

GERE: Generative Evidence Retrieval for Fact Verification

Jiangui Chen, Ruqing Zhang, Jiafeng Guo*, Yixing Fan, and Xueqi Cheng
CAS Key Lab of Network Data Science and Technology, Institute of Computing Technology,
Chinese Academy of Sciences, Beijing, China
University of Chinese Academy of Sciences, Beijing, China
{chenjiangui18z,zhangruqing,guojiafeng,fanyixing,cxq}@ict.ac.cn

ABSTRACT

Fact verification (FV) is a challenging task which aims to verify a claim using multiple evidential sentences from trustworthy corpora, e.g., Wikipedia. Most existing approaches follow a three-step pipeline framework, including document retrieval, sentence retrieval and claim verification. High-quality evidences provided by the first two steps are the foundation of the effective reasoning in the last step. Despite being important, high-quality evidences are rarely studied by existing works for FV, which often adopt the off-the-shelf models to retrieve relevant documents and sentences in an “index-retrieve-then-rank” fashion. This classical approach has clear drawbacks as follows: i) a large document index as well as a complicated search process is required, leading to considerable memory and computational overhead; ii) independent scoring paradigms fail to capture the interactions among documents and sentences in ranking; iii) a fixed number of sentences are selected to form the final evidence set. In this work, we propose *GERE*, the first system that retrieves evidences in a generative fashion, i.e., generating the document titles as well as evidence sentence identifiers. This enables us to mitigate the aforementioned technical issues since: i) the memory and computational cost is greatly reduced because the document index is eliminated and the heavy ranking process is replaced by a light generative process; ii) the dependency between documents and that between sentences could be captured via sequential generation process; iii) the generative formulation allows us to dynamically select a precise set of relevant evidences for each claim. The experimental results on the FEVER dataset show that *GERE* achieves significant improvements over the state-of-the-art baselines, with both time-efficiency and memory-efficiency.

CCS CONCEPTS

• **Information systems** → **Retrieval models and ranking**; • **Computing methodologies** → **Natural language processing**.

KEYWORDS

Fact Verification; Evidence Retrieval; Generative Retrieval

* Jiafeng Guo is corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '22, July 11–15, 2022, Madrid, Spain

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-8732-3/22/07...\$15.00

<https://doi.org/10.1145/3477495.3531827>

ACM Reference Format:

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 (SIGIR '22)*, July 11–15, 2022, Madrid, Spain. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3477495.3531827>

1 INTRODUCTION

With the growing online contents with false information, such as fake news, political deception and online rumors, how to automatically “fact check” the integrity of information is urgently needed for our society. Hence, many recent research efforts have been devoted to the fact verification (FV) task, which targets to automatically verify the truthfulness of a textual claim using multiple evidential sentences from trustworthy corpora, e.g., Wikipedia.

The majority of existing approaches adopts a three-step pipeline framework [3, 15, 20, 30, 38], where document retrieval and sentence retrieval component are employed to retrieve relevant documents and sentences respectively to provide evidences to the claim verification component. In essence, high-quality evidences are the foundation of claim verification to support effective determination of the veracity of claims. Existing methods for FV mainly make the claim verification the primary concern and directly adopt the off-the-shelf “index-retrieve-then-rank” approaches for evidence retrieval. Unfortunately, such approach has several shortcomings. Firstly, document retrieval needs a large document index to search over the given corpus, which requires considerable memory resources to store the whole data. The complicated search process for both the document and sentence retrieval also causes significant computation overhead. Secondly, the relevance is modeled independently on each document and sentence, missing cross-candidate interactions and local context information. Finally, a fixed number of top-ranked documents and sentences are selected for the final evidence set, limiting the flexibility in verifying different claims.

Therefore, in this paper, we propose to bypass the explicit retrieval process and introduce *GERE* (for *Generative Evidence Retrieval*), the first system that retrieves evidences in a generative way. Specifically, *GERE* exploits a transformer-based encoder-decoder architecture, pre-trained with a language modeling objective and fine-tuned to generate document titles and evidence sentence identifiers jointly. To our knowledge, this is the first attempt to use this generative class of retrieval models for FV. Finally, based on the evidences obtained by *GERE*, we train an existing claim verification model to verify the claim. *GERE* enables us to mitigate the aforementioned technical issues. Firstly, the memory footprint is greatly reduced, since the parameters of a sequence-to-sequence model scale linearly with the vocabulary size, not document count. The heavy ranking process is replaced with a light generative process, and thus we can skip the time-consuming step of searching over a

potentially massive space of options. Secondly, GERE considers the dependency information, which contributes to improving the consistency and eliminating duplication among the evidences. Finally, the generative formulation allows us to dynamically decide the number of relevant documents and sentences for different claims.

We conduct experiments on the large-scale benchmark Fact Extraction and VERification (FEVER) dataset. Experimental results demonstrate that GERE can significantly outperform the state-of-the-art baseline systems. The results also show that GERE leads to a much smaller memory footprint and time cost.

2 RELATED WORK

This section reviews previous studies on fact verification and generative retrieval models.

2.1 Fact Verification

The majority of the most successful FV framework is a three-step pipeline system, i.e., document retrieval, sentence retrieval and claim verification [15, 20, 21, 31, 38]. For document retrieval, the existing methods can be generally divided into three categories, i.e., mention-based methods [4, 7, 15, 27, 36], keyword-based methods [17, 20, 21] and feature-based methods [8, 32, 35]. However, these methods generally needs a large document index to search over the corpus, requiring a large memory footprint [2]. For sentence retrieval, three types of approaches are usually used, including traditional probabilistic ranking models [4, 25, 30], neural ranking models [7, 16, 20, 21], and pre-trained models [15, 27, 28, 36]. However, these approaches model the relevance independently, which lack of flexibility to select specific number of sentences for different claims. For claim verification, most recent studies formulate it as a graph reasoning task and pre-trained language models like BERT [10] have been widely used [14, 27, 38]. Yin and Roth [34] proposed a supervised training method named TwoWingOS to jointly conduct sentence retrieval and claim verification. Without loss of generality, this step heavily depends on the retrieved evidences, i.e., a precise set of evidences could lead to a better verification result [15, 18].

