Fact or Fiction? Improving Fact Verification with Knowledge Graphs through Simplified Subgraph Retrievals

Anonymous ACL submission

Abstract

Despite recent success in natural language processing (NLP), fact verification still remains a difficult task. Due to misinformation spreading increasingly fast, attention has been directed towards automatically verifying the correctness of claims. In the domain of NLP, this is usually done by training supervised machine learning models to verify claims by utilizing evidence from trustworthy corpora. We present efficient methods for verifying claims on a dataset where the evidence is in the form of structured knowledge graphs. We use the FACTKG dataset, which is constructed from the DBpedia knowledge graph extracted from Wikipedia. By simplifying the evidence retrieval process, from fine-tuned language models to simple logical retrievals, we are able to construct models that both require less computational resources and achieve better test-set accuracy.

1 Introduction

002

003

007

011

012

014

021

033

037

041

As the volume of information generated continues to grow, so does the risk of misinformation spreading, which has made automatic fact verification a crucial task in NLP (Cohen et al., 2011; Hassan et al., 2015; Thorne and Vlachos, 2018; Bekoulis et al., 2021). Traditionally, fact verification has been tackled in journalism by experts manually researching topics and writing articles about their findings. Some specific websites dedicated to this approach are *FactCheck.org* and *PolitiFact.com*. However, it is time consuming and labor intensive, and is not able to follow the pace of the creation of information in digital media (Cohen et al., 2011; Hassan et al., 2015).

One of the most popular datasets for fact verification is the *Fact Extraction and VERification* (FEVER) dataset (Thorne et al., 2018). It consists of claims supported by a corpus of Wikipedia articles. Models trained on the dataset need to extract the relevant evidence and use it to classify claims as *supported*, *refuted* or *not enough information*. Despite its popularity, several issues have been discovered. Due to the manual construction of claims, the structure of the language is inherently biased with respect to the classes, and therefore it is possible to achieve good performance without using the evidence at all (Schuster et al., 2019). It has also been shown that models trained on FEVER experience a significant drop in performance when the factual evidence is changed in a way that influences the validity of claims (Hidey et al., 2020). These issues can be improved by accordingly adjusting the validation and test dataset to contain less biased data (Schuster et al., 2019; Hidey et al., 2020), but we believe it is important to develop models on other datasets as well. 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

A less studied approach to process evidence is by structured data. In many real-world examples, data is available in large structured databases, rather than unstructured articles. This is relevant for domains such as social networks, logistics, management systems and database systems. The dataset *TabFact* (Chen et al., 2019) is created with this intent, consisting of claims with tabular evidence extracted from Wikipedia.

We will dedicate this article to increase the performance of models trained on the FACTKG dataset (Kim et al., 2023), a dataset created for fact verification with structured evidence in the form of knowledge graphs (KGs). The claims are created with evidence from DBpedia (Lehmann et al., 2015), a large KG extracted from Wikipedia. A KG consists of nodes and edges linked together to represent structural concepts. Nodes represent entities, such as persons, things or events, and edges represent relations, conveying how entities are related, as shown in Figure 1. For instance, a node can be the company Meyer Werft, and since it is located in the city Papenburg, they are connected with the edge location. We refer to Meyer Werft, location, Papenburg as a knowledge triple.

Since the task of fact verification with KGs re-



Figure 1: An example claim from FACTKG (Kim et al., 2023). The claim can be verified or refuted based on the DBpedia KG (Lehmann et al., 2015). This is Figure 1 from Kim et al. (2023).

mains relatively unexamined, we want to explore several different approaches to the problem. We use the following three model architectures:

• **Textual Fine-tuning:** Fine-tuning pretrained encoder models on text data for claim verification. We use BERT (Devlin et al., 2018) by concatenating the claims with subgraphs represented as strings.

• Hybrid Graph-Language Model: Using a modification of a *question answer graph neural network* (QA-GNN) (Yasunaga et al., 2021), which both uses a pretrained encoder model to embed the claim, and a graph neural network (GNN) to structurally process the subgraphs.

• LLM Prompting: Deploying state-of-the-art language models in a few-shot setting, without the need for additional finetuning. We use ChatGPT 40 (Achiam et al., 2023; Open AI, 2024) for this setting.

The textual finetuning serves as a simple and conventional method, while the QA-GNN can handle graph based data efficiently and is more specifically constructed for the task of interest. In contrast, the LLM prompting displays how well general purpose language models can perform on the task. It does not require any further training and does not use any evidence. Therefore, it will serve as a baseline and give insight to how difficult the task is.

