

# Zero-shot Fact Verification by Claim Generation

Liangming Pan<sup>1,2</sup> Wenhuchen Chen<sup>3</sup> Wenhan Xiong<sup>3</sup>

Min-Yen Kan<sup>2</sup> William Yang Wang<sup>3</sup>

<sup>1</sup>NUS Graduate School for Integrative Sciences and Engineering

<sup>2</sup>School of Computing, National University of Singapore, Singapore

<sup>3</sup>University of California, Santa Barbara, CA, USA

liangmingpan@u.nus.edu

{wenhuchen, xwhan, william}@cs.ucsb.edu

kanmy@comp.nus.edu.sg

## Abstract

Neural models for automated fact verification have achieved promising results thanks to the availability of large, human-annotated datasets. However, for each new domain that requires fact verification, creating a dataset by manually writing claims and linking them to their supporting evidence is expensive. We develop QACG, a framework for training a robust fact verification model by using automatically-generated claims that can be supported, refuted, or unverifiable from evidence from Wikipedia. QACG generates question-answer pairs from the evidence and then convert them into different types of claims. Experiments on the FEVER dataset show that our QACG framework significantly reduces the demand for human-annotated training data. In a zero-shot scenario, QACG improves a RoBERTa model’s  $F_1$  from 50% to 77%, equivalent in performance to 2K+ manually-curated examples. Our QACG code is publicly available.<sup>1</sup>

## 1 Introduction

Fact verification aims to validate a claim in the context of evidence. This task has attracted growing interest with the rise in disinformation in news and social media. Rapid progress has been made by training large neural models (Zhou et al., 2019; Liu et al., 2020b; Zhong et al., 2020) on the FEVER dataset (Thorne et al., 2018), containing more than 100K human-crafted (evidence, claim) pairs based on Wikipedia.

Fact verification is demanded in many domains, including news articles, social media, and scientific documents. However, it is not realistic to assume that large-scale training data is available for every new domain that requires fact verification. Creating training data by asking humans to write claims and

search for evidence to support/refute them can be extremely costly.

We address this problem by exploring the possibility of automatically *generating* large-scale (evidence, claim) pairs to train the fact verification model. We propose a simple yet general framework **Question Answering for Claim Generation (QACG)** to generate three types of claims from any given evidence: 1) claims that are supported by the evidence, 2) claims that are refuted by the evidence, and 3) claims that the evidence does Not have Enough Information (NEI) to verify.

To generate claims, we utilize *Question Generation (QG)* (Zhao et al., 2018; Liu et al., 2020a; Pan et al., 2020), which aims to automatically ask questions from textual inputs. QG has been shown to benefit various NLP tasks, such as enriching QA corpora (Alberti et al., 2019), checking factual consistency for summarization (Wang et al., 2020), and data augmentation for semantic parsing (Guo et al., 2018). To the best of our knowledge, we are the first to employ QG for fact verification.

As illustrated in Figure 1, given a passage  $P$  as the evidence, we first employ a *Question Generator* to generate a question-answer pair  $(Q, A)$  for the evidence. We then convert  $(Q, A)$  into a claim  $C$  (*QA-to-Claim*) based on the following logical assumptions: a) if  $P$  can answer  $Q$  and  $A$  is the correct answer, then  $C$  is a supported claim; b) if  $P$  can answer  $Q$  but  $A$  is an incorrect answer, then  $C$  is a refuted claim; c) if  $P$  cannot answer  $Q$ , then  $C$  is a NEI claim. The Question Generator and the QA-to-Claim model are off-the-shelf BART models (Lewis et al., 2020), finetuned on SQuAD (Rajpurkar et al., 2016) and QA2D (Demszky et al., 2018) datasets.

