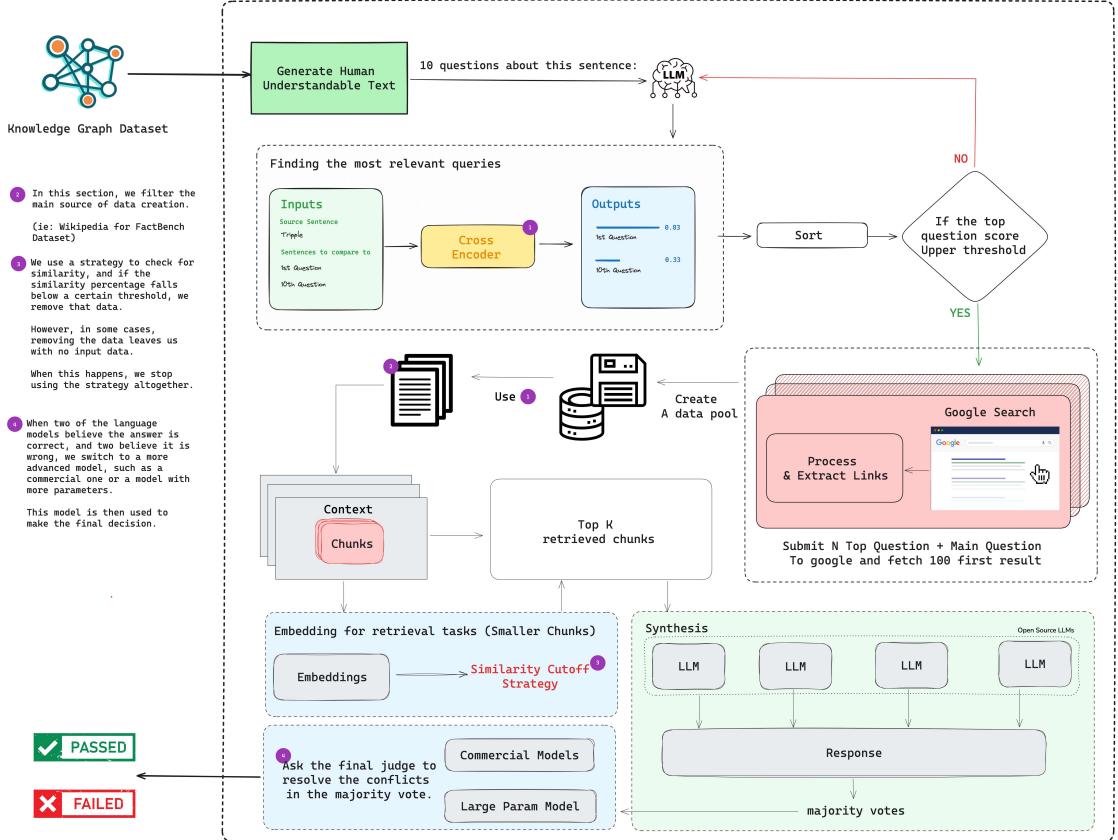


Figure 3.1: FactCheck: RAG-based KG fact verification system

This chapter aims to provide a complete analysis of each stage of the system as showed on Figure 3.1, 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 KG 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. Also, as we proceed, we will explore the ethical considerations and potential biases inherent in the system, highlighting

ing the importance of responsible AI development and the need for ongoing scrutiny and oversight in the deployment of advanced technologies.

Note that in this chapter we just outline each component and its role in the system, and in the next chapters, we will provide parameters we used in each component and the results of the system.

3.1 KG DATASET

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

3.1.1 DEFINITION AND PURPOSE

In the form of a graph, a KG is an ordered way to show information that models real-world things and how they relate to each other. The KG 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 KG fulfills multiple essential functions:

- **Semantic Representation:** Unlike traditional relational databases, KGs 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:** KGs 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.

3.2. QUERY GENERATION AND PROCESSING

3.1.2 STRUCTURE AND COMPONENTS

The KG Dataset in our system 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.
- **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.

Many KGs currently represent extracted facts in the form of Subject-Predicate-Object (SPO) triples which is in line with the standard prescribed by Resource Description Framework (RDF). In this case the subject and object are nodes, and the predicate is the edge connecting them.

3.1.3 ROLE IN THE OVERALL SYSTEM

As illustrated in the system diagram, the KG dataset is the thing that we want to verify the correctness of it. The system uses the KG to generate queries, retrieve relevant information, and synthesize responses to find the correctness of the KG.

3.2 QUERY GENERATION AND PROCESSING

The Query Generation and Processing step, following to the KG dataset, is a point wherein spo triples (*i.e.* facts) are converted into human-understandable text suitable for efficient processing by later components. This text is then used to generate questions. These questions alongside the text are then used to retrieve information from external sources. This section clarifies the techniques and methodologies utilized in this phase for subsequent actions in the information retrieval process.

3.2.1 HUMAN-UNDERSTANDABLE TEXT GENERATION

In the first step of this stage, we use a LLM (*e.g.* LLama3 [9], Gemma2 [44]) 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. This step can be done using other techniques, such as rule-based systems, but the LLMs are more efficient as they do not require manual intervention and can be easily adapted to different types of data.

Key aspects of this process include:

- **Contextual Awareness:** Integrating relevant context from the KG 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 KG datasets are human-readable, whereas others are not; for instance, the FactBench [11] dataset may be easily comprehended by concatenating the subject, predicate, and object of each triple. Table 3.1 provides examples of human-understandable text generation from various KGs datasets. This table shown that the generated text is more readable and more understandable than the raw triples. Also, we can figure out that the LLMs can be compatible with different types of KGs datasets.

3.2.2 QUESTION FORMULATION TECHNIQUES

A cornerstone of our system is its ability to generate questions about the input spo triple, as illustrated in the diagram. We generate 10 questions about the input sentence to ensure all the spo triple aspects are covered and disambiguated.

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.

3.2. QUERY GENERATION AND PROCESSING

Table 3.1: Examples of human-understandable text generation, illustrates entries from multiple KGs, detailing the transformation of raw triples into readable sentences.

KG			Source	Is Generated ?	Final Text
Subject	Predicate	Object			
<i>Albert Einstein</i>	<i>Birth Place</i>	<i>Ulm, Germany</i>	FactBench [11]	✗	<i>Albert Einstein birth place Ulm, Germany</i>
<i>Chris Benoit</i>	<i>deathPlace</i>	<i>Fayetteville, Georgia</i>	FactBench [11]	✗	<i>Chris Benoit death place Fayetteville, Georgia</i>
<i>Alexander_III_of_Russia</i>	<i>isMarriedTo</i>	<i>Maria_Feodorovna_Dag-mar_of_Denmark_</i>	YAGO [42]	✓	<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>

- **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.

Implementing this technique involves using LLMs to generate 10 questions about the input sentence, using the models’ language understanding capabilities to ensure the questions are relevant and contextually appropriate. The prompt template used to guide the question generation process is reported in Appendix A.2.

From now as we want to submit the generated questions to the Google Search engine, we named the generated questions as queries.

3.2.3 QUERY SCORING

In this step, we want to evaluate the similarity of the generated queries to the input sentence. This way we can filter out the irrelevant queries and just keep the relevant ones. For this purpose we use the Cross-Encoder component, which assesses the similarity of generated queries by taking multiple input pairs and assigning a similarity score to each, supporting query evaluation and refinement.

Cross-Encoder Component: The Cross-Encoder component is essential for evaluating the similarity of the generated queries. As indicated in the Figure 3.2, this module takes multiple inputs and produces relevance scores for each question.

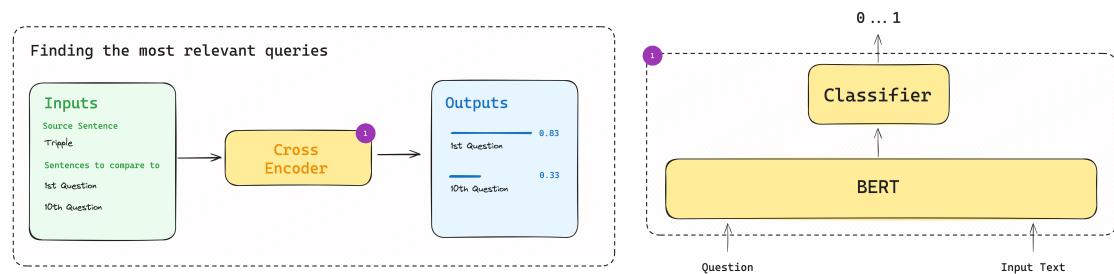


Figure 3.2: Cross-Encoder component, which assesses the similarity of generated questions by taking multiple input pairs and assigning a relevance score to each, supporting question evaluation and refinement.

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 3.2) 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.

3.2. QUERY GENERATION AND PROCESSING

In our case we use the *jinaai/jina-reranker-v1-turbo-en*¹ 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 JinaBERT [13] model as its foundation. The model employs a process known as knowledge distillation as shown in the Figure 3.3, to attain exceptional speed and efficiency, making it the optimal selection for our system.

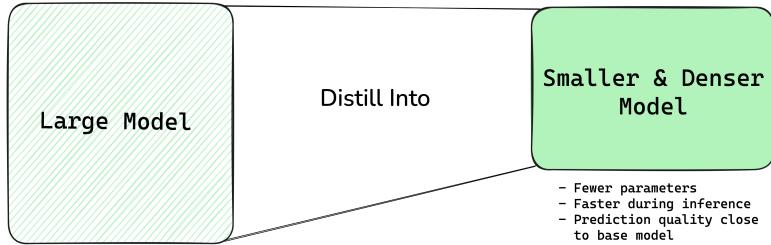


Figure 3.3: Knowledge Distillation Process for Enhanced Re-Ranking Efficiency

Similarity Threshold and Sorting: Following the Cross-Encoder’s scoring, the system 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 system.
- **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.

Table 3.2 provides an example of generated questions by the Gemma2 model with similarity scores assigned by the Jina Re-ranker Cross-Encoder. The top question is highlighted, indicating that it has the highest similarity score and is likely to be the most pertinent for further processing. If the top question’s score exceeds the predefined threshold, it proceeds to the next stage of the system. Otherwise, the system must loop back to generate new queries till the threshold is met.

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

Table 3.2: Example of generated questions by the Gemma2 model with relevance scores assigned by the Jina Re-ranker Cross-Encoder.

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?	0.9491
	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

3.3 INFORMATION RETRIEVAL MECHANISMS

The IR Mechanisms are an essential element of FactCheck, 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 system diagram. This integration serves to expand the information horizon beyond the confines of our internal KG Dataset. Key aspects of this integration include:

- **Query Submission:** The system submits the N top queries (where N is a predefined number) along with the human-understandable text 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

3.3. INFORMATION RETRIEVAL MECHANISMS

nature of the internal KG.

- **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.

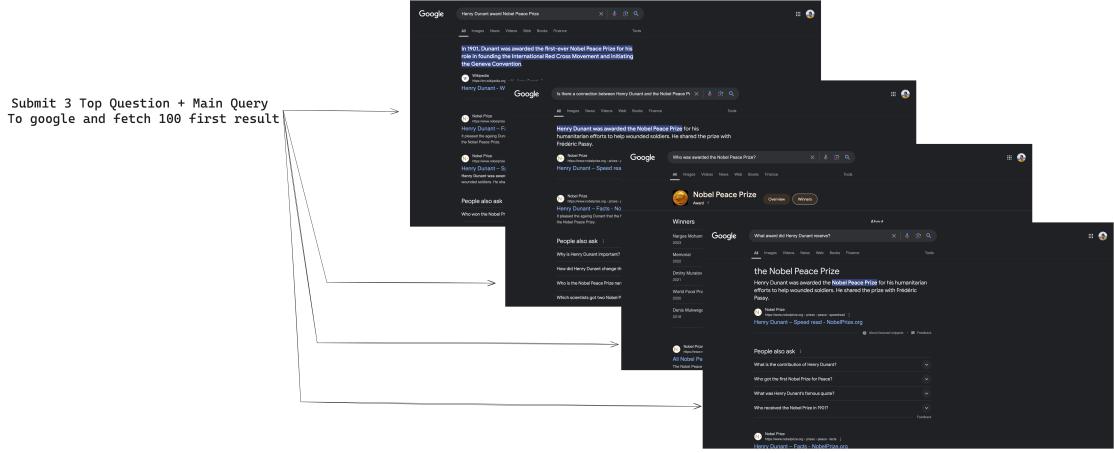


Figure 3.4: Fetching results from Google Search engine for top N questions and the main KG query.