Our main contribution is that we increase the accuracy and computational efficiency of models

trained on FACTKG. By utilizing efficient subgraph retrieval methods, we are able to substantially increase the test-set accuracy from 77.65% (Kim et al., 2023) to 93.49%. To the best of the authors knowledge, this is the best performance achieved so far on the dataset. Additionally, our models train quicker, taking only 1.5-10 hours, compared to the 2-3 days spent on the benchmark model from Kim et al. (2023), reported by the authors. The code and documentation can be found at https://anonymous.4open.science/ r/Fact-or-Fiction/README.md. 114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

162

2 Related Work

2.1 Fact Verification

The FEVER dataset is one of the most popular datasets used for fact verification (Thorne et al., 2018), and has influenced several model architectures. Graph-based Evidence Aggregating and Reasoning (GEAR) (Zhou et al., 2019) works by finding relevant articles with entity linking, giving them a relevance score, embedding the claim and sentences in the relevant evidence with a pre-trained BERT (Devlin et al., 2018), and then using a GNN to reason over the embeddings. The Neural Semantic Matching Network (NSMN) (Nie et al., 2019) used three homogenous neural networks used for document retrieval, sentence selection and claim verification. By using a transformer based architecture, Generative Evidence REtrieval (GERE) (Chen et al., 2022) combined the evidence retrieval and sentence identifying into a single step.

Several other datasets for fact verification have also been proposed. The *Fake News Challenge* (Hanselowski et al., 2018) were aimed towards predicting the relevance and agreement of a title and text. *VitaminC* (Schuster et al., 2021) focuses on representing changing evidence, and was created by constructing claims based on different revisions of Wikipedia articles. The dataset *FAVIQ* (Park et al., 2021) explored ambiguous parts of claims, while *TabFact* (Chen et al., 2019) used tabular data as evidence. There have also been proposed multimodular dataset for fact verification, combining claims and images (Zlatkova et al., 2019; Mishra et al., 2022).

2.2 The FactKG Dataset

The FACTKG dataset (Kim et al., 2023) consists of 108 000 English claims for fact verification, where the downstream task is to predict whether

101

102

104

105

106

107

108

109

110

111

112

the claims are true or false. The claims are constructed from the DBpedia KG (Lehmann et al., 2015), which is extracted from Wikipedia and represents how entities are related to each other.

163

164

165

166

168

169

171

172

173

174

175

178

179 180

181

182

183

187

189

190

191

192

196

197

198

204

205

207

The claims are constructed on either of the following five reasoning types:

- **One-hop:** To answer a one-hop claim, one only needs to traverse one edge in the KG. In other words, only one knowledge triple is needed to verify the validity of the claim.
- **Multi-hop:** As opposed to one-hop claims, one needs to traverse multiple steps in the KG to verify multi-hop claims.
- **Conjunction:** The claim includes a combination of multiple claims, which are often added together with the word *and*.
- Existence: These claims state that an entity has a relation, but does not specify which entity it relates to.
- Negation: The claim contains negations, such as *not*.

The dataset is split in a train-validation-test set of proportion 8:1:1. The train and validation set includes relevant subgraphs for each claim, but not the test set. All claims include a list of entities present in the claim and as nodes in the KG.

2.3 Question Answer Graph Neural Networks

The *question answer graph neural network* (QA-GNN) (Yasunaga et al., 2021) is a hybrid language and GNN model that both uses a pre-trained language model to process the text, and couples it with a GNN reasoning over a subgraph. It is given text and a subgraph as input. The text, consisting of a question and possible answers, is added as a node to the subgraph. The language model embeds the text, and assigns a relevance score to each node in the subgraph. The relevance scores are multiplied with the node features, before being sent into the GNN. The GNN output, text-node and the text embedding are concatenated before being put into the classification layer.

3 Methods

3.1 Efficient Subgraph Retrieval

We experiment with different ways of retrieving relevant subgraphs for the claim, focusing on com-







Figure 2: **Examples of the different subgraphs explored in this article**. Boxes and **bold letters** represent entities, while arrows and *italic letters* represent relations. This claim is meant for illustrative purposes and is not present in the FACTKG dataset.

putational efficiency. Each datapoint in the FAC-TKG dataset consists of a claim and a list of entities that appear both in the claim and the KG. Since the part of DBpedia used in FactKG is fairly large (1.53GB), it is necessary to only use a small subgraph of it as input to the models. The benchmark model from Kim et al. (2023) uses two language models to predict the relevant edges and the depth of the graph. We wish to simplify this step in order to reduce the model complexity, and propose non-trainable methods for subgraph retrieval.