We generate 100K (evidence, claim) pairs for each type of claim, which we then use to train a RoBERTa (Liu et al., 2019) model for fact verification. We evaluate the model on three test sets

<sup>1</sup><https://github.com/teacherpeterpan/Zero-shot-Fact-Verification>

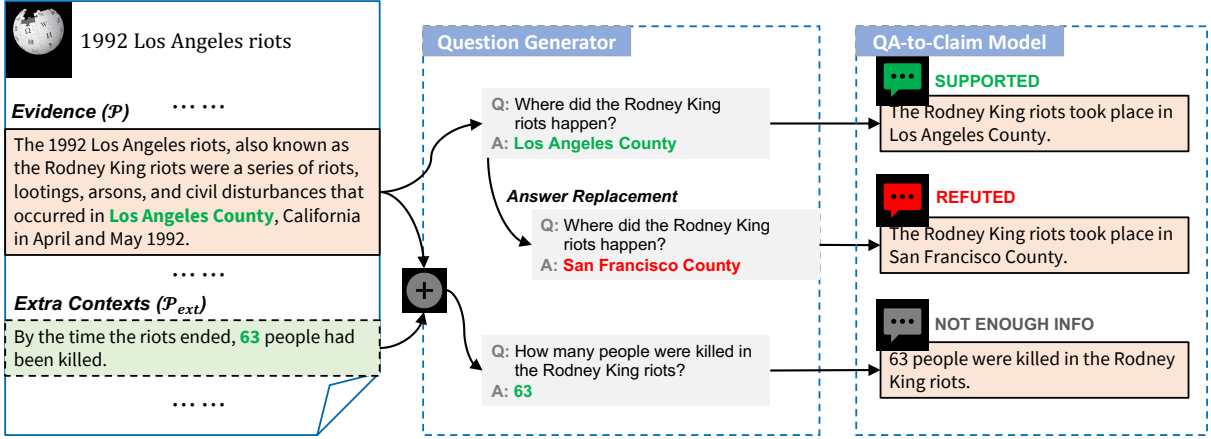


Figure 1: Overview of our QACG framework, consisting of two modules: 1) *Question Generator* generates questions from the evidence  $\mathcal{P}$  and the extra contexts  $\mathcal{P}_{ext}$  given different answers extracted from the passage (in green), and 2) *QA-to-Claim* converts question-answer pairs into claims with different labels.

based on the FEVER dataset. Although we do not use any human-labeled training examples, the model achieves over 70% of the  $F_1$  performance of a fully-supervised setting. By finetuning the model with only 100 labeled examples, we further close the performance gap, achieving 89.1% of fully-supervised performance. The above results show that pretraining the fact verification model with generated claims greatly reduces the demand for in-domain human annotation. When evaluating the model on an unbiased test set for FEVER, we find that training with generated claims also produces a more *robust* fact verification model.

In summary, our contributions are:

- To the best of our knowledge, this is the first work to investigate zero-shot fact verification.
- We propose QACG, a novel framework to generate high-quality claims via question generation.
- We show that the generated training data can greatly benefit the fact verification system in both zero-shot and few-shot learning settings.

## 2 Methodology

Given a claim  $\mathcal{C}$  and a piece of evidence  $\mathcal{P}$  as inputs, a *fact verification* model  $\mathcal{F}$  predicts a label  $\mathcal{Y} \in \{\text{supported}, \text{refuted}, \text{NEI}\}$  to verify whether  $\mathcal{C}$  is supported, refuted, or can not be verified by the information in  $\mathcal{P}$ .

For the *zero-shot* setting, we assume no human-annotated training example is available. Instead, we generate a synthetic training set based on our QACG framework to train the model.

### 2.1 Question Generator and QA-to-Claim

As illustrated in Figure 1, our claim generation model QACG has two major components: a *Question Generator*  $\mathcal{G}$ , and a *QA-to-Claim* model  $\mathcal{M}$ .

The **Question Generator** takes as input an evidence  $\mathcal{P}$  and a text span  $A$  from the given evidence and aims to generate a question  $Q$  with  $A$  as the answer. We implement this with the BART model (Lewis et al., 2020), a large transformer-based sequence-to-sequence model pretrained on 160GB of text. The model is finetuned on the SQuAD dataset processed by Zhou et al. (2017), where the model encodes the concatenation of the SQuAD passage and the answer text and then learns to decode the question. We evaluate the question generator using automatic and human evaluation and investigate its impact on fact verification in Appendix A.

The **QA-to-Claim Model** takes as inputs  $Q$  and  $A$ , and outputs the declarative sentence  $C$  for the  $(Q, A)$  pair, as shown in Figure 1. We also treat this as a sequence-to-sequence problem and fine-tune the BART (Lewis et al., 2020) model on the QA2D dataset (Demszky et al., 2018), which contains the human-annotated declarative sentence for each  $(Q, A)$  pair in SQuAD.