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 system, 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.

3.3.2 PROCESS AND EXTRACT LINKS

Following the retrieval of search results, the system incorporates an important 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*². *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
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:

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

3.3. INFORMATION RETRIEVAL MECHANISMS

- **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*³ library to handle the first issue, and for the second and third issues, we can use the *Scrapy*⁴ 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 system, 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*⁵ library to extract the content of the page. *newspaper4k* is a Python 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.

³<https://www.selenium.dev/>

⁴<https://scrapy.org/>

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

- **Image Extraction:** Identifies and extracts images associated with articles.

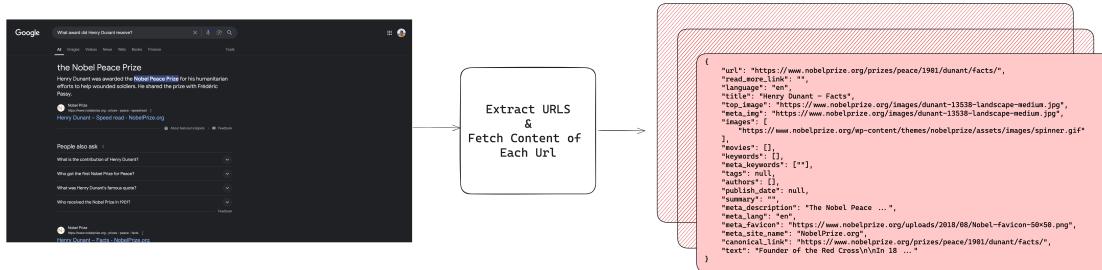


Figure 3.5: Extracted content from the crawled URLs using newspaper4k library.

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 system.

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 KG.
- **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 KG. 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.

3.4. EMBEDDING AND RETRIEVAL TASKS

On the other hand, we just keep the top 10 relevant information for the system using cross-encoders discussed in Section 3.2.3, this way we ensure about the quality of the data.

Specifically (1) We start with a larger pool of data; Then (2) We identify sources that are already the source of the KG and Finally (3) We remove any data points that come from these identified sources.

3.4 EMBEDDING AND RETRIEVAL TASKS

In the Embedding and Retrieval section, we aim to 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 chunks, as depicted in the system diagram.

3.4.1 EMBEDDING TECHNIQUES FOR SMALLER CHUNKS

The system 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.

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

For making the system to be more accurate and efficient, we need to filter out the irrelevant chunks. As highlighted in the system diagram, we use “Similarity Cutoff Strategy” strategy for filtering and prioritizing the information that is included. The architecture of this process has many similarities to the query scoring section.

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.

3.5. LLMS

- **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.

Note that this procedure can be disabled if the dataset is small or if the similarity cutoff is not needed.

3.5 LLMs

LLMs play a pivotal role in our FactCheck system, serving as the cornerstone for sophisticated information synthesis and response generation. As illustrated in the system 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 system 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 workload across the available models.
- **Version Management:** Maintaining and updating multiple LLMs to ensure they remain current and aligned with the latest advancements in NLP.

For running open-source LLMs, we use *Ollama*⁶, Ollama is an open-source project that simplifies the process of setting up, running, and using 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.cpp⁷ and a centralized hub of models/settings. Performance Considerations and Limitations of *Ollama*:

- 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 SYSTEM

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

- **Information Synthesis:** Integrating and coherently combining information from various sources.

⁶<https://ollama.com/>

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

3.6. VOTING SYSTEM AND CONFLICT RESOLUTION

- **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 KG.

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.

3.6 VOTING SYSTEM AND CONFLICT RESOLUTION

Our state-of-the-art system 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. We use majority voting over four models to select the final response, and in the cases where we have a tie, we use a more advanced approach to resolve the conflict.

3.6.1 MAJORITY VOTING SYSTEM

A significant aspect of our LLM integration, is the establishment of a majority voting mechanism for response creation. The majority vote decides the final response based on the outputs of multiple models.

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 system incorporates several mechanisms for resolving conflicts that may arise from divergent model outputs.

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 system diagram. This element’s role is resolving conflicts and ensuring coherence in the system’s outputs. Basically, this component is responsible for making the final decision based on the majority voting results.

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.

3.6. VOTING SYSTEM AND CONFLICT RESOLUTION

- **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 system 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 commercial 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 consistent model or the least consistent 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 system.

Algorithm 1 Resolve Ties in Majority Voting System

Require:

modelScores - A list of consistency scores for each model
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(modelScores, finalJudger, atLeast)
2:   if atLeast is True then
3:     maxScore  $\leftarrow$  maximum score in modelScores
4:     candidates  $\leftarrow$  indices with score equal to maxScore
5:   else
6:     minScore  $\leftarrow$  minimum score in modelScores
7:     candidates  $\leftarrow$  indices with score equal to minScore
8:   end if
9:   chosenIndex  $\leftarrow$  random choice from candidates
10:  return finalJudger[chosenIndex]
11: end procedure
```

3.7 PERFORMANCE REPORT

In Tables 3.3, 3.4, we provide a detailed performance report, highlighting key metrics that demonstrate the system's efficiency. Other sections of the system, such as RAG and conflict resolution mechanisms, will be evaluated in future chapters.

3.8 SYSTEM FLOW AND DECISION POINTS

The architecture of our system 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,

3.8. SYSTEM FLOW AND DECISION POINTS

Table 3.3: Performance of our LLM-based tasks in production, generated by Openlit with Gemma2 model.

Task	Avg. Request Time*	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

* The average time is calculated on the Macbook Pro with M2 max chip and 32GB of RAM.

Table 3.4: Performance of our information retrieval mechanisms.

Task	Avg. Time*
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

* The average time is calculated on the unix based server with 2 cores and 4GB of RAM.

and generate accurate responses. This section provides a comprehensive analysis of the system's flow and the key decision points that guide the processing of information.

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

Let's examine each of these stages in detail in table 3.5, focusing on the flow of information and the critical decision points within each.

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

- **Computational Efficiency:** Balancing the depth of processing at each stage with the need for timely responses.

Table 3.5: Comprehensive list of decision points in the system flow.

Component	Decision Point
<i>Input Processing and Query Generation</i>	
KG Dataset Integration	Extent and nature of KG 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.

- **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 system for accountability and debugging purposes.

3.9. ETHICAL CONSIDERATIONS AND LIMITATIONS

- **Scalability:** Ensuring that the system can handle increasing query volumes without significant degradation in performance or accuracy.

In conclusion, the system 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 system, 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.

3.9 ETHICAL CONSIDERATIONS AND LIMITATIONS

The development and deployment of advanced systems, 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 system 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 system.

3.9.2 LIMITATIONS

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

Table 3.6: Limitations of the system, categorized by scope of knowledge.

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

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 KG fact verification system.

4.1.1 FACTBENCH DATASET

FactBench is a multilingual dataset specifically designed for fact-checking in KGs¹ [11]. 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 Wikipedia² (DBpedia respectively) and Freebase³.

To obtain positive examples, the authors 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

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

²<https://www.wikipedia.org/>

³<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.

For generating incorrect facts (negative examples), the authors modified correct ones while adhering to domain and range constraints. 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. We consider a configuration where incorrect facts are a mix produced using different negative example generation strategies, resulting in a ground-truth accuracy of $\mu = 0.54$. 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, WordNet⁴, and GeoNames⁵. For our evaluation, we use YAGO2-sample⁶ [32], a subset of the full YAGO2 KG derived from AMIE horn clauses [10]. This sample consists of 1,386 beliefs spanning 16 unique predicates.

Key characteristics of the YAGO dataset in our evaluation include:

⁴<https://wordnet.princeton.edu/>

⁵<https://www.geonames.org/>

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

- High accuracy: The gold standard accuracy of the *YAGO2-sample* is 99.20%, indicating a very high-quality dataset.
- 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 which contains approximately 6.2M entities and 1.1B triplets. Following Marchesin et al.'s approach [28] to entity-oriented research, several filtering criteria were applied to ensure high-quality data for evaluation.

The analysis was restricted to subject entities that include both: 1) rdfs:label predicate and 2) rdfs:comment predicate

Additionally, they focused exclusively on A-Box triplets (assertional knowledge) while excluding T-Box triplets (terminological knowledge). The T-Box encompasses ontological entities and relationships, while A-Box contains the actual assertions that need verification. After applying these filters, their working dataset consisted 4.6M entities with 170M triplets.

From this filtered dataset, they conducted a comprehensive annotation study on 9,930 facts, which were carefully selected to represent diverse types of relationships and knowledge domains within *DBpedia*. To ensure annotation quality, Marchesin et al. implemented several measures:

- Multiple annotators per fact (minimum of three annotations per triplet)
- Expert consensus requirement for final labels
- Binary validation approach treating all incorrect facts equally regardless of error type

4.2. CANDIDATE MODELS

- Documented agreement rates between expert annotators (77% agreement with Cohen’s score of 0.51)
- Third-party resolution for 82% of initial disagreements

The dataset was carefully curated to ensure a manageable yet representative evaluation set, derived from the vast scale of *DBpedia*. By utilizing this curated subset of *DBpedia*, we benefit from a balance between the richness of a real-world KG and the practicality required for thorough empirical evaluation. For our evaluation system, we use subset of this annotated dataset, by removing facts with `<UNK>` labels, resulting in 9,344 facts.

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 KG veracity framework across diverse scenarios. Table 4.1 summarizes the key features of these datasets:

Table 4.1: Statistical summary of FactBench, YAGO, and DBpedia datasets

	FactBench	YAGO	DBpedia
Num. of Facts	2,800	1,386	9,344
Num. of Predicates	10	16	1,092
Avg. Facts per Entity	2.42	1.69	3.18
Gold Accuracy (μ)	0.54	0.99	0.85

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 KG. The *DBpedia* subset, curated specifically for entity-oriented search tasks, allows us to evaluate our framework in the context of query-dependent fact checking. 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 mod-

els⁷. They are text-to-text, decoder-only LLMs, available in English, with open weights for both pre-trained variants and instruction-tuned variants. Gemma 2 implements a similar architecture to the original Gemma model, with a few key differences. The model alternates between local sliding window attention with a 4096-token span and global attention with an 8192-token span in alternate layers. Logits are capped within a specified range to stabilize the values during attention and final layers, with `soft_cap` set to 50 for self-attention layers and 30 for the final layer. RMSNorm is used for normalization in transformer sub-layers, and Grouped-Query Attention (GQA) with two groups enhances inference speed without sacrificing performance. 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 [51]. For *Qwen2.5*, Alibaba Cloud⁸ 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 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

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

⁸https://www.alibabacloud.com/en?_p_lc=7

4.2. CANDIDATE MODELS

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 [45].

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. *Llama3.1* was pre-trained on 15 trillion tokens of data. In post-training The models produced by doing several rounds of alignment on top of the pre-trained model. Each round involves Supervised Fine-Tuning (SFT), Rejection Sampling (RS), and Direct Preference Optimization (DPO). Meta use synthetic data generation to produce the vast majority of our SFT examples, iterating multiple times to produce higher and higher quality synthetic data across all capabilities. Additionally, they invest in multiple data processing techniques to filter this synthetic data to the highest quality. This enables model to scale the amount of fine-tuning data across capabilities. Compared to previous versions of Llama, developers improved both the quantity and quality of the data we use for pre and post-training. These improvements include the development of more careful pre-processing and curation pipelines for pre-training data, the development of more rigorous quality assurance, and filtering approaches for post-training data [9, 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 AI⁹, 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 [17, 29].

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

Candidate Model Summary: The selection of candidate models for our system was guided by the need for diversity, efficiency, and reliability in processing fact verification tasks within KGs. We chose *Gemma2*, *Qwen2.5*, *Llama3.1*, and *Mistral* for their specific strengths in handling diverse linguistic queries, reasoning capabilities, and compatibility with RAG pipelines. Each of these models brings unique advantages to our verification framework, as summarized in Table 4.2.

Table 4.2: Summary of key strengths of selected candidate LLMs for KG fact verification.

Model	Key Strengths	Description
Gemma2	Dense Retrieval & Query Processing	Optimized for dense retrieval tasks, Gemma2 processes complex linguistic structures, making it well-suited for entity-rich query generation and document ranking.
Qwen2.5	Logical Reasoning & Prompt Efficiency	Excels in reasoning tasks with minimal prompting. Its accuracy in logical inference supports consistent veracity assessments, especially for ambiguous or conflicting evidence.
Llama3.1	Efficiency & Versatility	Offers a balance of efficiency and accuracy, with robust performance across fact-checking benchmarks. Llama3.1's lower computational demands ensure responsive processing without compromising output quality.
Mistral	Context Sensitivity & Interpretability	Known for nuanced, context-driven outputs and interpretability. Mistral's language generation capabilities provide clear, human-like explanations, making it ideal for understanding the facts like a human.

For model selection, we can choose either instruction-tuned and quantized

⁹<https://mistral.ai/>

4.3. EXPERIMENTAL SETUP

models with similar architectures or with different architectures. Here, we opted for models with varied architectures to make the ensemble more versatile and capable of handling a wide range of query scenarios. Using diverse models in the ensemble offers a balanced approach to complex fact verification tasks across KGs. This multi-model setup enhances adaptability and reliability, allowing the system to respond accurately to diverse verification scenarios, even when they differ in nature.

4.3 EXPERIMENTAL SETUP

4.3.1 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:

- **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. **Here by accuracy, we mean the predication performance of our system that measured against gold standard labels.**

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). Fully reported in ablation study in chapter 5 and performance report in section 3.7.
- **Cost Efficiency:** Evaluating the cost-effectiveness of the pipeline in terms of computational resources and infrastructure by reporting the average token used per query.

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 Calculate Model Consistency Per Model