We experiment with the following methods (examples can be found in Figure 2):

- **Direct:** Only includes knowledge triples where both nodes are present in the entity list.
- **Contextualized:** First, includes all direct subgraphs. Additionally, lemmatize the words in the claim and check if the nodes in the entity list have any relations corresponding to the lemmatized words in the claim. Include all knowledge triples where at least one node is

325

326

278

279

in the entity list and the relation is found in the claim.

• **Single-step:** Includes every knowledge triple one can be traversed in one step from a node in the entity list. In other words, include every knowledge triple that contains at least one node in the entity list.

3.2 Finetuning BERT

230

234

235

237

238

241

242

246

247

248

249

254

256

261

265

267

269

273

274

275

277

We use BERT (Devlin et al., 2018) as our pretrained language model. We first train a baseline model using only the claims and no subgraphs, and then with all of the different methods for retrieving subgraphs. The subgraphs are converted to strings, where each knowledge triple is represented with square brackets, and the name of the nodes and edges are the same as they appear in DBpedia. The order of the knowledge triples is determined by the order of the list of entities in the FactKG dataset and the order of the edges in DBpedia. The subgraphs are concatenated after the claims and a " I" separation token.

3.3 QA-GNN Architecture

In order to adapt the QA-GNN to the fact verification setting, we perform some slight modifications. Because the possible answers are always "true" or "false", we embed only the claims, instead of the question and answer combination. Additionally, we do not connect the embedded question or claim to the subgraph.

We use a pre-trained BERT (Devlin et al., 2018) as the language model to embed and calculate the relevance scores. In order to reduce the complexity of the model, we use a frozen BERT to calculate the embeddings for the nodes and the edges in the graph. This way, all of the embeddings in the graph can be pre-calculated. We use the last hidden layer representation of the CLS token, which is of length 768. The BERT that calculates the relevance scores and the embedding of the claim is not frozen. The relevance scores are computed as the cosine similarity between the claim embedding and the embedding of the text in the nodes.

We use a graph attention network (Veličković et al., 2017) for our GNN. Since the subgraphs are quite shallow, we only use two layers in the GNN, and apply batch norm (Ioffe and Szegedy, 2015). Each layer has 256 features, which is mean-pooled and concatenated with the BERT embedding and sent into the classification layer. We add dropout (Srivastava et al., 2014) to the GNN layers and the classification layer.

3.4 ChatGPT Prompting

We construct a prompt for ChatGPT 40 in order to answer a list of claims as accurately as possible. This is done by creating an initial prompt and validating the results on 100 randomly drawn claims from the validation set, and by trying different configurations of the prompt until we do not get a better validation set accuracy. We then use the best prompt with 100 randomly drawn unseen test-set questions, and attempt to ask 25, 50 and 100 claims at a time, to see if the amount of claims at a time influences the performance. We run the testing three times.

Since we do not have access to vast enough computational resources to run an LLM, this analysis is limited by only using 100 datapoints from the test set. In order to get access to a state-of-the-art LLM, we used the ChatGPT website with a "Chat-GPT Plus" subscription to perform the prompting. This model is not seeded, so the exact answers are not reproducible, but every prompt and answer are available in the software provided with this article. We used the ChatGPT 40 model 30th of May 2024. Every prompt was performed in the "temporary chat" setting, so the model did not have access to the history of previous experiments.

Due to the inability to use the entire test set and the lack of reproducibility, we do not directly compare this experiment to the other models. However, we still believe it serves as a valuable benchmark. Recently, the performance of LLMs has rapidly improved, which suggests that their applications will continue to broaden. Additionally, this approach is not fine-tuned, and may work as an interesting benchmark that can contextualize the results of the other models.

3.5 Benchmark Models

We will compare the results against the best benchmark models from Kim et al. (2023) and the best performing models known to the authors, found in Gautam (2024). These comparisons include both baselines that use only the claims and models that also incorporate subgraph evidence. Claim-Only Models:

• FactKG BERT Baseline: The baseline model from Kim et al. (2023) uses a fine-tuned BERT, training only on the claims.