### 2.2 Claim Generation

Given the pretrained question generator  $\mathcal{G}$  and the QA-to-Claim model  $\mathcal{M}$ , we then formally introduce how we generate claims with different labels.

**Supported claim generation.** Given an evidence  $\mathcal{P}$ , we use named entity recognition to identify all entities within  $\mathcal{P}$ , denoted as  $\mathcal{E}$ . For each

entity  $a \in \mathcal{E}$ , we treat each  $a$  in turn as an answer and generate a question  $q = \mathcal{G}(\mathcal{P}, a)$  with the question generator. The question-answer pair  $(q, a)$  are then sent to the QA-to-Claim model to generate the supported claim  $c = \mathcal{M}(q, a)$ .

**Refuted claim generation.** To generate a refuted claim, after we generate the question-answer pair  $(q, a)$ , we use *answer replacement* (shown in Figure 1) to replace the answer  $a$  with another entity  $a'$  with the same type such that  $a'$  becomes an incorrect answer to the question  $q$ . Using  $a$  as the query, we randomly sample a phrase from the top-5 most similar phrases in the pretrained Sense2Vec (Trask et al., 2015) as the replacing answer  $a'$ . The new pair  $(q, a')$  is then fed to the QA-to-Claim model to generate the refuted claim.

To avoid the case that  $a'$  is still the correct answer, we define rules to ensure that the  $a'$  has less lexical overlap with  $a$ . However, this problem is sometimes non-trivial and cannot be completely avoided. For example, for the QA pair: (“Who is the producer of Avatar?”; “James Cameron”), another valid answer  $a'$  is “Jon Landau”, who happens to be another producer of Avatar. However, we observe that such coincidences rarely happen: among the 100 randomly sampled claims, we only observed 2 such cases. Therefore, we leave them as the natural noise of the generation model.

**NEI claim generation.** We need to generate a question  $q'$  which is relevant but cannot be answered by  $\mathcal{P}$ . To this end, we link  $\mathcal{P}$  back to its original Wikipedia article  $\mathcal{W}$  and expand the evidence with additional contexts  $\mathcal{P}_{ext}$ , which are five randomly-retrieved sentences from  $\mathcal{W}$  that are not present in  $\mathcal{P}$ . In our example in Figure 1, one additional context retrieved is “By the time the riots ended, 63 people had been killed”. We then concatenate  $\mathcal{P}$  and  $\mathcal{P}_{ext}$  as the expanded evidence, based on which we generate a supported claim given an entity in  $\mathcal{P}_{ext}$  as the answer (e.g., “63”). This results in a claim relevant to but unverifiable by the original evidence  $\mathcal{P}$ .

### 3 Experiments

By applying our QACG model to each of the 18,541 Wikipedia articles in the FEVER training set, we generate a total number of 176,370 supported claims, 360,924 refuted claims, and 258,452 NEI claims. Our generated data is around five times the size of the human-annotated

claims in FEVER. We name this generated dataset as QACG-*Full*. We then create a balanced dataset QACG-*Filtered* by randomly sampling 100,000 samples for each class. Statistics of the FEVER and the generated dataset are in Appendix B.

**Evaluation Datasets.** We evaluate fact verification on three different test sets based on FEVER: **1) FEVER-S/R:** Since only the supported and refuted claims are labeled with gold evidence in FEVER, we take the claim-evidence pairs of these two classes from the FEVER test set for evaluation. **2) FEVER-Symmetric:** this is a carefully-designed unbiased test set designed by Schuster et al. (2019) to detect the robustness of the fact verification model. Note that only supported and refuted claims are present in this test set. **3) FEVER-S/R/N:** The full FEVER test set are used for a three-class verification. We follow Atanasova et al. (2020) to use the system of Malon (2019) to retrieve evidence sentences for NEI claims.

**Fact Verification Models.** As shown in Table 1, we take a BERT model (S1) and a RoBERTa model (S2) fine-tuned on the FEVER training set as the *supervised* models. Their corresponding *zero-shot* settings are Rows U5 and U6, where the models are trained on our generated QACG-*Filtered* dataset. Note that for binary classification (FEVER-S/R and FEVER-Symmetric), only the supported and refuted claims are used for training, while for FEVER-S/R/N, the full training set is used.