```

1: procedure MODELSTABILITYCAL(models)                                ▷ Containing binary results
2:   m_len ← length(models), stabilityScores ← []
3:   for i ← 0 to m_len – 1 do
4:     m1 ← models[i], mStabilities ← []
5:     for j ← 0 to m_len – 1 do
6:       if i ≠ j then
7:         m2 ← models[j], matchCount ← 0
8:         totPreds ← length(m1)
9:         for k ← 0 to totalPredictions – 1 do
10:          if m1[k] = m2[k] then
11:            matchCount ← matchCount + 1
12:          end if
13:        end for
14:        mStabilities ←  $\frac{\text{matchCount}}{\text{totPreds}}$                                 ▷ Append the stability score
15:      end if
16:    end for
17:    stabilityScores[i] ← mean(mStabilities)
18:  end for
19:  return stabilityScores                                              ▷ Dictionary with model stability scores
20: end procedure

```

4.3. EXPERIMENTAL SETUP

4.3.2 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 chapter 5. Table 4.3 summarizes the key system configurations used in our empirical evaluation.

Table 4.3: System configurations for empirical evaluation

Section	Parameter	Considerations
Human Understandable Text	Gemma2:9b	Other LLMs can be used, but using instruction-tuned models is recommended. This is skipped for <i>Fact-Bench</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 RelevanceThreshold	0.5	–
Num. of Selected Questions	3	–
Google Search	–	Used query params: <i>lr</i> = 'lang_en', <i>gl</i> = 'us', <i>hl</i> = 'en', <i>num</i> = '100'. The <i>lr</i> parameter is set to the language of the query, <i>gl</i> to the country, <i>hl</i> to the language, and <i>num</i> to the number of results.
Num. of Selected Documents	10	–
Document Selection	ms-marco-MiniLM-L-6-v2	Filtered out the documents from these origins: dbpedia, wikipedia, wikimedia, wikidata, quora, britannica, scholarpedia, newworldencyclopedia, everipedia, encyclopedia, wikibooks, wiktionary, wikiversity, wikisource, wikiquote, wikivoyage, academia, and nytimes
Embedding Model	bge-small-en-v1.5	–
Chunking Strategy	Sliding Window 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	-	Use model with higher-param for each model, for llama3.1:8b → 70b, gemma2:9b → 27b, qwen2.5:7b → 14b, and mistral:7b → mistral nemo:12b.
--------------	---	--

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

4.4 COMPARATIVE ANALYSIS

4.4.1 DISCUSSION OF RESULTS

Evaluation across three distinct datasets - *FactBench*, *YAGO*, and *DBpedia* - shows significant insights into the performance and efficiency of different language models for KG fact verification.

Based on Table 4.4, in the *FactBench* dataset, *Gemma2* emerged as the strongest individual performer, achieving an accuracy of 0.9014 and an F1 score of 0.9085. These results were further enhanced by the ensemble approach using *Qwen2.5:14b*, which improved the accuracy to 0.9057 and F1 score to 0.9145. On the *YAGO* dataset, *Mistral* demonstrated exceptional performance, reaching an accuracy of 0.9221 and an F1 score of 0.9594. The consistent high performance across models on this dataset suggests that well-structured, high-quality KGs can be

4.4. COMPARATIVE ANALYSIS

Table 4.4: Empirical evaluation results of the proposed system and candidate LLMs over the FactBench, YAGO, and DBpedia.

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	Consistency*	Avg. request time	Avg. input tokens	Avg. output tokens		Acc	F1
FactBench	Gemma2	0.8738	2.32s	1509.30	19.95		0.9014	0.9085
	Qwen2.5	0.8747	2.46s	1509.18	67.73		0.8746	0.8910
	LLama3.1	0.8296	2.87s	1509.24	104.65		0.8243	0.8378
	Mistral	0.8686	1.73s	1509.17	8.81		0.8507	0.8729
	Most (Qwen2.5:14b)	–	27.66s	1509.18	57.44		0.9057	0.9145
	Least (Llama3.1:70b)	–	12.70s	1749.14	8.73		0.9025	0.9124
YAGO	Gemma2	0.8882	2.15s	1508.83	17.06		0.8506	0.9191
	Qwen2.5	0.8884	2.50s	1509.35	72.87		0.8600	0.9246
	LLama3.1	0.8552	7.20s	1509.25	104.58		0.8333	0.9089
	Mistral	0.8920	1.67s	1509.31	8.03		0.9221	0.9594
	Most (Mistral-nemo:12b)	–	2.31s	1560.97	9.13		0.8701	0.9304
	Least (Llama3.1:70b)	–	11.60s	1560.97	10.19		0.8853	0.9391
DBpedia	Gemma2	0.8207	7.80s	1551.69	27.15		0.6821	0.7865
	Qwen2.5	0.8247	2.63s	1552.02	73.41		0.7236	0.8211
	LLama3.1	0.7573	3.06s	1551.98	110.82		0.6224	0.7377
	Mistral	0.8162	1.81s	1551.95	8.76		0.7201	0.8192
	Most (Qwen2.5:14b)	–	6.48s	1695.92	84.47		0.7014	0.8020
	Least (Llama3.1:70b)	–	12.82s	1701.77	8.99		0.7099	0.8089

* Consistency score for each individual model is calculated across all models, excluding the ensemble models.

effectively verified using our approach. The ensemble method with *Mistral-nemo:12b* maintained strong results while providing additional verification confidence in ambiguous cases. The *DBpedia* dataset proved more challenging, with overall lower performance across all models. *Qwen2.5* achieved the best individual results with an accuracy of 0.7236 and an F1 score of 0.8211. This performance difference highlights the impact of dataset complexity and structure on verification accuracy.

Table 4.5: Statistical analysis of output tokens and request times per query across FactBench, YAGO, and DBpedia datasets for each used model.

	Gemma2		Qwen2.5		Llama3.1		Mistral		Mistral-nemo		Qwen2.5:14b		Llama3.1:70b	
	Avg	Std	Avg	Std	Avg	Std	Avg	Std	Avg	Std	Avg	Std	Avg	Std
Output tokens	24.98	24.91	72.32	29.78	108.11	87.26	9.85	21.12	9.13	12.83	73.64	60.80	10.06	13.94
Request time	6.11s	51.38	2.59s	0.68	3.44s	19.59	1.79s	0.55	2.31	0.77	11.18s	68.12	13.94s	6.48

In terms of computational efficiency as showed in Tables 4.4, and 4.5 , Mistral demonstrated remarkable performance, generating minimal output tokens (average 8.81-9.85) while maintaining competitive accuracy. This efficiency is especially remarkable when compared to *Llama3.1*, which generated considerably more tokens (average 104.58-110.82) without any significant improvement in accuracy. Additionally, it suggests that *Llama3.1* often failed to follow instructions closely, opting to reason through every fact rather than verifying correctness.

Input token counts remained relatively consistent across models, ranging from approximately 1500 to 1700 tokens.

The processing time analysis reveals that *Mistral* consistently achieved the fastest request times (1.73-1.81s), while ensemble methods required longer processing times. However, it's important to note that ensemble methods were only employed for tie-breaking scenarios, affecting approximately 5–10% of the total queries. This selective application of ensemble methods effectively balances the trade-off between computational cost and accuracy improvement.

Table 4.6: Statistical analysis of request time per query across FactBench, YAGO, and DBpedia datasets.

	FactBench		YAGO		DBpedia	
	Avg	Std	Avg	Std	Avg	Std
Request time*	4.59s	20.17	4.55s	33.03	6.32s	39.04

* Time taken for the embedding phase and LLM request.

The analysis from Tables 4.4 and 4.6 show notable patterns in the performance specific to each dataset. Models generally achieved better results on more structured datasets like *FactBench* and *YAGO* compared to *DBpedia*. This pattern suggests that the clarity and consistency of the underlying KG significantly influence verification accuracy. Request times also varied across datasets, with *DBpedia* queries requiring longer processing times (6.32s) compared to *FactBench* (4.59s) and *YAGO* (4.55s), likely due to its greater complexity and size.

The ensemble approach proved particularly effective in resolving ambiguous cases. While the computational cost of ensemble methods is higher, their selective application only to uncertain cases (5-10% of queries) makes this trade-off acceptable in practice. The high consistency scores observed in ensemble methods (>0.91) suggest more reliable predictions for challenging cases, justifying the additional computational investment for these specific instances.

4.4.2 QUALITATIVE ERROR ANALYSIS

As depicted in Figure 4.1, our objective is to categorize errors by LLMs and text embedding model. This approach helps reveal common error types and patterns by clustering explanations into distinct groups. The process begins by gathering logs of incorrectly labeled data, referred to here as "wrongly labeled facts." These logs are analyzed, and we then prompt an LLM to generate explanations, or "reasons," for each error. This step provides context and may

4.4. COMPARATIVE ANALYSIS

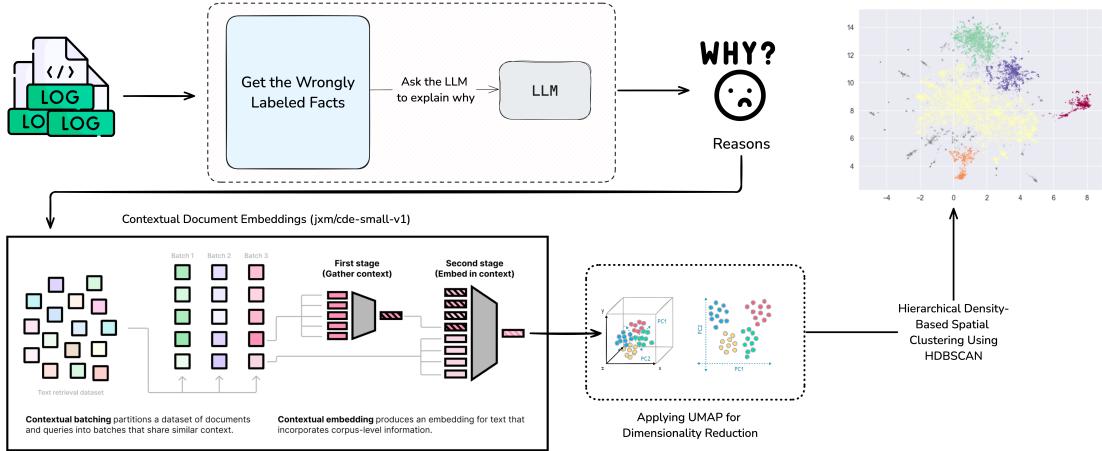


Figure 4.1: Collecting logs and leveraging LLM-generated reasoning, combined with contextual document embeddings (jxm/cde-small-v1) [30], to cluster errors using a hierarchical density-based spatial technique.

highlight underlying patterns or causes that contribute to these errors. The prompt template used for this reasoning process is detailed in Appendix A.4. After obtaining explanations from the LLM, we use a specialized text embedding model named "jxm/cde-small-v1"¹⁰. We selected cde-small-v1 because it is the highest-ranked small model (under 400 million parameters) on the MTEB leaderboard for text embedding models, as of October 1, 2024. This model transforms each explanation into a contextualized embedding, capturing both semantic meaning and specific instruction-driven nuances for each error's context [30]. Next, these embeddings undergo dimensionality reduction using Uniform Manifold Approximation and Projection (UMAP), a technique that projects high-dimensional data into two or three dimensions while preserving local and some global structure. UMAP's visualization helps identify potential clusters or groupings of similar errors, making it easier to observe patterns that might be difficult to see in higher dimensions. Once reduced in dimensionality, the embeddings are fed into Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN), a clustering algorithm well-suited for discovering clusters in data with varying density. HDBSCAN clusters the error embeddings based on their density, identifying groups of similar errors and isolating outliers. Following clustering, we identify some reasons from each dataset. These representative reasons are then provided to an LLM to assign

¹⁰<https://huggingface.co/jxm/cde-small-v1>

descriptive labels to each error category, which encapsulate the main types of errors across the dataset.

The labeling data for each cluster is as follows:

- **UnLabeled:** The information or context provided does not contain the claimed details, such as references to specific individuals, places, or events that are purportedly associated with the topic.
- **Relationship Errors:** Errors arise from misstatements regarding relationships between people, such as marital status or religious affiliations that conflict with the provided details.
- **Role Attribution Errors:** Errors are due to incorrect associations of individuals with particular roles, places, or teams that do not match the details in the context.
- **Geographic/Nationality Errors:** This category includes errors related to locations, national affiliations, or settings that do not align with the context or provide contradictory information.
- **Genre/Classification Errors:** Misclassifications of films, genres, or roles are highlighted here, especially when certain works are wrongly associated with people, studios, or genres.
- **Identifier/Biographical Errors:** These errors involve incorrect identifiers or biographical details, such as award titles, label names, or authorship that don't match the context.