Input Type	Model	One-hop	Conjunction	Existence	Multi-hop	Negation	Total
Claim Only	FACTKG BERT Baseline	69.64	63.31	61.84	70.06	63.62	65.20
Claim Only	FactGenius RoBERTa Baseline	71	72	52	74	54	68
	BERT (no subgraphs)	67.71	67.48	62.51	73.28	64.23	68.99
	FACTKG GEAR Benchmark	83.23	77.68	81.61	68.84	79.41	77.65
With Subgraphs	FactGenius RoBERTa-two-stage	89	85	95	75	87	85
	QA-GNN (single-step)	79.08	74.43	83.37	74.72	79.60	78.08
	BERT (single-step)	97.40	97.51	97.31	80.32	92.54	93.49

Table 1: **Test-set accuracy for the best models from this article and the best benchmark models.** The FACTKG models are from Kim et al. (2023), while the FactGenius models are from Gautam (2024). The fine-tuned BERT model performed the best, while the QA-GNN was the computationally most efficient model.

• **RoBERTa Baseline**: Similarly to above, the baseline from Gautam (2024) uses a finetuned language model with claims only, but uses RoBERTa (Liu et al., 2019) as the base model.

Models Utilizing Subgraphs:

327

328

332

334

336

337

338

340

342

344

346

347

348

350

354

356

361

362

- GEAR-Based Model: The benchmark model from Kim et al. (2023) is inspired by GEAR (Zhou et al., 2019), but has been adapted to handle graph-based evidence. It uses two finetuned language models to retrieve the subgraphs. One of them predicts relevant edges, the other predicts the depth of the subgraph.
- FactGenius: This model combines zero-shot LLM prompting with fuzzy text matching on the KG (Gautam, 2024). The LLM filters relevant parts of the subgraphs, which are then refined using fuzzy text matching. Finally, a fine-tuned RoBERTa is used to make the downstream prediction.

3.6 Further Experimental Details

Due to computational constraints, we tuned the hyperparameters one by one, instead of performing a grid search. All the training was performed on a RTX 2080 Ti GPU with 11GB VRAM. The BERT model has 109 483 778 parameters, which all were fine-tuned. The QA-GNN used a total of 109 746 945 parameters. The FACTKG dataset comes with a lighter version of DBpedia that only contains relevant entries, which was used for this paper. Further details can be found in Appendix A.

4 Results

4.1 Improved Performance and Efficiency

The test results for our best model configurations and the benchmark models can be found in Table 1. The best performing model is the fine-tuned BERT with single-step subgraphs. The fine-tuned BERT without any subgraphs were able to achieve slightly higher performance than the one from Kim et al. (2023), which we suggest is due to finding better hyperparameters.

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

385

386

387

388

390

391

392

393

394

395

396

397

398

399

400

401

Additionally, our models were much faster to train. While the GEAR model used 2-3 days to train on an RTX 3090 GPU (reported by the authors by email), our QA-GNN only used 1.5 hours. The training time of our fine-tuned BERT model was greatly influenced by the size of the subgraphs we used. With no subgraphs, it took about 2 hours to train, while with the one-hop subgraph it took 10 hours. FactGenius was reported to use substantially more computational resources, running the LLM inference on a A100 GPU with 80GB VRAM for 8 hours.

4.2 Successful Subgraphs Retrievals

We now look at the different configurations for the subgraph retrievals, which greatly influenced the performance of the models. Since the *direct* and *contextual* approach only includes subgraphs if a certain requirement is fulfilled, it will result in some of the claims having empty subgraphs. In the training and validation set, 49.0% of the graphs were non-empty for the *direct* approach, and 62.5% were non-empty for the *contextual* approach. The *single-step* method resulted in vastly bigger subgraphs.

While the QA-GNN could handle the big subgraphs efficiently, the fine-tuned BERT was severely slowed down when the size of the subgraphs got bigger. Therefore, we substituted any empty subgraphs with the *single-step* subgraph when using QA-GNN, but kept the empty graphs when using fine-tuned BERT. This means that some claims for the direct and contextual BERT models were predicted only using the bias in the language model and the claim.

Model	One-hop	Conjunction	Existence	Multi-hop	Negation	Total
BERT (no subgraphs)	67.71	67.48	62.51	73.28	64.23	68.99
BERT (direct)	80.24	83.30	59.05	77.62	74.58	79.64
BERT (contextual)	81.20	84.45	61.05	77.04	77.40	80.25
BERT (single-step)	97.40	97.51	97.31	80.32	92.54	93.49
QA-GNN (direct)	74.60	74.01	58.97	76.41	74.12	75.01
QA-GNN (contextual)	76.58	69.94	84.68	74.58	80.75	76.12
QA-GNN (single-step)	79.08	74.43	83.37	74.72	79.60	78.08

Table 2: **Test-set accuracy for different subgraph retrieval methods on FACTKG.** The *direct* approach only includes knowledge triples where both nodes appear in the claim, the *contextual* also includes edges appearing in the claim, and the *single-step* includes all knowledge triples where at least one node appears in the claim. The QA-GNN models used the single-step subgraph if the direct or contextual is empty, while the BERT models did not.