We employ four baselines that also do not need any human-annotated claims to compare with our method. *Random Guess* (U1) is a weak baseline that randomly predicts the class label. *GPT2 Perplexity* (U2) predicts the class label based on the perplexity of the claim under a pretrained GPT2 (Radford et al., 2019) language model, following the assumption that “misinformation has high perplexity” (Lee et al., 2020a). *MNLI-Transfer* (U3) trains a BERT model for natural language inference on the MultiNLI corpus (Williams et al., 2018) and applies it for fact verification. *LM as Fact Checker* (Lee et al., 2020b) (U4) leverages the implicit knowledge stored in the pretrained BERT language model to verify a claim. The implementation details are given in Appendix C.

#### 3.1 Main Results

Table 1 summarizes the fact verification performance, measured by the macro Precision ( $P$ ), Recall ( $R$ ), and F1 Score ( $F_1$ ).

Model		FEVER -Symmetric	FEVER-S/R	FEVER-S/R/N
		$P / R / F_1$	$P / R / F_1$	$P / R / F_1$
Supervised	S1. BERT-base (Devlin et al., 2019)	81.5 / 81.3 / 81.2	92.8 / 92.6 / 92.6	85.7 / 85.6 / 85.6
	S2. RoBERTa-large (Liu et al., 2019)	<b>85.5 / 85.5 / 85.5</b>	<b>95.2 / 95.1 / 95.1</b>	<b>88.0 / 87.9 / 87.8</b>
Zero-shot	U1. Random Guess	50.0 / 50.0 / 50.0	50.0 / 50.0 / 50.0	33.3 / 33.3 / 33.3
	U2. GPT2 Perplexity	52.7 / 52.7 / 52.7	55.6 / 55.6 / 55.6	35.3 / 35.3 / 35.3
	U3. MNLI-Transfer	62.2 / 55.5 / 58.7	63.6 / 60.5 / 61.8	41.4 / 39.6 / 40.7
	U4. LM as Fact Checker (Lee et al., 2020b)	71.2 / 64.5 / 67.8	77.9 / 65.6 / 70.2	64.3 / 54.6 / 49.8
	U5. QACG (BERT-base)	73.2 / 73.0 / 72.9	74.2 / 74.0 / 74.1	56.5 / 55.7 / 55.9
	U6. QACG (RoBERTa-large)	<b>77.3 / 77.0 / 77.1</b>	<b>78.1 / 78.1 / 78.1</b>	<b>64.6 / 62.0 / 62.6</b>

Table 1: Fact verification performance for supervised models and zero-shot models on three different settings.

**Comparison with supervised settings.** The zero-shot setting with RoBERTa-large (U6) attains 78.1  $F_1$  on the FEVER-S/R and 62.6  $F_1$  on the FEVER-S/R/N. The  $F_1$  gap to the fully-supervised RoBERTa-large (S2) is only 17.0 and 15.2 on these two settings, respectively. These results demonstrate the effectiveness of QACG in generating good (evidence, claim) pairs for training the fact verification model. The RoBERTa model (S2, U6) is more effective than the BERT model (S1, U5) for both the zero-shot and the supervised setting.

**Comparison with zero-shot baselines.** Our model (U6) achieves the best results among all the zero-shot baselines across all three test sets. We find that validating a claim by its perplexity (U2) only works slightly better than random guess (U1) (+3.43  $F_1$ ), showing that misinformation does not necessarily have high perplexity. Although natural language inference seems highly correlated with fact verification, directly transferring the model trained on the MNLI dataset (U3) only outperforms random guess by 9.30  $F_1$ . We believe this is due to the domain gap between FEVER (from Wikipedia) and the MNLI (from fiction, letters, etc.) dataset. As a generation framework, our model can avoid the domain gap issue by generating pseudo training data from the same domain (Wikipedia). Another reason is the “task gap” between NLI and fact verification, in which the former makes inference about the situation described in a sentence, while the latter focuses on claims about entities in Wikipedia.