The heatmaps in Figure 4.2 visualize the overlap in error patterns between different models across error categories and datasets. The overlap matrices reveal distinct patterns of agreement and disagreement between models when making errors, with darker colors indicating higher overlap percentages. For *DBpedia*, we observe moderate to high overlap (45-75%) between models across most error categories, suggesting similar challenges in handling complex factual relationships. The *FactBench* dataset shows lower overlap percentages (30-40% typical), indicating more independent error patterns between models. YAGO exhibits variable overlap, with particularly high agreement in unlabeled errors (75-85% overlap) but lower overlap in relationship errors (23-35%). These patterns suggest that while models often struggle with similar types of facts, they

4.4. COMPARATIVE ANALYSIS

Table 4.7: Dataset-wise error clustering based on LLM-generated reasoning, using Contextual Document Embeddings for embeddings, UMAP, and HDBSCAN.

Dataset	Model	UnLabeled	Relationship	Role Errors	Geo Errors	Classification	Identifiers	Total*
FactBench	Gemma2	4	36	45	176	13	1	275
	Qwen2.5	33	27	60	194	34	1	349
	Llama3.1	38	44	73	295	38	3	491
	Mistral	53	27	53	242	40	2	417
	Unique. Ratio (%)	0.62	0.72	0.44	0.52	0.63	0.57	0.53
YAGO	Gemma2	6	134	0	14	51	2	207
	Qwen2.5	7	109	0	13	63	2	194
	Llama3.1	8	98	0	19	104	2	231
	Mistral	7	54	0	10	34	3	108
	Unique. Ratio (%)	0.35	0.52	—	0.46	0.51	0.33	0.50
DBpedia	Gemma2	353	22	98	1729	459	299	2960
	Qwen2.5	339	19	91	1525	357	237	2568
	Llama3.1	382	28	109	2172	509	318	3518
	Mistral	325	20	94	1487	438	241	2605
	Unique. Ratio (%)	0.41	0.43	0.44	0.42	0.42	0.40	0.41

* Some errors may not be included in this analysis because we did not receive any responses for them. While we classify these as incorrect predictions, they are not considered in the error analysis section.

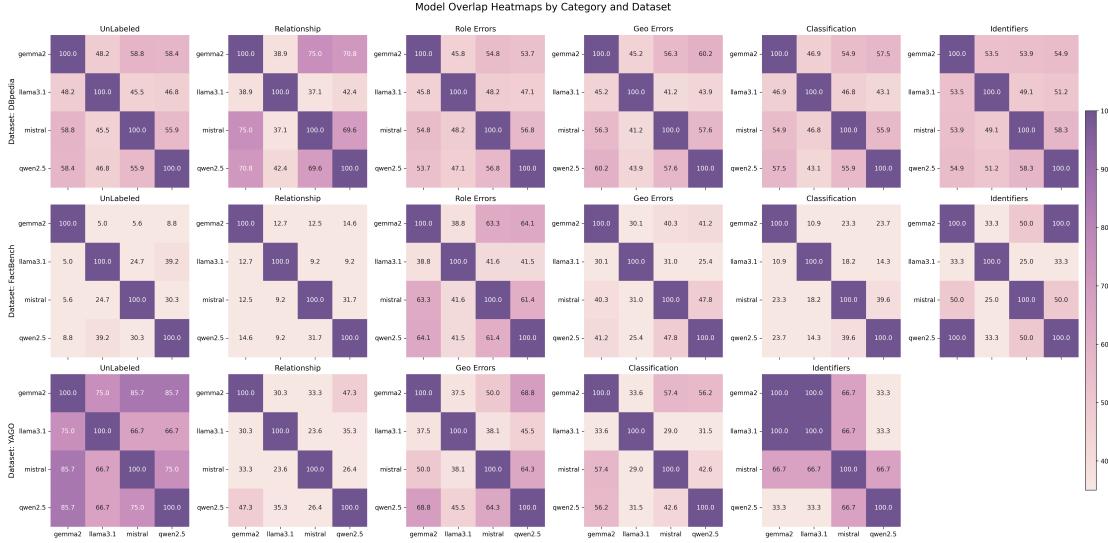


Figure 4.2: Model overlap heatmaps by category and dataset. Each cell shows the percentage overlap in errors between model pairs. Matrices are organized by error category (UnLabeled, Relationship, Role Errors, etc.) and dataset (DBpedia, FactBench, YAGO), revealing patterns in how models agree or disagree when making verification errors.

also make distinct errors, supporting the value of ensemble approaches. The lower overlap in *FactBench* errors particularly validates our multi-model verification strategy.

Results from Table 4.7 alongside Figures 4.3, 4.4, and 4.5 highlight both common challenges and model-specific characteristics in fact validation performance.

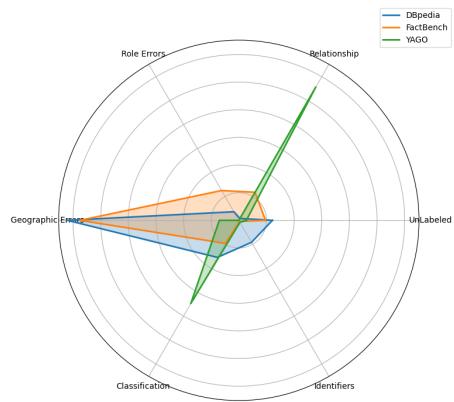


Figure 4.3: Normalized distribution of error clusters across datasets.

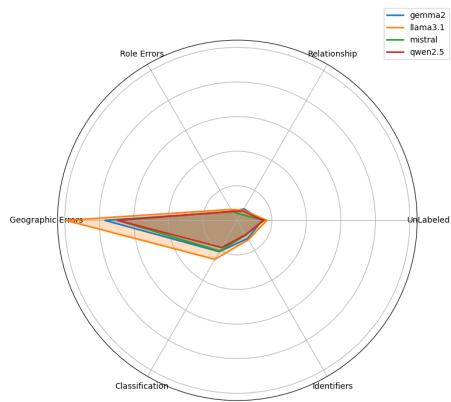


Figure 4.4: Distribution of error clusters across selected LLMs.

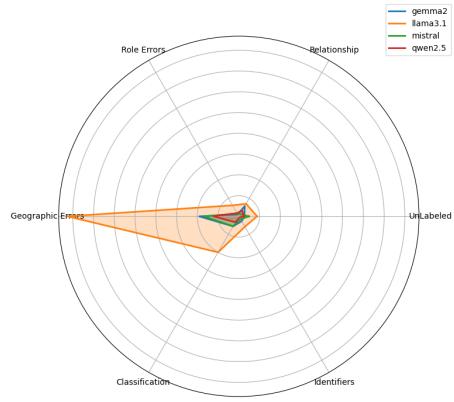


Figure 4.5: Distribution of tendency to be wrong across gemma2, qwen2.5, LLama3.1 and mistral models. The right chart illustrates the distribution of fully incorrect predictions (4/4) detailing the instances where all predictions made by the models were incorrect. The left chart depicts the distribution of just one wrong prediction (1/4).

Geographic and nationality-related errors emerged as the predominant challenge, accounting for 56.9% of total errors across all models and datasets. This pattern was particularly pronounced in the *DBpedia* dataset, where geographic errors constituted 58.5% of all errors, suggesting a systematic challenge in processing and validating location-based information. This pervasive difficulty across all models indicates a fundamental challenge in handling geographic relationships and facts.

The analysis of dataset-specific patterns revealed distinct characteristics and challenges. The *DBpedia* dataset proved to be the most challenging, generating

4.4. COMPARATIVE ANALYSIS

the highest error count and showing particular vulnerability to geographic and classification errors. In contrast, the *FactBench* dataset demonstrated a more balanced distribution of errors across categories, though still showing a predominance of geographic errors. The *YAGO* dataset exhibited a unique pattern, with relationship errors being the most frequent, followed by classification errors, and notably showing no role attribution errors a distinctive characteristic that sets it apart from other datasets.

When examining model-specific performance, *LLama3.1* consistently generated the highest error counts across datasets, showing particular vulnerability to geographic errors in the *DBpedia* dataset and elevated classification errors compared to other models. *Mistral*, on the other hand, demonstrated stronger overall performance, particularly in the *YAGO* dataset. *Gemma2* and *Qwen2.5* showed similar error patterns and counts, positioning themselves between *LLama3.1* and *Mistral* in terms of performance.

The hierarchical distribution of error types shows a consistent pattern across models. This consistency in error distribution suggests that these challenges are inherent to the task rather than model-specific limitations.

Despite varying error counts, models maintained similar error distribution patterns within each dataset, indicating that these challenges are systematic rather than model-specific. The analysis suggests several critical areas for future development in LLM fact validation capabilities. Primary attention should be directed toward enhancing geographic and location-based reasoning capabilities, given their dominant role in error generation. Additionally, improving classification tasks and the handling of insufficient context or ambiguous information could significantly enhance overall performance. The consistent patterns across models suggest that these improvements would benefit the field broadly rather than being model-specific enhancements.

4.4.3 DBPEDIA ANALYSIS IN DEPTH

Since *DBpedia* proved to be the most challenging dataset in our evaluation, we performed a deeper analysis based on this dataset. Here we focus on the partition-wise evaluation and topic-wise analysis to understand the impact of fact popularity on verification performance.

Partition-wise: We performed deeper analysis on stratification established by Marchesin et al. [28]. Each KG triple was assigned to one of seven partitions, numbered 1–7, where partition 1 represents the least popular/common knowledge and partition 7 represents the most popular/common knowledge. This stratification helps us understand how our verification system performs across different levels of fact popularity and complexity within *DBpedia*. The analysis can reveal whether the system’s accuracy varies between common, well-documented facts versus more obscure or specialized knowledge. This insight is valuable for identifying areas where the system needs improvement and understanding its real-world applicability across different types of knowledge.

Table 4.8: Partition-wise evaluation results of the proposed system and candidate LLMs over the DBpedia dataset.

Weight	Size	Gemma2		Qwen2.5		Llama3.1		Mistral		Qwen2.5:14b		Llama3.1:70b		
		Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	
Stratum 1	0.9120	2223	0.672	0.773	0.708	0.804	0.600	0.712	0.710	0.806	0.688	0.786	0.700	0.797
Stratum 2	0.0616	1695	0.697	0.796	0.720	0.816	0.627	0.738	0.726	0.820	0.711	0.807	0.720	0.814
Stratum 3	0.0177	1588	0.689	0.791	0.736	0.828	0.632	0.743	0.719	0.819	0.708	0.806	0.719	0.815
Stratum 4	0.0044	1327	0.689	0.797	0.737	0.835	0.628	0.747	0.709	0.815	0.706	0.811	0.714	0.817
Stratum 5	0.0029	1058	0.666	0.780	0.705	0.813	0.638	0.758	0.724	0.826	0.690	0.800	0.695	0.804
Stratum 6	0.0010	814	0.692	0.800	0.719	0.822	0.614	0.739	0.733	0.833	0.709	0.813	0.708	0.812
Stratum 7	0.0001	629	0.671	0.776	0.779	0.864	0.650	0.762	0.754	0.848	0.731	0.826	0.747	0.838

The weights declared in Table 4.8 are gathered by Marchesin et al. [28], stratum 1 contains the vast majority of the triplets in DBpedia. This makes sense because stratum 1 represents the lowest-utility facts - those that are rarely used in queries.

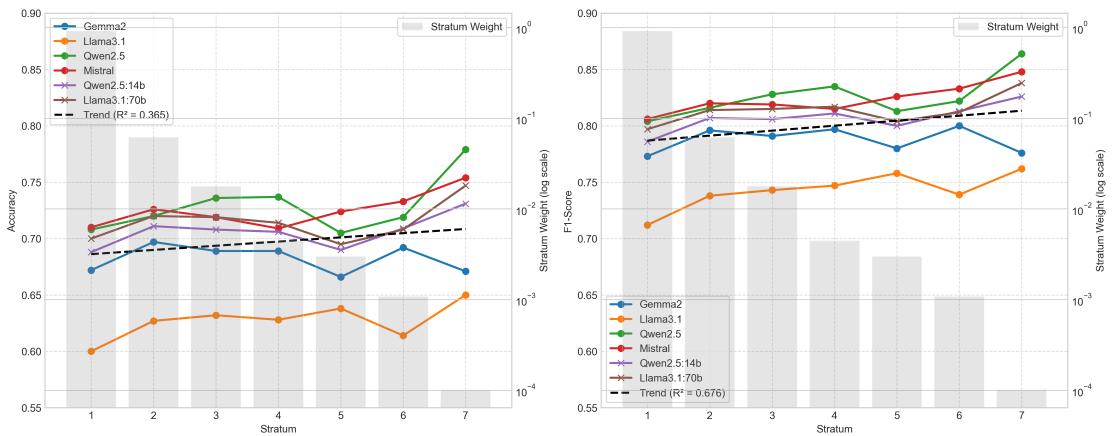


Figure 4.6: Partition-wise model performance comparison: Accuracy and F1-scores for KG fact verification on DBpedia dataset. Gray bars indicate stratum weights (log scale).

4.4. COMPARATIVE ANALYSIS

Table 4.8 and Figure 4.6, demonstrate that the models have a clear trend of improved performance on more common knowledge. *Qwen2.5* exhibits particularly strong performance, with accuracy increasing from 0.708 in Stratum 1 to 0.779 in Stratum 7, and F1-scores following a similar upward trajectory. This suggests that the model benefits from the richer context and more consistent representation of popular facts in the knowledge base.

The performance disparity between lower and higher strata highlights a common challenge in KG verification: the system's reliability varies with fact popularity. This insight is particularly valuable for real-world applications, where handling both common and specialized knowledge is crucial.

Building on our previous analysis of *DBpedia* results, we conducted a stratum-wise error analysis to better understand how error distributions vary across different data partitions. Our analysis of error patterns across knowledge strata showed in Figure 4.7 declares distinct performance characteristics among the four LLMs. The models demonstrated varying levels of effectiveness in handling knowledge from different popularity strata, with error rates showing notable patterns across the commonality spectrum.

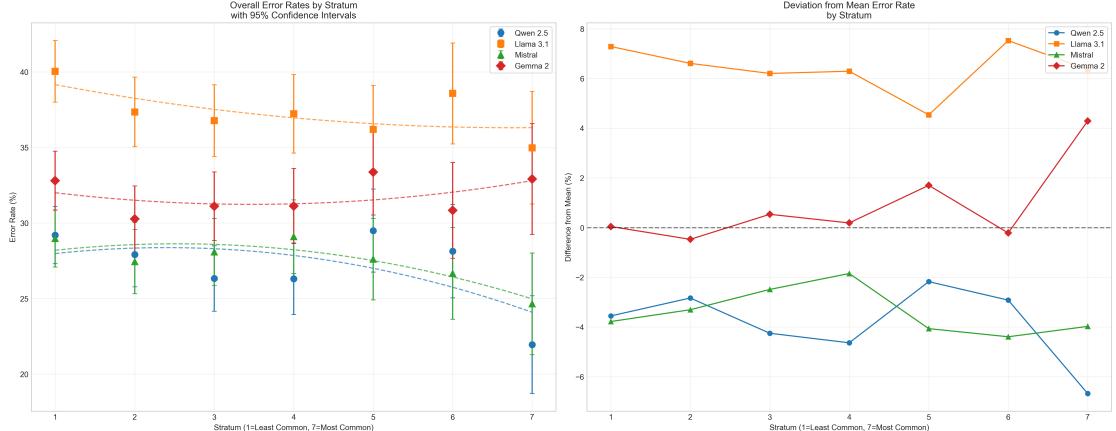


Figure 4.7: Error distribution analysis across different language models and frequency strata.

Llama3.1 exhibited the highest overall error rate ($37.31\% \pm 1.51\%$), significantly exceeding other models' error rates. Despite its higher error rate, *Llama3.1* maintained relatively consistent performance across strata (range: 34.98% - 40.04%), suggesting uniform handling of both common and rare knowledge. The model showed a slight negative correlation with stratum number ($r=-0.628$, $p=0.1309$), indicating a modest tendency to perform better with more common

knowledge, though this trend was not statistically significant.

In contrast, *Qwen2.5* and *Mistral* demonstrated notably lower error rates ($27.04\% \pm 2.38\%$ and $27.50\% \pm 1.41\%$ respectively), with *Mistral* showing the most pronounced negative correlation with stratum number ($r=-0.761$, $p=0.0470$). This statistically significant correlation indicates that *Mistral*'s performance improves substantially as knowledge becomes more common. *Qwen2.5* showed the widest range of error rates (21.94% - 29.49%), suggesting more variable performance across different knowledge types.

Gemma2 maintained an intermediate position with a mean error rate of $31.77\% \pm 1.13\%$ and showed the most stable performance across strata (range: 30.27% - 33.36%). Uniquely among the models, *Gemma2* exhibited a slight positive correlation with stratum number ($r=0.237$, $p=0.6083$), though this trend was not statistically significant. In general, it has the most uniform error distribution across strata.

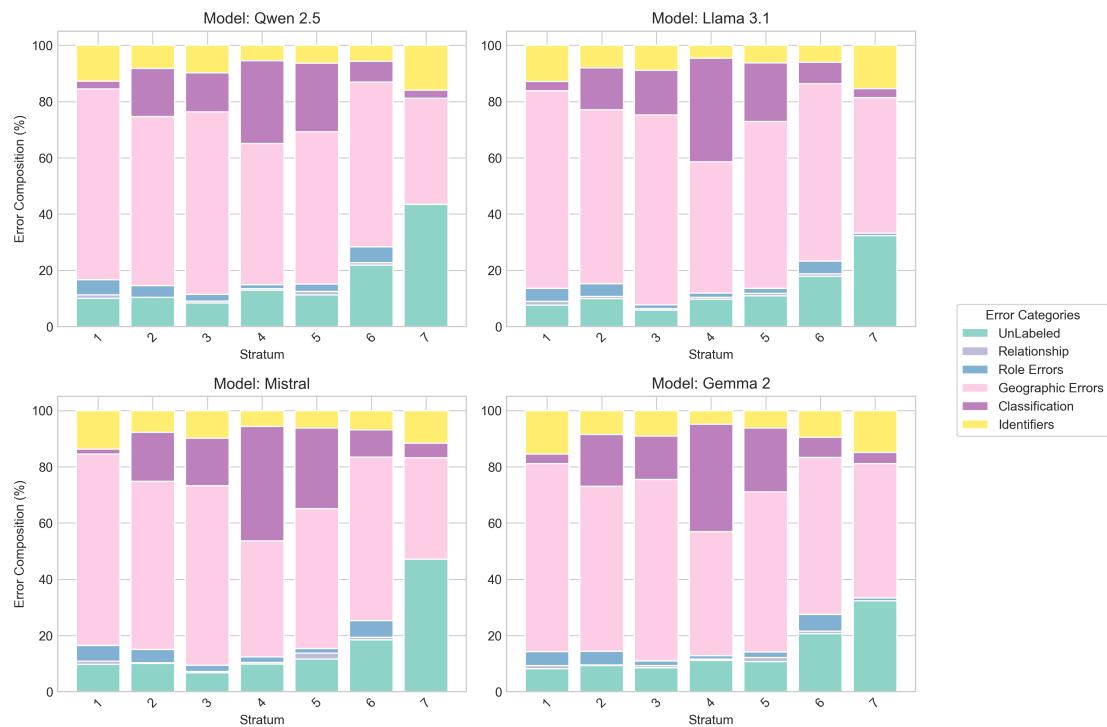


Figure 4.8: distribution of error categories across different language models and frequency strata

Analysis of error categories in Figure 4.8 reported separate patterns across strata, with certain error types becoming more prevalent in specific knowledge domains. The proportion of unLabeled errors increased notably in the most

4.4. COMPARATIVE ANALYSIS

common knowledge strata (6-7), while geographic errors showed higher prevalence in less common knowledge strata (1-3). Classification errors maintained relatively consistent proportions across all strata, suggesting that this type of error is less influenced by knowledge commonality.

These findings suggest that while newer models like *Qwen2.5* achieve lower overall error rates, they may be more sensitive to knowledge popularity, performing notably better with common knowledge. In contrast, models like *Gemma2* offer more consistent performance across knowledge types, potentially making them more reliable for applications requiring uniform handling of both common and rare knowledge.