Model	One-hop	Conjunction	Existence	Multi-hop	Negation	Total
	P / R / F1	P / R / F1	P/R/F1	P / R / F1	P / R / F1	P / R / F1
BERT (no subgraphs)	71.89 / 51.66 / 60.12	75.44 / 34.20 / 47.06	59.52 / 73.63 / 65.82	85.19 / 60.90 / 71.03	58.88 / 73.13 / 65.24	75.25 / 54.00 / 62.88
QA-GNN (direct)	76.19/67.04/71.32	80.11 / 51.22 / 62.49	56.19 / 74.10 / 63.91	80.04 / 74.80 / 77.33	70.97 / 73.80 / 72.36	77.21 / 69.01 / 72.88
QA-GNN (contextual)	84.79/61.29/71.15	80.27 / 38.29 / 51.85	81.83 / 88.38 / 84.98	82.31 / 67.17 / 73.98	77.26 / 82.26 / 79.68	84.10/62.78/71.89
QA-GNN (single-step)	82.51 / 70.55 / 76.06	78.89 / 53.95 / 64.08	79.69 / 88.70 / 83.95	78.44 / 73.09 / 75.67	77.06 / 79.10 / 78.07	81.41 / 71.19 / 75.96
BERT (contextual)	83.05 / 75.51 / 79.10	88.60 / 72.56 / 79.78	59.68 / 63.42 / 61.49	84.10 / 70.67 / 76.80	75.84 / 74.46 / 75.15	83.30 / 74.28 / 78.53
BERT (direct)	83.89/71.86/77.41	88.69 / 69.32 / 77.82	58.97 / 54.16 / 56.46	83.38 / 72.91 / 77.80	69.99 / 78.11 / 73.82	83.76 / 72.12 / 77.51
BERT (single-step)	96.27 / 98.29 / 97.27	96.06 / 98.13 / 97.09	96.45 / 98.12 / 97.28	85.31 / 76.59 / 80.72	92.01 / 91.71 / 91.86	93.75 / 92.79 / 93.27

Table 3: Precision (P), Recall (R), and F1 scores for different models and subgraph types on the test-set.

The results can be found in Table 2 and Table 3. We see a clear improvement in BERT when using the direct subgraphs over none, a small improvement when using the contextual subgraphs, and a big improvement when using the single-step method. The same is true for the QA-GNN, but the differences in performance are smaller. The models score the lowest on multi-hop claims.

Since we used non-trainable subgraph retrieval methods and a frozen BERT for embedding the nodes and edges in the subgraphs, we performed this processing before training the models. During training, the models used a lookup table to get the subgraphs and the word embeddings, which significantly decreased the training time. The retrieval of all the subgraphs took about 15 minutes, and the embedding of all the words appearing in them took about 1 hour. We also tried training a QA-GNN without frozen embeddings, but it ran so slow that we were not able to carry out the training with our available computational resources.

4.3 Competitive ChatGPT Performance

The results for the ChatGPT prompting can be found in Table 4. The accuracy is substantially lower than from our best models, but better than the baselines using only the claims. The accuracy is fairly consistent over the three runs, and we do not see a big difference between the amount of questions asked at a time.

Model	Accuracy (mean ± std)
ChatGPT 25 questions	73.67 ± 0.5
ChatGPT 50 questions	76.33 ± 3.3
ChatGPT 100 questions	73.00 ± 1.4

Table 4: **Test-set accuracy for different configurations of ChatGPT prompting.** The metrics are averaged over three runs. The prompts included 25, 50 or 100 claims at a time, but the same claims were used in all of the configurations.

We started with an initial prompt asking for just the truth values for a list of claims, and updated it to also include some training examples and to ask for explanations. Several configurations of the prompt were tested, and it was also improved based on feedback from ChatGPT. 431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

We saw the biggest improvement when we asked for a short explanation of the answers, instead of just the truth values. Without asking for explanations, the amount of answers were often longer or shorter than the amount of questions, but this never happened when explanations were included. We added numbers to the questions to further help with this issue. We also saw a slight improvement by formulating the prompt with bullet point lists and by including some example inputs and outputs from the training set. The final prompt can be found in Figure 3.