**Model Robustness.** We observe a large performance drop when the supervised model is evaluated on the FEVER-Symmetric test set for both the BERT model (−11.4  $F_1$ ) and the RoBERTa model (−9.6  $F_1$ ). However, the models trained with our generated data (U2, U3) drop only 1.2 and 1.0  $F_1$  drop. This suggests that the wide range of different claims we generate as training data helps eliminate

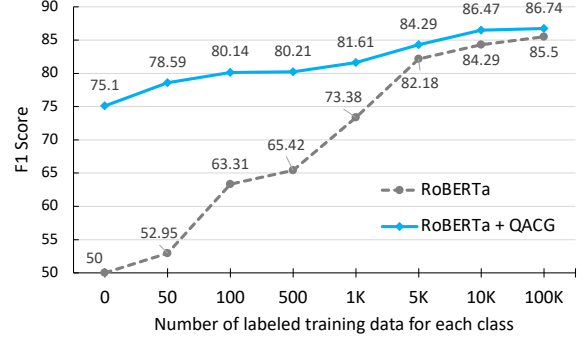


Figure 2: The few-shot learning experiment. The figure shows the  $F_1$  score on FEVER-Symmetric for progressively larger training dataset sizes.

some of the annotation artifacts present in FEVER, leading to a more robust fact verification model.

### 3.2 Few-shot Fact Verification

We then explore QACG’s effectiveness in the few-shot learning setting where only a few human-labeled (evidence, claim) pairs are available. We first train the RoBERT-large fact verification model with our generated dataset QACG-Filtered. Then we fine-tune the model with a limited amount of human-labeled claims in FEVER. The blue solid line in Figure 2 shows the  $F_1$  scores on FEVER-Symmetric after finetuning with different numbers of labeled training data. We compare this with training the model from scratch with the human-labeled data (grey dashed line).

Our model performs consistently better than the model without pretraining, regardless of the amount of labeled training data. The improvement is especially prominent in data-poor regimes; for example, our approach achieves 78.6  $F_1$  with only 50 labeled claims for each class, compared with 52.9  $F_1$  without pretraining (+25.7). This only leaves a 7.9  $F_1$  gap to the fully-supervised setting (86.5  $F_1$ ) with over 100K training samples. The results show pretraining fact verification with QACG



Evidence	Generated Claim
<b>Budapest</b> is cited as one of the most beautiful cities in <b>Europe</b> , ranked as the most liveable Central and <b>Eastern European</b> city on <b>EIU</b> ’s quality of life index, ranked as “the world’s <b>second</b> best city” by <b>Conde Nast Traveler</b> , and “ <b>Europe</b> ’s <b>7th</b> most idyllic place to live” by <b>Forbes</b> .	<b>SUPPORTED claims</b> <b>Budapest</b> is ranked as the most liveable city in central Europe. Budapest ranks <b>7th</b> in terms of idyllic places to live in Europe. <b>REFUTED claims</b> Budapest ranks in <b>11th</b> in terms of idyllic places to live in Europe. Budapest is ranked the most liveable city in <b>Asia</b> . <b>NEI claims</b> Budapest is one of the largest cities in the European Union. Budapest is the capital of Hungary.
<b>Alia Bhatt</b> received critical acclaim for portraying emotionally intense characters in the road drama <b>Highway</b> (2014), which won her the <b>Filmfare Critics Award for Best Actress</b> , and the crime drama <b>Udta Punjab</b> (2016), which won her the <b>Filmfare Award for Best Actress</b> .	<b>SUPPORTED claims</b> Bhatt won the <b>Filmfare Award for Best Actress</b> in <b>Udta Punjab</b> . Bhatt received the Filmfare Critics Award for her role in <b>Highway</b> . <b>REFUTED claims</b> Alia Bhatt won the <b>Best Original Screenplay</b> award in <b>Highway</b> . <b>2 States</b> (2014) won Alia Bhatt the Filmfare Award for Best Actress. <b>NEI claims</b> Alia Bhatt made her acting debut in the 1999 thriller <b>Sangharsh</b> . Bhatt played her first leading role in <b>Karan Johar</b> ’s romantic drama.

Table 2: Examples of evidence and claims generated by QACG, categorized by class labels. In the evidence, the identified answers for question generation are highlighted in blue. For claims, the correct answers are highlighted in blue for SUPPORTED claims and the replaced wrong answers are in red for REFUTED claims.