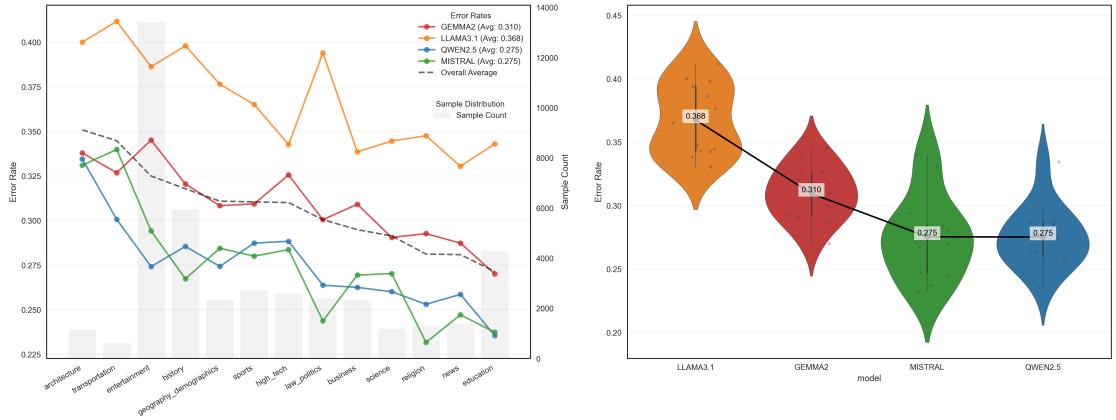


Figure 4.9: Comparative analysis of language model error rates across knowledge domains. (Left) domain-specific error rates across 13 knowledge categories, with overlaid sample distribution bars. (Right) distribution of error rates for each model

Topic-wise: By employing BERTopic from the work of Marchesin et al. [28] to generate interpretable clusters, we obtained 64 distinct clusters, which were manually scrutinized and aggregated to form a set of 13 broad topics. These topics encompass “Architecture”, “Business”, “Education”, “Entertainment”, “Geography”, “High Tech”, “History”, “Law and Politics”, “News”, “Religion”, “Science”, “Sports”, and “Transportation”. Some facts were not assigned to any topic or were assigned to multiple topics.

As showed in Figure 4.9 “Education” and “News” domains consistently show lower error rates across all models. “Architecture” and “Transportation” categories present higher error rates, particularly for *Llama3.1*. High variance is observed in the “Business” and “Law and Politics” domains, suggesting these may be more challenging areas for current models. Also we figured out that

Qwen2.5 and *Mistral* have compact distributions, indicating consistent performance, but *Llama3.1* exhibits the broadest distribution, suggesting more variable performance across domains

Correlation analysis reveals varying relationships between sample count and error rates. *Gemma2* shows a strong positive correlation (0.410), suggesting lower reliability in less-sampled domains. *Llama3.1* exhibits a weak positive correlation (0.204). *Mistral* shows negligible correlation (0.006), indicating consistent performance regardless of sample size. *Qwen2.5* demonstrates a slight negative correlation (-0.101), hinting at better performance in less-sampled domains. This results suggest that *Mistral* may be the most reliable model across different knowledge domains.

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 four 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. For our ablation study, we focus specifically on the *FactBench* dataset.

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.
3. **Performance Evaluation:** For each configuration, we assess the system's performance using predefined metrics (detailed in Sections 5.1.2 and 4.3.1).
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 4.3.1:

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

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

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 Time:** Measured in seconds per query to assess computational efficiency, measured on the Macbook Pro with M2 Max chip and 32GB of RAM.

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

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

5.2. DOCUMENT SELECTION

information need expressed in the query. We consider both unsupervised and supervised approaches. Using these methods, we aim to find the most relevant documents from the data pool we have collected through web scraping 3.3.3.

5.2.1 UNSUPERVISED METHODS

BM25: BM25 [37] 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

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.[16] that leverages contrastive learning to train dense retrieval models. Rather than relying on term matching, Contriever learns to map semanti-

cally 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.¹ 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

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

5.2. DOCUMENT SELECTION

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 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*². 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 [27] spanning 26 languages, while also exhibiting superior NDCG@10 scores on English-language tasks in the BEIR benchmark [46]. Notably, it secured the top position on the AirBench leaderboard upon its release³.

These capabilities make the model particularly valuable for multilingual information retrieval, agentic 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 trained on the MS MARCO dataset but distilled to a smaller size for efficiency. This comparison in Table 5.1

Table 5.1: Performance comparison of various distilled MS MARCO models based on BERT architecture, measured across NDCG@10 on TREC DL 2019 and MRR@10 on MS MARCO Dev benchmarks.

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

highlights the trade-offs between model size, retrieval effectiveness, and pro-

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

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

cessing speed (documents per second). As model size increases from *TinyBERT-L-2-v2* to *MiniLM-L-12-v2*, there is a noticeable improvement in retrieval metrics (NDCG@10 and MRR@10), indicating higher relevance in retrieved documents. However, this comes at the cost of reduced inference speed, with larger models processing fewer documents per second. For our project, we use *ms-marco-MiniLM-L-6-v2*⁴ Cross-Encoder model. This model, trained on the extensive MS MARCO dataset comprising approximately 500,000 authentic search queries from the Bing search engine [35], demonstrates superior performance within a two-stage Retrieve & Re-rank framework. In this paradigm, an initial retrieval phase employs either lexical search methods or 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

Figure 5.1 visualizes similarity patterns in document selection across models, using Jaccard Similarity to quantify the overlap in documents identified as relevant by different models. Higher values indicate greater agreement between models in identifying similar documents, providing insights into consistency and variations in retrieval behavior across the models under evaluation. We

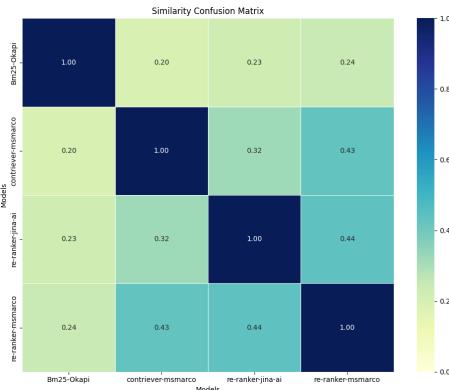


Figure 5.1: Document retrieval confusion matrix based on Jaccard similarity between documents retrieved by each model.

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

5.2. DOCUMENT SELECTION

can figure 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, 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), highlighting more consistent retrieval patterns. However, they still exhibit differences in document selection. In general, we can find out with different methods we have different result over the same query. 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.

Table 5.2: Performance evaluation of various document retrieval methods on the Fact-Bench dataset, using the Gemma2 model.

Retrieval Method			
Method	Acc	F1	Latency*
<i>Unsupervised</i>			
Bm25	0.8882	0.8940	0.4614s
contriever-msmarco	0.8932	0.8988	25.903s
<i>supervised</i>			
jina-reranker-v2-base-multilingual	0.9004	0.9065	9.8958s
ms-marco-MiniLM-L-6-v2	0.9014	0.9077	0.8172s

* It is measured in seconds on average per query.

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 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 accuracy and F1 score relative to the other models under consideration with acceptable latency.

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) [41] 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) [53].

5.3. EMBEDDING MODELS

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 [24]. This series of models has demonstrated superior performance in multiple benchmarks, including the Massive Text Embedding Benchmark (MTEB) [31] and the LoCo long-context retrieval benchmark [38]. In particular, the tuned models of this 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 various 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 on 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 [41, 40].

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) [15] 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 wide range of 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 OpenAI⁵ and Cohere⁶ on English tasks and achieves high scores on multilingual benchmarks.

Jina-embeddings-v3 also features flexible Matryoshka Representation Learning (MRL) [19], allowing users to reduce the embedding size from 1024 to as low as 16 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*⁷ 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 [19], 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 em-

⁵<https://openai.com/>

⁶<https://cohere.com/>

⁷https://huggingface.co/dunzhang/stella_en_1.5B_v5

5.3. EMBEDDING MODELS

beddings. 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 [49]. 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 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 InfoNCE contrastive loss with only in-batch negatives, while other hyperparameters remain consistent with the English E5 models.
- **Fine-tuning:** Following pre-training, the model was fine-tuned using high-quality labeled datasets from the E5-mistral paper [48]. 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 [54] 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 a little degradation in performance for low-resource languages, a common limitation of multilingual models, but

still outperforms many other models in this category. The use of contrastive learning and instruction tuning enables the model to generate highly effective embeddings for information retrieval 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-a-embeddings-v3. It is a compact English-specific model with just 33.4M parameters, making it highly efficient for deployment in resource-constrained environments. The model architecture is BERT-like which goes through three-stage of training. *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 model fine-tuned using a process of contrastive learning, where sentences are embedded to prioritize semantic similarity. This technique enhances retrieval tasks by training the model to produce high similarity scores for semantically related sentences while keeping unrelated pairs distant in embedding space. The fine-tuning emphasizes retrieval for short queries to long passages, optimized with a contrastive loss function and often utilizes mined hard negatives to improve differentiation between similar and unrelated sentence pairs.

Despite its small size, *bge-small-en-v1.5* punches above its weight on several English benchmarks and outperforms the base-sized BERT and RoBERTa models on most tasks while being more compact. 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.

5.3. EMBEDDING MODELS

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

Table 5.3: Comparison of characteristics of embedding models

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

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.

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, we should carefully evaluate our 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.

Table 5.4: Performance evaluation of various embedding models on the FactBench dataset, using the Gemma2 model.

Model	ms-marco-MiniLM-L-6-v2			<i>Embedding Model</i>
	Acc	F1	Latency	
stella_en_1.5B_v5	0.8961	0.9028	17.692s	
multilingual-e5-large-instruct	0.8954	0.9018	5.0038s	
bge-small-en-v1.5	0.9014	0.9077	1.6958s	
jina-embeddings-v3*	0.8852	0.9097	4.8745s	
gte-large-en-v1.5*	0.8971	0.9174	5.8571s	

* 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.

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 F1 and accuracy metrics. The low latency of 1.6958 seconds per query also makes it an attractive choice for real-time applications. 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 section 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.

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

