# FactGenius: Combining Zero-Shot Prompting and Fuzzy Relation Mining to Improve Fact Verification with Knowledge Graphs

**Sushant Gautam**

Simula Metropolitan Center for Digital Engineering
Oslo, Norway
sushant@simula.no

## Abstract

Fact-checking is a crucial natural language processing (NLP) task that verifies the truthfulness of claims by considering reliable evidence. Traditional methods are often limited by labour-intensive data curation and rule-based approaches. In this paper, we present FactGenius, a novel method that enhances fact-checking by combining zero-shot prompting of large language models (LLMs) with fuzzy text matching on knowledge graphs (KGs). Leveraging DBpedia, a structured linked data dataset derived from Wikipedia, FactGenius refines LLM-generated connections using similarity measures to ensure accuracy. The evaluation of FactGenius on the FactKG, a benchmark dataset for fact verification, demonstrates that it significantly outperforms existing baselines, particularly when fine-tuning RoBERTa as a classifier. The two-stage approach of filtering and validating connections proves crucial, achieving superior performance across various reasoning types and establishing FactGenius as a promising tool for robust fact-checking. The code and materials are available at https://github.com/SushantGautam/FactGenius.

## 1 Introduction

Fact-checking is a critical task in natural language processing (NLP) that involves automatically verifying the truthfulness of a claim by considering evidence from reliable sources (Thorne et al., 2018). This task is essential for combating misinformation and ensuring the integrity of information in digital communication (Cotter et al., 2022). Traditional fact-checking methods rely heavily on manually curated datasets and rule-based approaches, which can be labour-intensive and limited in scope (Papadopoulos et al., 2024).

Recent advancements in large language models (LLMs) have shown promise in enhancing fact-checking capabilities (Choi and Ferrara, 2024).

LLMs, with their extensive pre-training on diverse textual data, possess a vast amount of embedded knowledge (Yang et al., 2024). However, their outputs can sometimes be erroneous or lacking in specificity, especially when dealing with complex reasoning patterns required for fact-checking. External knowledge, such as knowledge graphs (KGs) (Hogan et al., 2021), can aid in fact-checking.

In this paper, we propose FactGenius, a novel approach that combines zero-shot prompting of LLMs with fuzzy relation-mining technique to improve reasoning on knowledge graphs. Specifically, we leverage DBpedia (Lehmann et al., 2015), a structured source of linked data, to enhance the accuracy of fact-checking tasks.

Our methodology involves using the LLM to filter potential connections between entities in the KG, followed by refining these connections through Levenshtein distance-based fuzzy matching. This two-stage approach ensures that only valid and relevant connections are considered, thereby improving the accuracy of fact-checking.

We evaluate our method using the FactKG dataset (Kim et al., 2023b), which comprises 108,000 claims constructed through various reasoning patterns applied to facts from DBpedia. Our experiments demonstrate that FactGenius significantly outperforms existing baselines (Kim et al., 2023a), particularly when fine-tuning RoBERTa (Liu et al., 2019) as a classifier, achieving superior performance across different reasoning types.

In summary, the integration of LLMs with KGs and the application of fuzzy matching techniques represent a promising direction for advancing fact-checking methodologies. Our work contributes to this growing body of research by proposing a novel approach that effectively combines these elements, yielding significant improvements in fact-checking performance.
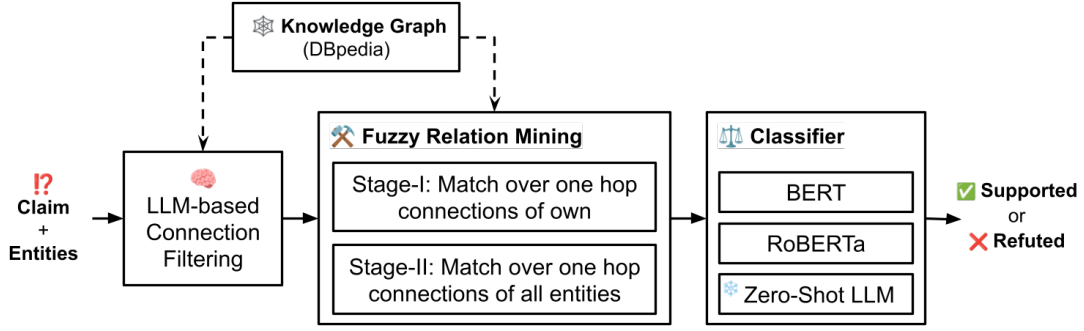
Figure 1: Overall pipeline of FactGenius: The process starts with LLM-based Connection Filtering using a knowledge graph (see Section 3.3.1). In Fuzzy Relation Mining (see Section 3.3.2), Stage-I matches one-hop connections of entities, and optionally, Stage-II includes all entities' connections. The classifier (BERT, RoBERTa, or Zero-Shot LLM; see Section 3.4) then determines if the claim is supported or refuted.