422

423

494

425

426

427

428

429

430

Task: Determine the truth value (T in the DBpedia knowledge solely on your pre-trained H Instructions:	True or False) of the following claims based on information verifiable from Wikipedia, as represented graph. Provide your answers without using real-time internet searches or code analysis, relying cnowledge.
• You will evaluate the	following claims, presented one per line.
• Base your answers so	olely on your knowledge as of your last training cut-off.
• Provide answers in P	ython list syntax for easy copying.
• Respond with True f	or verifiable claims, and False otherwise.
• Include a brief explanation	nation for each answer, explaining your reasoning based on your pre-training.
• If the claim is vague False.	or lacks specific information, please make an educated guess on whether it is likely to be True or
Output Format: Format ye explanation). Example Claims: 1. The Atatürk Monument i 2. Yes, Eamonn Butler's al 3. I have heard 300 North I 4. The band Clinton Gregor Example output:	our responses as a list in Python. Each entry should be a tuple, formatted as (claim number, answer is located in Izmir, Turkey, where the capital is Ankara. ma mater is the University of Texas System! _aSalle was completed in 2009. ry created an album in the rock style
[(1, True, "The Atatürk (2, False, "Eamonn B educational background do (3, True, "300 North L (4, False, "Clinton Gro	Monument is indeed located in Izmir, and the capital of Turkey is Ankara."), utler did not attend the University of Texas System; he is a British author and economist whose es not include this institution."), aSalle in Chicago was indeed completed in 2009."), egory is primarily known as a country music artist, not rock."),
] Here are the actual claims	s you should answer:

Figure 3: Final prompt used to get truth values from ChatGPT 40. The actual questions are not included, but were in the format of the Example Claims. The Example Claims are from the training set, and the Example Output is copy pasted from an actual ChatGPT answer.

5 Discussion

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463 464

465

466

467

468

We were able to train better and more efficient models by simplifying the subgraph retrieval methods, both by using a fine-tuned BERT and a slightly modified QA-GNN model. While the QA-GNN models trained the fastest, the fine-tuned BERT clearly performed the best, significantly outperforming the benchmark models. This suggests that the simple logical subgraph retrievals worked better than the complex trained approaches in previous work. We suggest that the performance gain in the claim-only benchmark was due to slightly better hyperparameters.

All of the models performed better the bigger the subgraphs were. This suggests that the model architectures are able to utilize the relevant parts of the subgraphs, without needing an advanced subgraph retrieval step. This is emphasized by our fine-tuned BERT model achieving a 93.49% test set accuracy by only using the single-step subgraphs, while the GEAR model from Kim et al. (2023), which trained two language models to perform graph retrieval, achieved a 77.65% test-set accuracy.

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

When examining the precision and recall metrics in Table 3, we see that most of the models has a higher precision than recall, except for the best performing model, the single-step BERT. However, the single-step BERT does have a lower recall for the multi-hop claims, which it performs significantly worse on. Therefore, the models mostly has a higher precision than recall when their performance is not so good, suggesting they are slightly more likely to predict "false" on claims that they are not confident about.

A limitation of our subgraph retrieval methods is that they never include nodes that are more than one step away from an entity node, while the trained approach from Kim et al. (2023) is dynamic and may include more. This might make the hypothesis that the simple subgraph retrieval methods will

perform worse on *multi-hop* claims than the dy-489 namically trained, however, we see the exact op-490 posite behavior. The best BERT and OA-GNN 491 models score 80.32% and 74.72% at the multi-492 hop claims respectively, while the dynamic GEAR 493 model scores 68.84%, even lower than the models 494 not using the subgraphs at all. We do however see 495 that the best performing BERT model clearly per-496 forms the worst on the multi-hop claims compared 497 to the other claim types, indicating that even bigger 498 subgraphs might be beneficial. 499

> While the sample size for the ChatGPT metrics were small, it does indicate that non-fine-tuned LLMs can achieve adequate few-shot performance compared to a fine-tuned claim-only BERT. The performance does not seem to be substantially compromised when the amount of questions prompted increases, so with a bigger access to computational resources, it might be possible to prompt the full test-set at once. The removal of fine-tuning greatly improves the ease of use if one only needs to verify a few claims. While we are hesitant to make any conclusion with the small sample size, we believe that the results serve as an approximate benchmark of how difficult the dataset is.

502 503

507

508

509

510

511

512

513

514

515

516

517

518

519

520

522

524

525

526

530

531

532

533

534

535

536

538

6 Conclusion and Future Work

Our results show that with simple, yet efficient methods for subgraph retrieval, our models were able to improve fact verification with knowledge graphs with respect to both performance and efficiency. The fine-tuned BERT model with singlestep subgraphs clearly achieves the best performance, while the QA-GNN models are more efficient to train.