```

Code 5.1: Small2Big chunking method

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.

5.4. CHUNKING STRATEGIES

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 use the *SentenceWindowNodeParser* to parse documents into single sentences per node. Each node also contains a "window" with the

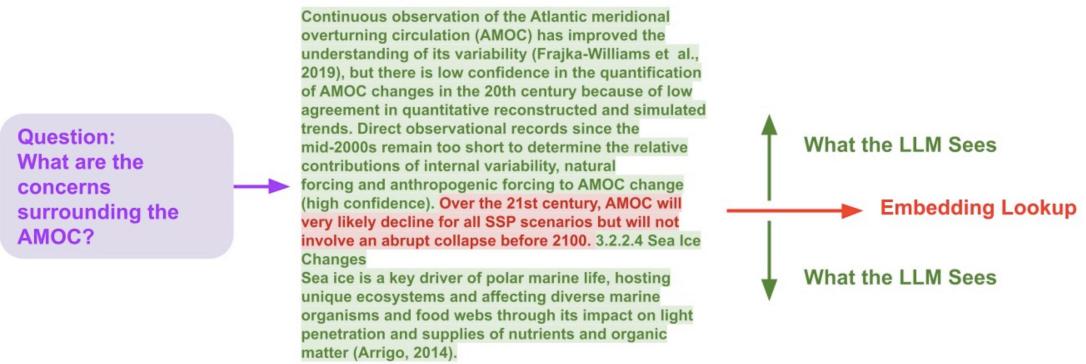


Figure 5.2: Node sentence window replacement technique as described by Liu [25].