## 2 Literature Review

Fact-checking has become an increasingly vital aspect of natural language processing (NLP) due to the proliferation of misinformation in digital communication (Guo et al., 2022). Traditional approaches to fact-checking have typically relied on manually curated datasets and rule-based methods, which, while effective in controlled environments, often struggle with scalability and adaptability to new types of misinformation (Saquete et al., 2020; Guo et al., 2022). The labour-intensive nature of these methods also poses significant challenges in rapidly evolving information landscapes (Nakov et al., 2021; Zeng et al., 2021).

To address challenges in understanding machine-readable concepts in text, FactKG introduces a new dataset for fact verification with claims, leveraging knowledge graphs, encompassing diverse reasoning types and linguistic patterns, aiming to enhance reliability and practicality in KG-based fact verification (Kim et al., 2023b). Similarly, the Fact Extraction and VERification (FEVER) dataset (Thorne et al., 2018) pairs claim with Wikipedia sentences that support or refute them, providing a benchmark for fact-checking models. The authors employed a combination of natural language inference models and information retrieval systems to assess claim veracity. The GEAR framework (Zhou et al., 2019) improves fact verification by using a graph-based method to aggregate and reason over multiple pieces of evidence, surpassing previous methods by enabling evidence to interact.

Recent advancements in large language models (LLMs) have demonstrated considerable potential in enhancing fact-checking processes (Kim et al., 2023a; Choi and Ferrara, 2024). LLMs have been pre-trained on vast and diverse corpora (Yang et al., 2024), enabling them to generate human-like text and possess a broad knowledge base (Choi and Ferrara, 2024). However, despite their impressive capabilities, LLMs can produce outputs that are erroneous or lack the specificity required for complex fact-checking tasks (Choi and Ferrara, 2024). This is particularly evident when intricate reasoning and contextual understanding are necessary to verify claims accurately (Chai et al., 2023). Several studies have explored the integration of LLMs with external knowledge sources to improve their performance in fact-checking tasks (Cui et al., 2023; Ding et al., 2023).

The incorporation of knowledge graphs (KGs) into fact-checking frameworks has also garnered attention. KGs, such as DBpedia (Lehmann et al., 2015), provide structured and linked data that can enhance the contextual understanding of LLMs.

Knowledge graphs have been used to improve various NLP tasks by providing additional context and relationships between entities, as demonstrated by initiatives for knowledge-aware language models (Li et al., 2023; Logan Iv et al., 2019) and KG-BERT (Yao et al., 2019).

Approximate string matching (Navarro, 2001), also called fuzzy string matching, is a technique used to identify partial matches between text strings. Fuzzy matching techniques (Navarro, 2001) have been applied to enhance the integration of LLMs and KGs (Wang et al., 2024).

Levenshtein distance-based similarity measure (Levenshtein et al., 1966) helps in identifying strings which have approximate matches which can be useful to find relevant connections between entities by accommodating minor discrepancies in data representation This approach has been beneficial in refining the outputs of LLMs, ensuring that only valid and contextually appropriate connections are considered (Guo et al., 2023).

Our proposed method, FactGenius, builds on these advancements by combining zero-shot prompting of LLMs with a fuzzy relation-mining technique to improve reasoning over KGs. This methodology leverages DBpedia as a structured source of linked data to enhance fact-checking accuracy. By using LLMs to filter potential connections between entities and refining these connections through fuzzy matching, FactGenius aims to address the limitations of existing fact-checking models.