<b>Evidence:</b>	Roman Atwood is best known for his vlogs, where he posts updates about his life.
<b>Claim:</b>	Roman Atwood is <u>a content creator</u> .
<b>Evidence:</b>	In 2004, Slovenia <u>entered NATO and the European Union</u> .
<b>Claim:</b>	Slovenia <u>uses the euro</u> .
<b>Evidence:</b>	He has traveled to <u>Chad and Uganda</u> to raise awareness about conflicts in the regions.
<b>Claim:</b>	Ryan Gosling has been to <u>a country in Africa</u> .

Table 3: Examples of claims in FEVER that require commonsense or world knowledge (underlined).

greatly reduces the demand for in-domain human-annotated data. Our method can provide a “warm start” for fact verification system when applied to a new domain where training data are limited.

### 3.3 Analysis of Generated Claims

Table 2 shows representative claims generated by our model. The claims are fluent, label-cohesive, and exhibit encouraging language variety. However, one limitation is that our generated claims are mostly *lack of deep reasoning over the evidence*. This is because we finetune the question generator on the SQuAD dataset, in which more than 80% of its questions are shallow factoid questions.

To better understand whether this limitation brings a domain gap between the generated claims and the human-written claims, we randomly sampled 100 supported claims and 100 refuted and analyze whether reasoning is involved to verify those claims. We find that 38% of the supported

claims and 16% of the refuted claims in FEVER require either commonsense reasoning or world knowledge to verify. Table 3 show three typical examples. Therefore, we believe this domain gap is the main bottleneck of our system. Future studies are required to generate more complex claims which involves multi-hop, numerical, and commonsense reasoning, such that we can apply our model to more complex fact checking scenario.

## 4 Conclusion and Future Work

We utilize the question generation model to ask different questions for given evidence and convert question-answer pairs into claims with different labels. We show that the generated claims can train a well-performing fact verification model in both the zero-shot and the few-shot learning setting. Potential future directions could be: 1) generating more complex claims that require deep reasoning; 2) extending our framework to other fact checking domains beyond Wikipedia, *e.g.*, news, social media; 3) leveraging generated claims to improve the robustness of fact checking systems.

## Acknowledgments

This research is supported by the National Research Foundation, Singapore under its International Research Centres in Singapore Funding Initiative. The UCSB authors are not supported by any of the projects above. They thank Google, Amazon, Facebook, and JP Morgan for their generous support.

## Ethical Considerations

We discuss two potential issues of claim generation, showing how our work sidesteps these issues. While individuals may express harmful or biased claims, our work only focuses on generating factoid claims from a corpus. In this work, we take Wikipedia as the source for objective fact. Practicing this technique thus requires the identification of an appropriate source of objective truth to generate claims from. Another potential misuse of claim generation is to generate `refuted` claims and subsequently spread such misinformation. We caution practitioners to treat the generated claims with care. In our case, we use the generated claims only to optimize for the downstream fact verification task. We advise against releasing generated claims for public use — especially on public websites, where they may be crawled and then subsequently used for inference. As such, we will release the model code but not the output in our work. Practitioners can re-run the training pipeline to replicate experiments accordingly.

## References

- Chris Alberti, Daniel Andor, Emily Pitler, Jacob Devlin, and Michael Collins. 2019. Synthetic QA corpora generation with roundtrip consistency. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 6168–6173.
- Pepa Atanasova, Dustin Wright, and Isabelle Augenstein. 2020. Generating label cohesive and well-formed adversarial claims. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3168–3177.
- Dorottya Demszky, Kelvin Guu, and Percy Liang. 2018. Transforming question answering datasets into natural language inference datasets. *CoRR*, abs/1809.02922.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, pages 4171–4186.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. In *Annual Conference on Neural Information Processing Systems (NeurIPS)*, pages 13042–13054.
- Han Guo, Ramakanth Pasunuru, and Mohit Bansal. 2018. Soft layer-specific multi-task summarization with entailment and question generation. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 687–697.
- Alon Lavie and Abhaya Agarwal. 2007. METEOR: an automatic metric for MT evaluation with high levels of correlation with human judgments. In *Proceedings of the Second Workshop on Statistical Machine Translation (WMT@ACL)*, pages 228–231.
- Nayeon Lee, Yejin Bang, Andrea Madotto, and Pascale Fung. 2020a. Misinformation has high perplexity. *CoRR*, abs/2006.04666.
- Nayeon Lee, Belinda Z. Li, Sinong Wang, Wen-tau Yih, Hao Ma, and Madian Khabsa. 2020b. Language models as fact checkers? *CoRR*, abs/2006.04102.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 7871–7880.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. *Text Summarization Branches Out*.
- Bang Liu, Haojie Wei, Di Niu, Haolan Chen, and Yancheng He. 2020a. Asking questions the human way: Scalable question-answer generation from text corpus. In *International World Wide Web Conference (WWW)*.
- 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. *CoRR*, abs/1907.11692.
- Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2020b. Fine-grained fact verification with kernel graph attention network. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 7342–7351.
- Christopher Malon. 2019. Team papelo: Transformer networks at FEVER. *CoRR*, abs/1901.02534.
- Liangming Pan, Yuxi Xie, Yansong Feng, Tat-Seng Chua, and Min-Yen Kan. 2020. Semantic graphs for generating deep questions. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1463–1475.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 311–318.