This indicates that complex models can work well with simple subgraph retrieval methods. Since the single-step subgraphs mostly contain information not relevant for the claims, the models are themselves able to filter away irrelevant information, and complex subgraph retrieval methods may hence not be necessary for accurate fact verification. Additionally, since the best performing model performed the poorest with multi-hop claims, future research could explore simple subgraphs retrieval methods allowing for bigger depths than one. Additionally, future work should also be directed towards running similar experiments on other datasets.

We also encourage researchers that have access to bigger computational resources to further explore the performance of LLMs for fact verification. A core limitation of our ChatGPT prompting was the inability to use the full test-set, and we consider this crucial for further development. We also think it would be especially interesting to make LLM and KG hybrid models. Since our results indicate that simple single-step subgraph retrievals outperform more complex methods, a promising path of future research would be to simply use both the claims and the single-step subgraphs as input to the LLM. If possible, the LLM could also be fine-tuned on the dataset. We also encourage future work to create fully reproducible results with LLMs, which we were unable to do. 539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

7 Limitations

Our experiments with ChatGPT were done on a small sample of test questions, with a model that was not possible to seed, and therefore is not reproducible. Due to the small sample size, we are not able to directly compare the performance to other approaches. The lack of reproducibility, which stems from the state-of-the-art model that was available to the author is not fully publicly available, makes it impossible for other researchers to completely verify our findings. Additionally, the process for creating prompts were not standardized, we simply tried different configurations based on our own experience with using LLMs until we could not increase the validation accuracy further. Due to these limitations, one should therefore be very hesitant to make any confident conclusions based on the experiments we performed with ChatGPT.

Because our intention was to specifically explore different language models' abilities of fact verification with knowledge graphs on the FACTKG dataset, we did not conduct any experiments on other datasets. It is possible that our results will not be consistent with other datasets.

Additionally, our selection of models and hyperparameter settings could be more diverse. Due to computational constraints, we did not perform a grid search for hyperparameters, but tuned hyperparameters one by one. Which parameters we searched for were not decided in advance. A predefined grid search might lead to a fairer and more reproducible approach. We did not experiment with different orderings of the knowledge triples for the fine-tuned BERT models, which could influence the performance.

References

588

592

594 595

596

597

598

610

611

612

613

615

617

619

625

631

638

641

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Giannis Bekoulis, Christina Papagiannopoulou, and Nikos Deligiannis. 2021. A review on fact extraction and verification. *ACM Computing Surveys (CSUR)*, 55(1):1–35.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. Natural language processing with Python: analyzing text with the natural language toolkit. " O'Reilly Media, Inc.".
- Jiangui Chen, Ruqing Zhang, Jiafeng Guo, Yixing Fan, and Xueqi Cheng. 2022. Gere: Generative evidence retrieval for fact verification. In *Proceedings of the* 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 2184–2189.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. 2019. Tabfact: A largescale dataset for table-based fact verification. arXiv preprint arXiv:1909.02164.
- S Cohen, C Li, J Yang, and C Yu. 2011. Computational journalism: A call to arms to database researchers, 148-151. In 5th Biennial Conference on Innovative Data Systems Research, CIDR.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Sushant Gautam. 2024. Factgenius: Combining zeroshot prompting and fuzzy relation mining to improve fact verification with knowledge graphs. *arXiv preprint arXiv*:2406.01311.
- Andreas Hanselowski, Avinesh PVS, Benjamin Schiller, Felix Caspelherr, Debanjan Chaudhuri, Christian M Meyer, and Iryna Gurevych. 2018. A retrospective analysis of the fake news challenge stance detection task. *arXiv preprint arXiv:1806.05180*.
- Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. 2020. Array programming with numpy. *Nature*, 585(7825):357–362.
- Naeemul Hassan, Bill Adair, James T Hamilton, Chengkai Li, Mark Tremayne, Jun Yang, and Cong Yu. 2015. The quest to automate fact-checking. In *Proceedings of the 2015 computation+ journalism symposium*. Citeseer.
- Christopher Hidey, Tuhin Chakrabarty, Tariq Alhindi, Siddharth Varia, Kriste Krstovski, Mona Diab, and

Smaranda Muresan. 2020. Deseption: Dual sequence prediction and adversarial examples for improved fact-checking. *arXiv preprint arXiv:2004.12864*.