## 3 Methodology

FactGenius leverages the capabilities of a Large Language Model (LLM) to filter possible connections between entities in a Knowledge Graph (KG), particularly utilizing DBpedia (Lehmann et al., 2015) as a structured source of linked data.

Since the output of LLMs can be erroneous, the connections are further refined and enriched using Levenshtein distance (Levenshtein et al., 1966) and are also validated to ensure that such connections exist. This process is crucial for tasks such as fact-checking, where establishing valid and relevant connections between entities can validate or refute claims. Finally, the classifier, which can be fine-tuned over pre-trained models like BERT (Devlin et al., 2019) or RoBERTa (Liu et al., 2019), or a Zero-Shot LLM, determines whether the claim is supported or refuted. The overall pipeline is shown in Figure 1.

### 3.1 Dataset

The FactKG dataset (Kim et al., 2023b) is used which comprises 108,000 claims constructed through various reasoning patterns applied to facts sourced from DBpedia (Lehmann et al., 2015). Each data point consists of a natural language claim in English, the set of DBpedia entities mentioned in the claim, and a binary label indicating the claim's veracity (Supported or Refuted). The distribution across labels and five different reasoning types is shown in Table 1. The relevant relation paths starting from each entity in the claim are known which aids in the evaluation and development of models for claim verification tasks.

The dataset is accompanied by processed DBpedia, an undirected knowledge graph (KG). The dataset provides researchers with a valuable resource for exploring the intersection of natural language understanding and knowledge graph reasoning.

Table 1: Data distribution across labels and five reasoning types.

| Set | Train | Valid | Test |
|---|---|---|---|
| **Total Rows** | 86367 | 13266 | 9041 |
| True (Supported) | 42723 | 6426 | 4398 |
| False (Refuted) | 43644 | 6840 | 4643 |
| **One-hop** | 15069 | 2547 | 1914 |
| **Conjunction** | 29711 | 4317 | 3069 |
| **Existence** | 7372 | 930 | 870 |
| **Multi Hop** | 21833 | 3555 | 1874 |
| **Negation** | 12382 | 1917 | 1314 |

### 3.2 Claim Only Classifier

In this setting, where the models are given only the claim and tasked with predicting the label, it is expected that the model will heavily depend on stored evidence within its trained weights or identify patterns within the structure of the provided claims.

#### 3.2.1 Zero-shot Claim Only Baseline

A baseline is established using the Meta-Llama-3-8B-Instruct[1] (Meta, 2024) model with zero-shot promoting for claim verification, asking it to verify the claim without evidence. Through instruction prompt engineering, it is ensured that the model responds with either 'true' or 'false'. A retry mechanism is implemented to handle potential failures in LLM responses. A prompt example is shown in Figure 2.

#### 3.2.2 RoBERTa as Claim Only Fact Classifier

RoBERTa-base[2] is fine-tuned with claims as input, training it to predict Supported or Refuted. This is to compare with the BERT baseline reported in previous works (Kim et al., 2023b).

---

[1]huggingface.co/meta-llama/Meta-Llama-3-8B-Instruct
[2]huggingface.co/FacebookAI/roberta-base

```
[{
"role":"system", "content":
"You are an intelligent fact checker trained on
Wikipedia. You are given a single claim and your task
is to decide whether all the facts in the given claim are
supported by the given evidence using your knowledge.
Choose one of {True, False}, and output the one-sentence
explanation for the choice. "
},{
"role":"user", "content":
'''
## TASK:
Now let's verify the Claim based on the evidence.
Claim:
< < < Well, The celestial body known as 1097 Vicia has a
mass of 4.1kg.> > >

#Answer Template:
"True/False (single word answer),
One-sentence evidence."
'''
}]
```