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100, 000+ questions for machine comprehension of text. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2383–2392.
- Tal Schuster, Darsh J. Shah, Yun Jie Serene Yeo, Daniel Filizzola, Enrico Santus, and Regina Barzilay. 2019. Towards debiasing fact verification models. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3417–3423.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and verification. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, pages 809–819.
- Andrew Trask, Phil Michalak, and John Liu. 2015. sense2vec - A fast and accurate method for word sense disambiguation in neural word embeddings. *CoRR*, abs/1511.06388.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 5008–5020.
- Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT)*, pages 1112–1122.
- Yao Zhao, Xiaochuan Ni, Yuanyuan Ding, and Qifa Ke. 2018. Paragraph-level neural question generation with maxout pointer and gated self-attention networks. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3901–3910.
- WanJun Zhong, Jingjing Xu, Duyu Tang, Zenan Xu, Nan Duan, Ming Zhou, Jiahai Wang, and Jian Yin. 2020. Reasoning over semantic-level graph for fact checking. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 6170–6180.
- 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. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 892–901.
- Qingyu Zhou, Nan Yang, Furu Wei, Chuanqi Tan, Hangbo Bao, and Ming Zhou. 2017. Neural question generation from text: A preliminary study. In *CCF International Conference of Natural Language Processing and Chinese Computing (NLPCC)*, pages 662–671.

## A Evaluation of Question Generation

To implement the question generator, we finetune the pretrained BART model provided by HuggingFace library on the SQuAD dataset. The codes are based on the SimpleTransformers<sup>2</sup> library. The success of our QACG framework heavily rely on whether we can generate fluent and answerable questions given the evidence. Therefore, we separately evaluate the question generator using both automatic and human evaluation and investigate its impact to zero-shot fact verification.

### A.1 Automatic Evaluation

We employ BLEU-4 (Papineni et al., 2002), METEOR (Lavie and Agarwal, 2007), and ROUGE-L (Lin, 2004) to evaluate the performance of our implementation. We compare the BART model with several state-of-the-art QG models, using their reported performance on the Zhou split of SQuAD.

Table 4 shows the evaluation results comparing against all baseline methods. The BART model achieves a BLEU-4 of 21.32, outperforming NQG++, S2ga-mp-gsa, and CGC-QG by large margins. This is as expected since these three baselines are based on Seq2Seq and do not apply language model pretraining. Compared with the current state-of-the-art model UniLM, the BART model achieves comparable results, with slightly lower BLEU-4 but higher METEOR.

Model	B4	MR	$R_L$
NQG++ (Zhou et al., 2017)	13.5	18.2	41.6
S2ga-mp-gsa (Zhao et al., 2018)	15.8	19.7	44.2
CGC-QG (Liu et al., 2020a)	17.6	21.2	44.5
UniLM (Dong et al., 2019)	<b>23.8</b>	25.6	<b>52.0</b>
BART (Lewis et al., 2020)	21.3	<b>27.1</b>	43.6

Table 4: Performance evaluation of the *Question Generator* with different model implementations. We adopt the BART model in our QACG framework.  $B4$ : BLEU-4,  $MR$ : METEOR,  $R_L$ : ROUGE-L.