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 *MetadataReplace-*
mentNodePostProcessor. This is most useful for large documents/indexes, as it helps to retrieve more fine-grained details. We will examine two window sizes: 3 and 6.

5.4.4 ADVANTAGES AND LIMITATIONS

Table 5.5 provides an analysis of the advantages and limitations of each text segmentation method based on their methodological approach. Each of the three chunking methods presented in this section 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.

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

Table 5.5: Advantages and Limitations of different chunking strategies for RAG systems.

Method	Advantages	Limitations
Fixed Chunking	<ul style="list-style-type: none"> * Simple to implement and understand * Consistent chunk sizes facilitate uniform processing 	<ul style="list-style-type: none"> * Fixed chunk sizes may not align with natural breaks in the text * Larger chunks can introduce irrelevant information and increase computational costs
Small2Big	<ul style="list-style-type: none"> * Allows for fine-grained retrieval with the option to expand context * Adapts to different levels of specificity required by queries 	<ul style="list-style-type: none"> * More complex to implement and manage * Increased storage requirements due to redundancy in the hierarchy
Sentence Window	<ul style="list-style-type: none"> * Preserves semantic units (sentences) and their immediate context * Adapts to the natural structure of the text 	<ul style="list-style-type: none"> * Variable chunk sizes may complicate processing and indexing * Optimal window size may vary depending on the document type and content

5.4.5 EVALUATION

The Table 5.6 illustrates that as the chunk size increases, there is a minor downturn in the average latency. Consider that the pipeline's average response time increases as the chunk size increases, which is expected as the model has to process more tokens. Interestingly, the faithfulness (*i.e.* measuring how closely response is aligned with the source material) seems to reach its zenith at chunk_size of 1024, whereas average relevancy shows a consistent improvement with larger chunk sizes, also peaking at 1024. This suggests that a chunk size of 1024 might strike an optimal balance between response time and the quality of the responses, measured in terms of faithfulness and relevancy. In the sliding window method, the window size of 3 outperforms the window size of 6 in terms of both accuracy and F1 score, while maintaining nearly the same average response time. We can figure out that the 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.

5.5. SIMILARITY CUT-OFF

Table 5.6: Performance evaluation of various chunking strategy on the FactBench dataset, using the Gemma2 model.

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, <i>Chunking Strategy</i>				
Method	Parameters	Acc	F1	Latency
Original	Chuck Size: 256	0.8914	0.8914	0.04473s
	Chuck Size: 512	0.8932	0.8993	0.02670s
	Chuck Size: 1024	0.8946	0.8993	0.02378s
small2big	Chuck Size: 1024	0.8889	0.8953	0.19188s
Sliding Window	Window Size: 3	0.9014	0.9080	0.03076s
	Window Size: 6	0.9014	0.9077	0.03534s

We chose the Sliding Window with window size 3 method as it provides the surrounding context of the text and has the highest F1 score and accuracy.

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. The Algorithm 3 shows the similarity cut-off postprocessor implementation.

Algorithm 3 Similarity Cutoff Postprocessor (re-rank score)

```

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

```

Table 5.7: Performance evaluation of similarity cut-off method on the FactBench dataset, using the Gemma2 model.

ms-marco-MiniLM-L-6-v2, BAAI/bge-small-en-v1.5, Sliding Window (ws 3), <i>Similarity Cut-off</i>			
Method	Acc	F1	Latency Diff.*
w/o similarity cut-off (baseline)	0.8971	0.9036	–
similarity cut-off (original score)	0.9018	0.9080	-0.22237s
similarity cut-off (re-ranked score)	0.9014	0.9080	-0.350783

* The latency is compared to the baseline without the similarity cut-off on average per query.

Based on the results in Table 5.7, we decided to apply a similarity cut-off with the original score to provide the model with higher-quality data for further evaluations. This approach yields the highest accuracy and F1 score, though it is slightly slower than the re-ranker mode because it removes fewer irrelevant nodes. A drawback of re-ranking is that it may eliminate all relevant nodes based on low similarity score, requiring a re-run without the similarity cut-off, which increases processing time (not shown in the Table 5.7). In separate evaluations (not reported in Table 5.7), we also tested a normalized similarity cut-off using re-ranked scores scaled between the original minimum and maximum values. However, this normalized approach 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.

Table 5.8: Performance evaluation of different Top_k retrieval strategies on the FactBench dataset using the Gemma2 model.

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	5.21177s
Top_k 6	0.9032	0.9101	7.02713s

* The latency represents the average time per query for a complete run.

5.7. EVALUATION

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 on Table 5.8 it outperforms the top_k 3, so we will use top_k 6 for subsequent evaluations. The difference in latency is significant, presenting a trade-off between more data and response time. Since the quality of data for additional facts is uncertain, we aim to include more data in the context.

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 of the pipeline. As presented in Tables 5.9 and 5.10, along with Figure 5.3, 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 the provided configurations the model selected for *At_Most* methodology is the *Gemma2:21B*.

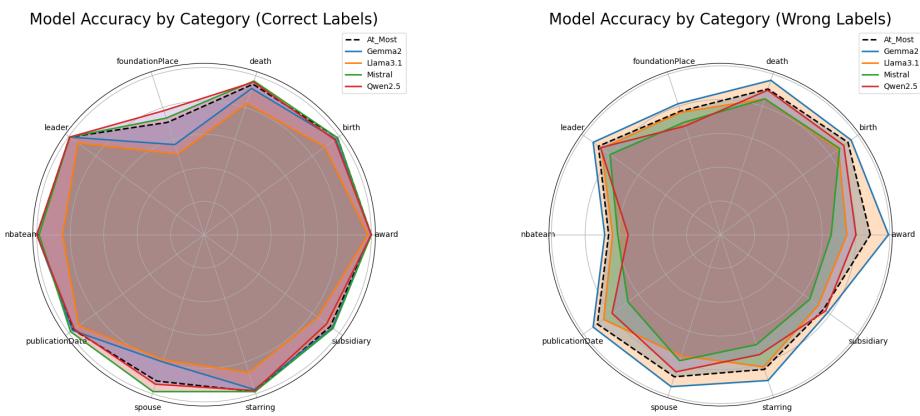


Figure 5.3: Category-wise performance of different models in identifying Positive Labels (left) and Negative Labels (right) on the FactBench dataset.

Table 5.9: Category-wise performance evaluation results of various models on the FactBench dataset.

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 5.10: Performance evaluation of various models on the FactBench dataset.

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

Model	Consistency ^a	Avg. Request Time ^b	Avg. tokens per request ^c	ACC	F1
Gemma2	0.8720	5.6826s	1605.29	0.9032	0.9101
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)	0.9176	16.815s	1604.196	0.8964	0.9071

^a Consistency score for each individual model is calculated across all models, excluding the ensemble models.

^b Each query uses two requests, so the average request duration is calculated based on the two requests.

^c The average is calculated based on the total number of input and output tokens combined.

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.

Error Type Categorization: After analyzing the failure cases, we categorized the errors into four main types:

- **Insufficient or Irrelevant Context:** In some cases, the provided context information does not directly support or refute the given triple. The LLMs

5.8. FAILURE ANALYSIS

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.
- **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.

Table 5.11: Example of failure cases and error analysis observed in the FactBench dataset using generated results and explanations.

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."

Based on Figure 5.4, 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.

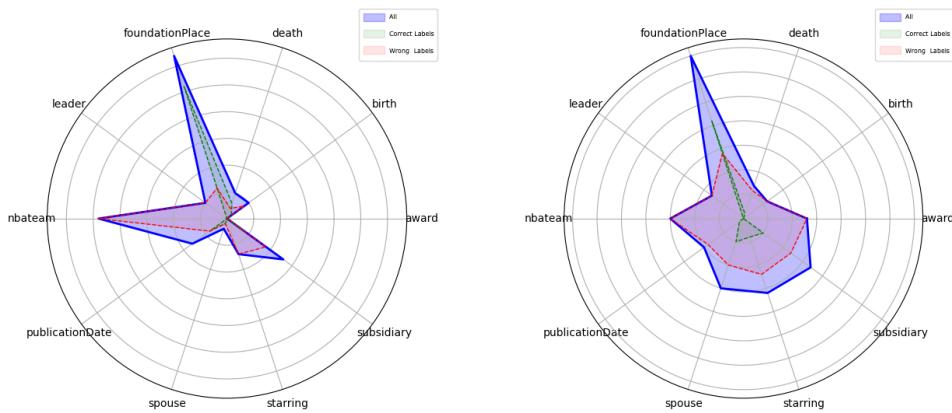


Figure 5.4: Prediction accuracy on the FactBench dataset, focusing on incorrect predictions. The right chart illustrates the Distribution of Fully Incorrect Predictions (4/4), detailing the instances where all predictions made by the models were incorrect. The left chart depicts the Distribution of Partially Incorrect Predictions (3/4).

In general, by spotting Figure 5.4, 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

LLMs show that they have changed the role of automated fact-checking, yet given the complexity and limitations of these models, while they are not 100% accurate, it is foreseeable that they will be able to act completely as human annotators in the future.

This thesis has presented FactCheck, a novel approach to KG fact verification using RAG. Through extensive experimentation and analysis, we have demonstrated the effectiveness of combining multiple LLMs with sophisticated IR techniques to verify facts in KGs. FactCheck’s prediction performance that measured against gold standard labels is 90% on *FactBench*, 87% on YAGO, and 70% on *DBpedia*. These results emphasize using LLMs to address fact verification in KGs.

One of our observations is that differences in LLMs architectures make them interpret the evidence differently and reach different conclusions despite accessing the same evidence. Additionally, another observation is the existence of particular challenges in geographic and nationality facts. This issue points to deeper issues in how language models process contextual information, indicating that improvements in basic reasoning capabilities may prove more valuable than simple increases in model scale.

Also, in fact-verification tasks, when a binary evaluation is requested (*i.e.* determining whether a statement is correct or incorrect), LLMs often explain when they classify a statement as incorrect. However, when predicting a statement is correct, they typically do not explain further if not explicitly requested.

The empirical evidence challenges the conventional wisdom that larger mod-

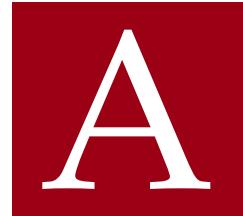
els invariably yield better results - our experiments with embedding models and text chunking techniques demonstrate that carefully optimized smaller models can match or exceed the performance of their larger counterparts when supported by robust retrieval mechanisms. This finding has profound implications for practical deployments, particularly in resource-constrained environments where computational efficiency is paramount.

Our analysis highlights a key challenge in current verification methods: ensemble techniques improve reliability by using consensus, but they also require a lot of computing power, which can be a problem for large-scale applications. This suggests a critical direction for future research, the development of more efficient verification strategies that maintain accuracy while reducing computational demands.