Figure 2: Example prompt given to Llama3-Instruct without evidence for zero-shot fact-checking.
< < < ... > > > signs are added just to indicate that the content inside is different for each prompt.

```
[{
"role":"system", "content":
"You are an intelligent graph connection finder. You
are given a single claim and connection options for the
entities present in the claim. Your task is to filter
the Connections options that could be relevant to connect
given entities to fact-check Claim1. ~ ( tilde ) in the
beginning means the reverse connection. "
},{
"role":"user", "content":
'''
Claim1:
<<<Well, The celestial body known as 1097 Vicia has a
mass of 4.1kg.>>>

## TASK:
- For each of the given entities given in the DICT
structure below:
        Filter the connections strictly from the given
options that would be relevant to connect given entities
to fact-check Claim1.
- Think clever, there could be multi-step hidden
connections, if not direct, that could connect the
entities somehow.
- Prioritize connections among entities and arrange them
based on their relevance. Be extra careful with ~ signs.
- No code output. No explanation. Output only valid
python DICT of structure:

{
<<<
"1097_Vicia": ["...", "...", ... ],
#  options  (strictly  choose  from):  discovered,
formerName,   epoch,   periapsis,   apoapsis,   ...,
Planet/temperature

"4.1": ["...", "...", ... ],
#   options   (strictly   choose   from):   ~length,
~ethnicGroups,              ~percentageOfAreaWater,
~populationDensity, ~engine, ..., ~number
}
>>>
'''
}]
```

Figure 3: Example prompt given to Llama3-Instruct to filter potential connections between entities based on a given claim.

## 3.3 Graph Filtering

The graph filtering is divided into two main stages:

### 3.3.1 Filtering Possible Connections

This stage involves utilizing an LLM, particularly the Llama3-Instruct model, to identify and filter

potential connections between entities based on a given claim. The detailed steps are as follows:

**Data Preparation** Entity sets and their possible connections are extracted from the KG (DBpedia). Each entity and its associated possible connections form the initial input for the LLM.

**LLM Integration** The LLM is tasked with identifying relevant connections for each entity in the specific claim. The process involves:

1. Encoding each entity and its possible connections into a structured format suitable for the LLM.

2. Utilizing the LLM's inference capabilities to filter out irrelevant connections based on the context provided by the claim.

3. Generating a filtered set of connections in a structured format, which is then evaluated for completeness and relevance.

An example of the prompt used with LLM in Stage-I is shown in Figure 3. Prompts are crafted through iterative testing and refinement, aiming to optimize results and performance.

**Handling Invalid LLM Response** A retry mechanism is implemented to handle potential failures in LLM responses. If the LLM output is inadequate (e.g., empty or nonsensical), the request is retried up to a specified maximum number of attempts, typically 10. Throughout this experiment, however, we did not encounter any cases where the retry exceeded this limit.

### 3.3.2 Fuzzy Relation Mining

The LLM-filtered connections are then validated against the KG to ensure their existence and relevance. This involves:

1. **Stage-I**: Checking each connection filtered using LLM against the KG to confirm its validity. For each connection in the entities, perform fuzzy matching using Levenshtein distance to match entities in the first-hop relation of the graph. This approach accommodates speckling and reverse connection errors.

2. **Stage-II**: Matching potential connections fuzzily, while considering reverse relationships and similarities across all the one-hop connections in the knowledge graph of all entities within the claims.

The details are explained in Algorithm 1.

---

**Algorithm 1** Relationship Mining with Validation

---

1: **Input:** $A$ - dictionary of entities with their connections, $G$: Graph
2: **Output:** $probable\_connections$- dictionary of entities with updated and validated connections

3: **procedure** VALIDATERELATION($A$)
4:    Initialize: **probable_connections**: {}

5:    — **Stage-I** —

6:    **for each** entity, connections in $A$ **do**
7:       Retrieve: all **one_hop_connections** for **entity** $G$
8:       **for each** connection in **connections do**
9:          Fuzzily match from **one_hop_connections**
10:          Filter matches with a similarity score greater than 90
11:          Update **entity** in **probable_connections**
12:       **end for**
13:    **end for**

14:    — **Stage-II (optional)** —

15:    **all_connections** = all connections in **probable_connections**
16:    **for each** entity, connections in **probable_connections do**
17:       Retrieve: all **one_hop_connections** for **entity** $G$
18:       **for each** connection in **all_connections do**
19:          Fuzzily match from **one_hop_connections**
20:          Filter matches with a similarity score greater than 90
21:          Update **entity** in **probable_connections**
22:       **end for**
23:    **end for**

24: **end procedure**

---

### 3.4 With Evidence Classifier

In this configuration, the model is supplied with both the claim and graphical evidence as input, and it then makes predictions regarding the label. FactGenius utilizes graph filtering, as explained in Section 3.3, to ensure retention of the most relevant and accurate connections.