### A.2 Impact of Answerability

Given the evidence  $P$  and the answer  $A$ , the generated question  $Q$  must be answerable by  $P$  and

<sup>2</sup><https://github.com/ThilinaRajapakse/simpletransformers>

Model	Answerable Rate	FV Performance $P / R / F_1$
NQG++	63.0%	62.2 / 62.4 / 62.3
BART	89.5%	76.3 / 76.0 / 76.1

Table 5: *Answerable Rate*: the ratio of answerable questions generated by the NQG++ and the BART model. *FV Performance*: the zero-shot fact verification performance on the FEVER-Symmetric.

take  $A$  as its correct answer. This is the premise of generating a correct SUPPORTED claim. Therefore, we specially evaluate this *answerability* property via human ratings. We randomly sample 100 generated question-answer pairs with their corresponding evidence and ask two workers to judge the answerability of each sample. We do this for both the NQG++ model and the BART model. To investigate the impact of question quality on the fact verification performance, we separately use the NQG++ and BART as the question generator to generate claims and train the RoBERTa model. The performance is summarized in Table 5.

We find that the ratio of answerable questions generated by the BART model is 89.5%, significantly outperforms the 63.5% achieved by the NQG++ model. When switching the question generator to NQG++, the fact verification  $F_1$  drops to 62.3 (−22.1% compared with BART). This shows that answerability plays an important role in ensuring the validity of the generated claims and has a huge impact on the fact verification performance.

## B Dataset Statistics

Table 6 shows the basic data statistics of the FEVER, FEVER-Symmetric, and our generated dataset by QACG. We use the balanced dataset QACG-Filtered sampled from QACG-Full to train the fact verification model in the zero/few-shot setting. Compared with the original FEVER dataset, our generated QACG-Filtered dataset has a balanced number of claims for each class. Moreover, because QACG can generate three different types of claims for the same given evidence (shown in Figure 1), it results in a more “unbiased” dataset in which the model must rely on the (*evidence*, *claim*) pair rather than the *evidence* itself to make an inference of the class label.

## C Model Implementation Details

**BERT-base and RoBERTa-large (S1, S2, U5, U6).** We use the bert-base-uncased

Dataset	Supported	Refuted	NEI
FEVER Train	80,035	29,775	35,517
FEVER Test	6,666	6,666	6,666
FEVER-Symmetric	710	710	—
QACG Full	176,370	360,924	258,452
QACG Filtered	100,000	100,000	100,000

Table 6: Basic statistics of the FEVER dataset and the dataset generated by QACG.

(110M parameters) and the roberta-large (355M parameters) model provided by HuggingFace library to implement the BERT model and the RoBERTa model, respectively. The model is fine-tuned with a batch size of 16, learning rate of 1e-5 and for a total of 5 epochs, where the epoch with the best performance is saved.

**GPT2 Perplexity (U2).** To measure the perplexity, we use the HuggingFace implementation of the medium GPT-2 model (gpt2-medium, 345M parameters). We then rank the claims in the FEVER test set by their perplexity under the GPT-2 model. We then predict the label for each claim based on the assumption that misinformation has high perplexity. However, manually setting the perplexity threshold is difficult. Since the FEVER test set contains an equal number of claims for each class, we predict the claims in the top 1/3 of the ranking list as *refuted*, and the bottom 1/3 as *supported*. The rest claims are set as *NEI*. Therefore, the number of predicted labels for each class is also equal.

**MNLI-Transfer (U3).** We use the HuggingFace – BERT base model (110M parameters) fine tuned on the Multi-Genre Natural Language Inference (MNLI) corpus<sup>3</sup>, a crowd-sourced collection of 433K sentence pairs annotated with textual entailment information. We then directly apply this model for fact verification in the FEVER test set. The class label *entailment*, *contradiction*, and *neutral* in the NLI task is mapped to *supported*, *refuted*, and *NEI*, respectively, for the fact verification task.

**LM as Fact Checker (U4).** Since there is no public available code for this model, we implement our own version following the settings described in Lee et al. (2020b). We use HuggingFace’s bert-base as the language model to predict the masked named entity, and use the NLI model described in U3 as the entailment model.

<sup>3</sup><https://huggingface.co/textattack/bert-base-uncased-MNLI>