Finally, as KGs continue to grow in importance for real-world applications, these insights provide crucial guidance for developing more efficient and reliable verification systems that balance accuracy with practical constraints. Based on our findings and identified limitations, several promising directions for future research emerge:

1. **Enhanced Context Processing:** Develop more sophisticated methods for handling cases with insufficient or irrelevant context. Implement better techniques for identifying and resolving contradictions in retrieved information.
2. **Model Integration:** Explore additional strategies for combining model outputs beyond majority voting. Investigate dynamic model selection based on query characteristics. Implement more sophisticated tie-breaking mechanisms.
3. **Retrieval Optimization:** Improve query generation for better coverage of fact verification requirements. Develop more effective filtering mechanisms for irrelevant information. Enhance the similarity cut-off strategy for more precise document selection.
4. **Scalability Improvements:** Optimize computational resource usage for handling larger KGs. Develop more efficient document processing and embedding techniques. Implement parallel processing capabilities for faster verification.

5. **Explainability and Transparency:** Develop better methods for explaining verification decisions. Implement confidence scoring mechanisms. Create visualization tools for the verification process.
6. **Domain Adaptation:** Create specialized verification strategies for different types of facts. Develop domain-specific knowledge integration mechanisms. Implement adaptive learning capabilities for new domains.



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.4. REASONING 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:

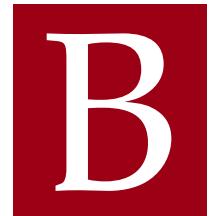
A.4 REASONING PROMPT

Prompt template for generate reason
Context information is below.

{context_str}

Given the context information without your knowledge,
you evaluate that the fact '{fact}' is {correct/wrong}, now explain
your reasoning for your evaluation. write you're answer in this format
'The claim is {correct/wrong} because: 'REASON''

without any further description or writing fact again, consider that you
can't change you're mind and you should reason why it's {correct/wrong}



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.

Table B.1: Evaluation of text segmentation using a chunk Size of 512, text chunks derived from the entry "Henry Dunant award Nobel Peace Prize".

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: https://digitalcommons.wku.edu/ijlc/vol5/iss1/5	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.2: Evaluation of text segmentation using a Small to Big technique (base chunk size 1024), text chunks derived from the entry "Henry Dunant award Nobel Peace Prize".

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.3: Evaluation of text segmentation using a Sliding Window with window size 3, text chunks derived from the entry "Henry Dunant award Nobel Peace Prize".

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.

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] Mehdi Ali et al. *The KEEN Universe: An Ecosystem for Knowledge Graph Embeddings with a Focus on Reproducibility and Transferability*. 2020. arXiv: [2001.10560](https://arxiv.org/abs/2001.10560) [cs.LG]. URL: <https://arxiv.org/abs/2001.10560>.
- [3] Kurt Bollacker et al. “Freebase: a collaboratively created graph database for structuring human knowledge”. In: *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*. SIGMOD ’08. Vancouver, Canada: Association for Computing Machinery, 2008, pp. 1247–1250. ISBN: 9781605581026. doi: [10.1145/1376616.1376746](https://doi.org/10.1145/1376616.1376746). URL: <https://doi.org/10.1145/1376616.1376746>.
- [4] Hyung Won Chung et al. *Scaling Instruction-Finetuned Language Models*. 2022. doi: [10.48550/ARXIV.2210.11416](https://doi.org/10.48550/ARXIV.2210.11416). URL: <https://arxiv.org/abs/2210.11416>.
- [5] Alexis Conneau et al. “Unsupervised Cross-lingual Representation Learning at Scale”. In: *CoRR* abs/1911.02116 (2019). arXiv: [1911.02116](https://arxiv.org/abs/1911.02116). URL: [http://arxiv.org/abs/1911.02116](https://arxiv.org/abs/1911.02116).
- [6] Jacob Devlin et al. *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*. 2019. arXiv: [1810.04805](https://arxiv.org/abs/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](https://arxiv.org/abs/2408.14317) [cs.CL]. URL: <https://arxiv.org/abs/2408.14317>.

REFERENCES

- [8] Xin Dong et al. "Knowledge vault: a web-scale approach to probabilistic knowledge fusion". In: *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '14. New York, New York, USA: Association for Computing Machinery, 2014, pp. 601–610. ISBN: 9781450329569. doi: 10.1145/2623330.2623623. URL: <https://doi.org/10.1145/2623330.2623623>.
- [9] Abhimanyu Dubey et al. *The Llama 3 Herd of Models*. 2024. arXiv: 2407.21783 [cs.AI]. URL: <https://arxiv.org/abs/2407.21783>.
- [10] 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.
- [11] 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>.
- [12] Google-Blog. *Introducing the Knowledge Graph: things, not strings*. Accessed: 2024-12-09. 2012. URL: <https://blog.google/products/search/introducing-knowledge-graph-things-not/>.
- [13] 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>.
- [14] Qingyu Guo et al. *A Survey on Knowledge Graph-Based Recommender Systems*. 2020. arXiv: 2003.00911 [cs.IR]. URL: <https://arxiv.org/abs/2003.00911>.
- [15] 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>.
- [16] Gautier Izacard et al. *Unsupervised Dense Information Retrieval with Contrastive Learning*. 2022. arXiv: 2112.09118 [cs.IR]. URL: <https://arxiv.org/abs/2112.09118>.
- [17] Albert Q. Jiang et al. *Mistral 7B*. 2023. arXiv: 2310.06825 [cs.CL]. URL: <https://arxiv.org/abs/2310.06825>.

- [18] M. Abdul Khaliq 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>.
- [19] Aditya Kusupati et al. *Matryoshka Representation Learning*. 2024. arXiv: 2205.13147 [cs.LG]. URL: <https://arxiv.org/abs/2205.13147>.
- [20] 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>.
- [21] Jens Lehmann et al. “Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia”. In: *Semantic web* 6.2 (2015), pp. 167–195.
- [22] Douglas B. Lenat. “CYC: a large-scale investment in knowledge infrastructure”. In: *Commun. ACM* 38.11 (Nov. 1995), pp. 33–38. ISSN: 0001-0782. DOI: [10.1145/219717.219745](https://doi.org/10.1145/219717.219745). URL: <https://doi.org/10.1145/219717.219745>.
- [23] Patrick Lewis et al. *Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks*. 2021. arXiv: 2005.11401 [cs.CL]. URL: <https://arxiv.org/abs/2005.11401>.
- [24] 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>.
- [25] Jerry Liu. *Tweet on Information Retrieval and LLMs*. Accessed: 2024-10-10. 2023. URL: <https://twitter.com/jerryjliu0/status/1708147687084986504>.
- [26] LlamaIndex Team. *Evaluating the Ideal Chunk Size for a RAG System Using LlamaIndex*. Accessed: 2024-12-11. 2024. URL: <https://www.llamaindex.ai/blog/evaluating-the-ideal-chunk-size-for-a-rag-system-using-llamaindex-6207e5d3fec5>.
- [27] 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>.

REFERENCES

- [28] Stefano Marchesin, Gianmaria Silvello, and Omar Alonso. "Utility-Oriented Knowledge Graph Accuracy Estimation with Limited Annotations: A Case Study on DBpedia". In: *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 12.1 (Oct. 2024), pp. 105–114. doi: 10.1609/hcomp.v12i1.31605. URL: <https://ojs.aaai.org/index.php/HCOMP/article/view/31605>.
- [29] Mistral AI Team. *Announcing Mistral 7B*. Accessed: 2024-10-17. 2023. URL: <https://mistral.ai/news/announcing-mistral-7b/>.
- [30] John X. Morris and Alexander M. Rush. *Contextual Document Embeddings*. 2024. arXiv: 2410.02525 [cs.CL]. URL: <https://arxiv.org/abs/2410.02525>.
- [31] 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. URL: <https://aclanthology.org/2023.eacl-main.148>.
- [32] 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. URL: <https://aclanthology.org/D17-1183>.
- [33] Reham Omar et al. *A Universal Question-Answering Platform for Knowledge Graphs*. 2023. arXiv: 2303.00595 [cs.AI]. URL: <https://arxiv.org/abs/2303.00595>.
- [34] Heiko Paulheim. "How much is a Triple? Estimating the Cost of Knowledge Graph Creation". In: *Proceedings of the ISWC 2018 Posters & Demonstrations, Industry and Blue Sky Ideas Tracks co-located with 17th International Semantic Web Conference (ISWC 2018), Monterey, USA, October 8th - to - 12th, 2018*. Ed. by Marieke van Erp et al. Vol. 2180. CEUR Workshop Proceedings. CEUR-WS.org, 2018. URL: https://ceur-ws.org/Vol-2180/ISWC%5C_2018%5C_Outrageous%5C_Ideas%5C_paper%5C_10.pdf.

- [35] 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>.
- [36] Daniel Ringler and Heiko Paulheim. "One Knowledge Graph to Rule Them All? Analyzing the Differences Between DBpedia, YAGO, Wikidata co." In: Sept. 2017, pp. 366–372. ISBN: 978-3-319-67189-5. doi: 10.1007/978-3-319-67190-1_33.
- [37] Stephen Robertson et al. "Okapi at TREC-3." In: Jan. 1994, pp. 0-.
- [38] 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>.
- [39] 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>.
- [40] 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>.
- [41] Jianlin Su et al. *RoFormer: Enhanced Transformer with Rotary Position Embedding*. 2023. arXiv: 2104.09864 [cs.CL]. url: <https://arxiv.org/abs/2104.09864>.
- [42] 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>.
- [43] Zafar Habeeb Syed, Michael Röder, and Axel-Cyrille Ngonga Ngomo. "FactCheck: Validating RDF Triples Using Textual Evidence". In: *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. CIKM '18. Torino, Italy: Association for Computing Machinery, 2018, pp. 1599–1602. ISBN: 9781450360142. doi: 10.1145/3269206.3269308. url: <https://doi.org/10.1145/3269206.3269308>.
- [44] 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>.

REFERENCES

- [45] Qwen Team. *Qwen2.5: A Party of Foundation Models*. Sept. 2024. URL: <https://qwenlm.github.io/blog/qwen2.5/>.
- [46] 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>.
- [47] Denny Vrandecic and Markus Krotzsch. “Wikidata: a free collaborative knowledgebase”. In: *Commun. ACM* 57.10 (Sept. 2014), pp. 78–85. issn: 0001-0782. doi: 10.1145/2629489. URL: <https://doi.org/10.1145/2629489>.
- [48] Liang Wang et al. *Improving Text Embeddings with Large Language Models*. 2024. arXiv: 2401.00368 [cs.CL]. URL: <https://arxiv.org/abs/2401.00368>.
- [49] Liang Wang et al. *Multilingual E5 Text Embeddings: A Technical Report*. 2024. arXiv: 2402.05672 [cs.CL]. URL: <https://arxiv.org/abs/2402.05672>.
- [50] Xiaohua Wang et al. *Searching for Best Practices in Retrieval-Augmented Generation*. 2024. arXiv: 2407.01219 [cs.CL]. URL: <https://arxiv.org/abs/2407.01219>.
- [51] An Yang et al. “Qwen2 Technical Report”. In: *arXiv preprint arXiv:2407.10671* (2024).
- [52] 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>.
- [53] 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>.
- [54] 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>.
- [55] Xuan Zhang and Wei Gao. *Towards LLM-based Fact Verification on News Claims with a Hierarchical Step-by-Step Prompting Method*. 2023. arXiv: 2310.00305 [cs.CL]. URL: <https://arxiv.org/abs/2310.00305>.

Acknowledgments

I want to thank everyone who believed in me and encouraged me to follow my passion. First, I am truly grateful to my supervisors, Prof. Stefano Marchesin and Prof. Gianmaria Silvello, for their guidance, support, and mentorship in shaping this thesis.

Next, I want to thank my family, especially my parents. A big thanks to K. Abedini (*i.e.* Ziedi) for suggestions and helping with the colors in the plots and encouraging me to finish this work. I also appreciate my long-time friend, M. Sohrabi, for always supporting and encouraging me.

Special thanks to A.A. Dehbaneh for the insightful discussions. I also want to thank my Italian friends, M. Cazzaro, M. Martinelli, and N. Boscolo, for their support and motivation.