### 3.5 Zero-shot LLM as Fact Classifier

This involves utilizing Llama-3-Instruct as a fact classifier, to predict Supported or Refuted for the given input claim and evidence. A retry mechanism is implemented to handle potential failures in LLM responses. A prompt example with evidence is shown in Figure 4.

### 3.6 Fine-tuning pre-trained models

Pre-trained BERT-base-uncased[3] and RoBERTa-base are finetuned with claim and evidence as inputs to predict whether the claim is supported or refuted. In addition, an ablation evaluates the contributions of each stage of our approach. This involved sequentially removing Stage-II and measuring the performance of the system after the removal. The results of the ablation study

---
[3]huggingface.co/google-bert/bert-base-uncased

allowed us to quantify the impact of both stages on the overall performance of the model. Accuracy as an evaluation metric across all reasoning types was employed to quantify the performance improvements resulting from the ablation study.

```
[{
"role":"system", "content":
"You are an intelligent fact-checker. You are given a
single claim and supporting evidence for the entities
present in the claim, extracted from a knowledge graph.
Your task is to decide whether all the facts in the given
claim are supported by the given evidence.
Choose one of {True, False}, and output the one-sentence
explanation for the choice. "
},{
"role":"user", "content":
'''
## TASK:
Now let's verify the Claim based on the evidence.
Claim:
< < < Well, The celestial body known as 1097 Vicia has a
mass of 4.1kg.> > >

Evidences:
< < <
1999_Hirayama >- mass -> ""4.1""
1097_Vicia >- mass -> ""9.8"""
> > >

#Answer Template:
"True/False (single word answer),
One-sentence evidence."
'''
}]
```

Figure 4: Example prompt given to Llama3-Instruct with evidence for zero-shot fact-checking.

### 3.7 Implementation

Our FactGenius system implementation leverages several advanced tools and frameworks to ensure efficient and scalable processing. The Llama3-Instruct inference server is set up using vLLM (vLLM Project, 2024; Kwon et al., 2023), running on an NVIDIA A100 GPU (80 GB vRAM) to facilitate rapid inference. This server runs standalone, integrating seamlessly with the FactGenius pipeline.

For model fine-tuning and evaluation, we employ the Hugging Face Transformers library, utilizing the `Trainer` class for managing the training process. This setup allows for the fine-tuning of pre-trained models like BERT and RoBERTa on our pipeline. Hyper-parameters such as batch size, learning rate, and training epochs are configured to optimize performance, with computations accelerated by PyTorch.

## 4 Results

To evaluate the performance of our proposed methods, we conducted a series of experiments comparing different strategies for fact-checking. The results are summarized in Table 2.
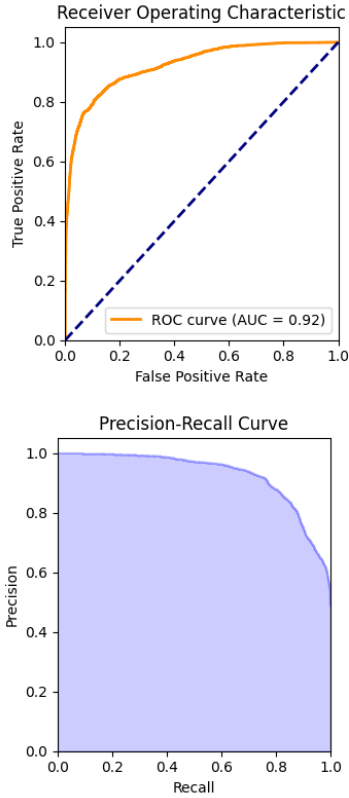


Figure 5: ROC curve (top) illustrates classifier discrimination ability with AUC 0.92, while Precision-Recall curve (bottom) reveals precision across recall levels for the best model in the test dataset.

### 4.1 Baseline and Claim Only Models

This achieved an accuracy of 0.61, demonstrating the inherent knowledge embedded within the Llama3 about the facts in our corpus. Adding evidence to the Llama3-Instruct model's instructions significantly improved its accuracy from 0.61 to 0.68. This reflects the trivial phenomenon that, incorporating relevant evidence can enhance fact-checking performance in a zero-shot learning scenario where the performance is mostly dependent on knowledge embedded in the model.