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

- Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.
- Sergey Ioffe and Christian Szegedy. 2015. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr.
- Jiho Kim, Sungjin Park, Yeonsu Kwon, Yohan Jo, James Thorne, and Edward Choi. 2023. Factkg: Fact verification via reasoning on knowledge graphs. *arXiv preprint arXiv:2305.06590*.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Shreyash Mishra, S Suryavardan, Amrit Bhaskar, Parul Chopra, Aishwarya N Reganti, Parth Patwa, Amitava Das, Tanmoy Chakraborty, Amit P Sheth, Asif Ekbal, et al. 2022. Factify: A multi-modal fact verification dataset. In *DE-FACTIFY*@ AAAI.
- Yixin Nie, Haonan Chen, and Mohit Bansal. 2019. Combining fact extraction and verification with neural semantic matching networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 6859–6866.
- Open AI. 2024. Hello gpt 40. https://openai.com/ index/hello-gpt-40/, Accessed 30.05.2024.
- Jungsoo Park, Sewon Min, Jaewoo Kang, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2021. Faviq: Fact verification from information-seeking questions. *arXiv preprint arXiv:2107.02153*.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. Get your vitamin c! robust fact verification with contrastive evidence. *arXiv preprint arXiv:2103.08541*.

Tal Schuster, Darsh J Shah, Yun Jie Serene Yeo, Daniel Filizzola, Enrico Santus, and Regina Barzilay. 2019. Towards debiasing fact verification models. arXiv preprint arXiv:1908.05267.

699

705

711

712

713

714

715

716

717

718

719

720

721

723

730

731

732

733

734 735

736

738

740

741

742

744

745

747

748

- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958.
 - James Thorne and Andreas Vlachos. 2018. Automated fact checking: Task formulations, methods and future directions. *arXiv preprint arXiv:1806.07687*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. Fever: a large-scale dataset for fact extraction and verification. *arXiv preprint arXiv:1803.05355*.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45.
- Michihiro Yasunaga, Hongyu Ren, Antoine Bosselut, Percy Liang, and Jure Leskovec. 2021. Qagnn: Reasoning with language models and knowledge graphs for question answering. *arXiv preprint arXiv:2104.06378*.
- Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019. Gear: Graph-based evidence aggregating and reasoning for fact verification. *arXiv preprint arXiv:1908.01843*.
- Dimitrina Zlatkova, Preslav Nakov, and Ivan Koychev. 2019. Fact-checking meets fauxtography: Verifying claims about images. *arXiv preprint arXiv:1908.11722.*

A Hyperparameter Details

We used an AdamW optimizer (Loshchilov and Hutter, 2017) and a linear learning rate scheduler with 50 warm up steps. We used the model from the epoch with lowest accuracy loss. The hyperparameters were tuned in a line search, first testing different learning rates, and testing all the other hyperparameters with the best learning rate. We searched for learning rates in $\{1e - 3, 5e - 4, 1e - 4, 5e - 5, 1e - 5\}$ for all models. We initially set the batch size to 32, except for the BERT models with large subgraphs, which were set to 4 due to

Model	Learning Rate	Batch Size	Best Epoch
BERT (no subgraphs)	1e-4	32	6
BERT (direct)	1e-4	32	7
BERT (contextual)	5e-5	8	7
BERT (single-step)	5e-5	4	7
QA-GNN (direct)	1e-4	128	8
QA-GNN (contextual)	5e-5	64	17
QA-GNN (single-step)	1e-5	128	20

Table 5: **Final hyperparameters for the different models.** The direct QA-GNN model used GNN and classifier dropout rates of 0.3 and 0.3, while both the two other QA-GNN used 0.1 and 0.5.

memory constraints. After finding the learning rate, we tried batch sizes in $\{32, 64, 128, 256\}$. For the QA-GNN model, we initially set the GNN dropout and the classifier dropout to 0.3, and tried values in $\{0, 0.1, 0.3, 0.5, 0.6\}$. We also tried to freeze some of the layers in the base model, and to use a RoBERTa (Liu et al., 2019) instead of BERT (Devlin et al., 2018), but neither of these approaches improved the validation loss. 749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

The final hyperparameters can be found in Table A.

B Scientific Artifacts

We conducted the experiments using several python libraries, including PyTorch version 2.0.1 (Paszke et al., 2019) with CUDA version 11.7, Hugging-Face Transformers (Wolf et al., 2020), NumPy (Harris et al., 2020), SpaCy (Honnibal and Montani, 2017) and NLTK (Bird et al., 2009). We will make all the code used for this paper publicly available.