### 4.2 Comparison of Different Models

We compared the performance of RoBERTa, under the claim-only scenario. RoBERTa outperformed the reported accuracy of BERT (Kim et al., 2023b), achieving an accuracy of 0.68, which is on par with the 12-shot ChatGPT model reported in the KG-GPT paper (Kim et al., 2023a). This suggests that RoBERTa is a highly effective model for fact-checking tasks.

### 4.3 Incorporating Evidence with FactGenius

We employed a zero-shot prompting technique to filter relevant connections from the evidence, followed by fuzzy matching across multiple stages. This approach enabled us to retrieve evidence by searching over the graph. However, directly applying zero-shot prompting with Llama3-Instruct, even with evidence, did not yield superior performance. When using fine-tuned BERT as a classifier, the performance was comparable to the 12-shot KG-GPT model. However, fine-tuning RoBERTa led to a significant performance boost, achieving an accuracy of 0.85, the highest among all models tested, even surpassing the GEAR baseline(Zhou et al., 2019), which enhances fact verification by using a graph-based approach to aggregate and reason over multiple pieces of evidence.

### 4.4 Two-Stage Approach

To assess the contribution of our two-stage approach, we first apply only the first-stage graph filtering method (Stage-I) to filter the relationships, which achieved an accuracy of 0.83. Incorporating the second stage (Stage-II) further improved the performance to 0.85. The second stage particularly enhanced performance across all reasoning types, with notable improvements in conjunction and negation tasks. Interestingly, for the existence reasoning type, the BERT classifier performed on par with the best models, indicating its robustness for this specific task.

### 4.5 Overall Performance

FactGenius with a fine-tuned classifier, demonstrated superior performance across all reasoning types compared to previously reported accuracies. This validates the effectiveness of our multi-stage evidence retrieval and classifier fine-tuning approach in improving fact-checking accuracy. Refer to Figure 5 for the ROC and Precision-Recall curves illustrating the classifier performance of the best FactGenius variant with two-stage filtering relationship mining and fine-tuned RoBERTa classifier and to Table 3 and 4 for the classification report and confusion matrix, respectively.

Table 2: Comparing our method with other strategies and methods in terms of reported accuracies in the test set. * indicates results obtained from KG-GPT paper (Kim et al., 2023a).

| Input type | Model | Variants | One-hop | Conjunction | Existence | Multi-hop | Negation | Total |
|---|---|---|---|---|---|---|---|---|
| Claim Only | baseline | Llama3-Instruct-zero-shot | 0.61 | 0.67 | 0.59 | 0.61 | 0.53 | 0.61 |
| | Fact-KG | BERT* | 0.69 | 0.63 | 0.61 | 0.70 | 0.63 | 0.65 |
| | KG-GPT | ChatGPT (12-shot)* | - | - | - | - | - | 0.68 |
| | baseline | RoBERTa | 0.71 | 0.72 | 0.52 | 0.74 | 0.54 | 0.68 |
| With Evidence | Fact-KG | GEAR* | 0.83 | 0.77 | 0.81 | 0.68 | 0.79 | 0.77 |
| | KG-GPT | KG-GPT (12-shot)* | - | - | - | - | - | 0.72 |
| | FactGenius | Llama3-Instruct-zero-shot | 0.72 | 0.75 | 0.76 | 0.62 | 0.52 | 0.68 |
| | | BERT-two-stage | 0.75 | 0.67 | 0.94 | 0.66 | 0.79 | 0.72 |
| | | RoBERTa-single-stage | 0.87 | 0.82 | 0.94 | 0.75 | 0.84 | 0.83 |
| | | **RoBERTa-two-stage** | **0.89** | **0.85** | **0.95** | **0.75** | **0.87** | **0.85** |

Table 3: Classification report for the best model across the test dataset.

| | Precision | Recall | F1 | Support |
|---|---|---|---|---|
| **Refuted** | 0.81 | 0.93 | 0.86 | 4643 |
| **Supported** | 0.91 | 0.77 | 0.83 | 4398 |
| **Accuracy** (average) | 0.86 | 0.85 | 0.85 | 9041 |

Table 4: Confusion matrix for the best model across the test dataset.

| | | Predicted | |
|---|---|---|---|
| | | **Refuted** | **Supported** |
| **Actual** | **Refuted** | 4315 | 328 |
| | **Supported** | 1031 | 3367 |

## 5 Discussion

The improved performance of FactGenius, particularly in Conjunction, Existence, and Negation reasoning, can be attributed to its innovative combination of zero-shot prompting with large language models and fuzzy text matching on knowledge graphs.

The two-stage approach, which involves an initial filtering phase followed by a validation phase, significantly enhances accuracy. However, the model shows moderate performance improvement in Multi-hop reasoning, indicating the need for more advanced techniques to handle its complexity.

The two-stage approach of filtering and validating connections proved to be particularly effective. In the first stage, the LLM helps to narrow down potential connections based on the context provided by the claim. This initial filtering significantly reduces the search space, making the subsequent validation stage more efficient. The second stage further refines these connections through fuzzy matching, ensuring that only the most relevant and accurate connections are retained. The comparative study confirmed the importance of each stage, showing that the second stage particularly enhances performance in conjunction and negation reasoning tasks.

As having an LLM inference server is a crucial component of this framework, we employed vLLM (vLLM Project, 2024) to streamline rapid inference with a single NVIDIA A100 GPU. In our experiment, the LLM inference speed was around 15 queries per second, including retries in case of failure. This rate is feasible, considering that LLM inference is continually optimized with the latest technologies. Embedding LLM in a framework has proven to be a wise choice.

## 6 Conclusion

In this paper, we introduced FactGenius, a novel method that combines zero-shot prompting of large language models with fuzzy relation mining for superior reasoning on knowledge graphs. This approach addresses several key challenges in traditional fact-checking methods. First, the integration of LLMs allows for the leveraging of extensive pre-trained knowledge, which is crucial for understanding and verifying complex claims through structured data from DBpedia. Second, the use of fuzzy text matching with Levenshtein distance ensures that minor discrepancies in entity names or relationships do not hinder the relationship selection process, thus improving robustness.

Our experiments on the FactKG dataset demonstrated that FactGenius significantly outperforms traditional fact-checking methods and existing baselines, particularly when fine-tuning RoBERTa as a classifier. The two-stage approach of filtering and validating connections proved crucial for achieving high accuracy across different reasoning types. This underscores the potential of FactGenius to improve fact-checking accuracy without requiring complex stages and components.

The findings from this study suggest that integrating LLMs with structured knowledge graphs and fuzzy matching techniques holds great promise for advancing fact-checking capabilities. Future work could explore the application of this approach to other domains and datasets, as well as the potential for incorporating additional sources of structured data to further enhance performance.

By improving the accuracy and efficiency of fact-checking, FactGenius contributes to the broader effort of combating misinformation and ensuring the reliability of information in digital communication.

## Acknowledgement

## References

Ziwei Chai, Tianjie Zhang, Liang Wu, Kaiqiao Han, Xiaohai Hu, Xuanwen Huang, et al. 2023. GraphLLM: Boosting Graph Reasoning Ability of Large Language Model. *arXiv*.

Eun Cheol Choi and Emilio Ferrara. 2024. FACT-GPT: Fact-Checking Augmentation via Claim Matching with LLMs. In *WWW '24: Companion Proceedings of the ACM on Web Conference 2024*, pages 883–886. Association for Computing Machinery, New York, NY, USA.

Kelley Cotter, Julia R. DeCook, and Shaheen Kanthawala. 2022. Fact-Checking the Crisis: COVID-19, Infodemics, and the Platformization of Truth. *Social Media + Society*, 8(1):20563051211069048.

Jiaxi Cui, Zongjian Li, Yang Yan, Bohua Chen, and Li Yuan. 2023. ChatLaw: Open-Source Legal Large Language Model with Integrated External Knowledge Bases. *arXiv*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *ACL Anthology*, pages 4171–4186.

Yan Ding, Xiaohan Zhang, Saeid Amiri, Nieqing Cao, Hao Yang, Andy Kaminski, et al. 2023. Integrating action knowledge and LLMs for task planning and situation handling in open worlds. *Auton. Robot.*, 47(8):981–997.

Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A Survey on Automated Fact-Checking. *Transactions of the Association for Computational Linguistics*, 10:178–206.

Zishan Guo, Renren Jin, Chuang Liu, Yufei Huang, Dan Shi, Supryadi, et al. 2023. Evaluating Large Language Models: A Comprehensive Survey. *arXiv*.

Aidan Hogan, Eva Blomqvist, Michael Cochez, Claudia D'amato, Gerard De Melo, Claudio Gutierrez, et al. 2021. Knowledge Graphs. *ACM Comput. Surv.*, 54(4):1–37.

Jiho Kim, Yeonsu Kwon, Yohan Jo, and Edward Choi. 2023a. KG-GPT: A General Framework for Reasoning on Knowledge Graphs Using Large Language Models. *ACL Anthology*, pages 9410–9421.

Jiho Kim, Sungjin Park, Yeonsu Kwon, Yohan Jo, James Thorne, and Edward Choi. 2023b. FactKG: Fact Verification via Reasoning on Knowledge Graphs. *ACL Anthology*, pages 16190–16206.

Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, et al. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *SOSP '23: Proceedings of the 29th Symposium on Operating Systems Principles*, pages 611–626. Association for Computing Machinery, New York, NY, USA.

Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N. Mendes, et al. 2015. DBpedia – A large-scale, multilingual knowledge base extracted from Wikipedia. *Semantic Web*, 6(2):167–195.

Vladimir I Levenshtein et al. 1966. Binary codes capable of correcting deletions, insertions, and reversals. In *Soviet physics doklady*, volume 10, pages 707–710. Soviet Union.

Xinze Li, Yixin Cao2, Liangming Pan, Yubo Ma, and Aixin Sun. 2023. Towards Verifiable Generation: A Benchmark for Knowledge-aware Language Model Attribution. *arXiv*.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, et al. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv*.

Robert L. Logan Iv, Nelson F. Liu, Matthew E. Peters, Matt Gardner, and Sameer Singh. 2019. Barack's Wife Hillary: Using Knowledge-Graphs for Fact-Aware Language Modeling. *arXiv*.

Meta. 2024. Meta Llama 3. [Online; https://llama.meta.com/llama3].

Preslav Nakov, David Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, et al. 2021. Automated Fact-Checking for Assisting Human Fact-Checkers. In *Proceedings of the Thirtieth International Joint Conference onArtificial Intelligence, {IJCAI-21}*, pages 4551–4558. International Joint Conferences on Artificial Intelligence Organization.

Gonzalo Navarro. 2001. A guided tour to approximate string matching. *ACM Comput. Surv.*, 33(1):31–88.

Stefanos-Iordanis Papadopoulos, Christos Koutlis, Symeon Papadopoulos, and Panagiotis C. Petrantonakis. 2024. VERITE: a Robust benchmark for multimodal misinformation detection accounting for unimodal bias. *Int. J. Multimed. Info. Retr.*, 13(1):1–15.

Estela Saquete, David Tomás, Paloma Moreda, Patricio Martínez-Barco, and Manuel Palomar. 2020. Fighting post-truth using natural language processing: A review and open challenges. *Expert Syst. Appl.*, 141:112943.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a Large-scale Dataset for Fact Extraction and VERification. *ACL Anthology*, pages 809–819.

vLLM Project. 2024. vLLM. [Online; https://github.com/vllm-project/vllm].

Yu Wang, Nedim Lipka, Ryan A. Rossi, Alexa Siu, Ruiyi Zhang, and Tyler Derr. 2024. Knowledge Graph Prompting for Multi-Document Question Answering. *AAAI*, 38(17):19206–19214.

Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, et al. 2024. Harnessing the Power of LLMs in Practice: A Survey on ChatGPT and Beyond. *ACM Trans. Knowl. Discovery Data*, 18(6):1–32.

Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. KG-BERT: BERT for Knowledge Graph Completion. *arXiv*.

Xia Zeng, Amani S. Abumansour, and Arkaitz Zubiaga. 2021. Automated fact-checking: A survey. *Lang. Linguist. Compass*, 15(10):e12438.

Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, et al. 2019. GEAR: Graph-based Evidence Aggregating and Reasoning for Fact Verification. *ACL Anthology*, pages 892–901.