2.2 Generative Retrieval Models

Recently, generative models have attracted an increasing amount of attention in the information retrieval (IR) field. Different from the commonly-adopted “index-retrieve-then-rank” blueprint in previous works [6, 22], generative models focus on predicting relevant documents given a query based on a generation model [1, 23, 24]. For example, Metzler et al. [19] envisioned a model-based IR approach that replaces the long-lived “retrieve-then-rank” paradigm into a single consolidated model. Cao et al. [2] proposed to retrieve entities by generating their unique names in an autoregressive fashion. In this work, we make the first attempt to adapt a pre-trained sequence-to-sequence model, which has been shown to retain factual knowledge, to the evidence retrieval task in FV.

3 OUR APPROACH

In this section, we present the Generative Evidence RETrieval (GERE), a novel generative framework for evidence retrieval in FV.

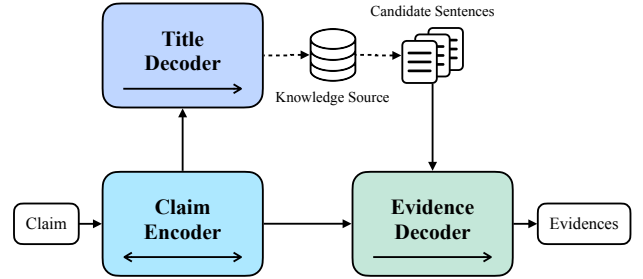


Figure 1: The overview of the GERE framework.

3.1 Model Overview

Suppose $\mathcal{K} = \{d_0, d_1, \dots\}$ denotes a large-scale text corpus where d_i denotes an individual document. d_i is composed of a sequence of sentences, namely $d_i = \{t_i, s_i^0, s_i^1, \dots, s_i^m\}$ with t_i denoting the document title of d_i and each s_i^j denoting the j -th sentence in d_i . Given a claim $c = \{c_1, c_2, \dots, c_O\}$ and a corpus \mathcal{K} , GERE aims to find a set of relevant documents $\hat{D} \subset \mathcal{K}$ and evidential sentences $\hat{E} \subset \hat{D}$ in a generative way, such that $\hat{D} = D$ and $\hat{E} = E$. This goal is different from the goal targeted by existing methods, which aim to retrieve documents $\hat{D} \subset \mathcal{K}$ and evidential sentences $\hat{E} \subset \hat{D}$ such that $D \subseteq \hat{D}$ and $E \subseteq \hat{E}$.

Basically, the GERE contains the following three dependent components: (1) Claim Encoder, a bidirectional encoder to obtain the claim representation; (2) Title Decoder, a sequential generation process to produce document titles; (3) Evidence Decoder, a sequential generation process to produce evidence sentence identifiers based on the relevant documents. The overall architecture of GERE is illustrated in Figure 1. GERE can be easily adapted to different claim verification models to provide evidences for FV.

3.2 Claim Encoder

The goal of the claim encoder is to map the input claim into a compact vector that can capture its essential topics. Specifically, the encoder represents the claim c as a series of hidden vectors, i.e.,

$$H_{enc} = \text{Encoder}(c_1, c_2, \dots, c_O), \quad (1)$$

where H_{enc} denotes the claim representation. In this work, we adopt the bidirectional Transformer-based encoder of BART [12] as the claim encoder due to its superiority in many natural language generation tasks.

3.3 Title Decoder

The goal of the title decoder is to generate a sequence of document titles $[t_1, t_2, \dots, t_{|\hat{D}|}]$ with respect to $|\hat{D}|$ relevant documents for each claim. Inspired by that “the relevance of any additional relevant document clearly depends on the relevant documents seen before” [5], the generation of a new title is decided by both the claim and previous generated titles.

Specifically, the title decoder predicts the n -th token $w_{m,n}$ in the m -th title t_m as follows,

$$p(w_{m,n} | w_{\leq m, < n}, c) = \text{Decoder}(w_{\leq m, < n}, H_{enc}),$$

where we adopt the BART [12] decoder as the title decoder. In this way, it is feasible to achieve dynamic predictions of relevant documents for different claims, i.e., a precise set of relevant documents. The title generation objective over the training corpus \mathcal{D} can be further formalized as,

$$\mathcal{L}_{title} = \arg \max_{\theta} \sum_{c \in \mathcal{D}} \sum_m \sum_n \log p(w_{m,n} | w_{\leq m, < n}, c; \theta),$$

where θ denotes the model parameters.

At the inference time, we adopt Beam Search [29] for token prediction, since it might be prohibitively expensive to compute a score for every document in Wikipedia with $\sim 5M$ documents. Naturally, the generated output might not always be a valid document title if allowing to generate any token from the vocabulary at every decoding step. To solve this problem, we employ a constrained Beam Search strategy [2] to force each generated title to be in a predefined candidate set, i.e., the titles of all the documents in \mathcal{K} . Specifically, we define our constrain in terms of a prefix tree where nodes are annotated with tokens from the predefined candidate set. For each node in the prefix tree, its children indicate all the allowed continuations from the prefix defined traversing the tree from the root to it. Note when generating a document title (e.g., Wikipedia title), the prefix tree is relatively small and it can be computed and stored into memory in advance.

3.4 Evidence Decoder

The evidence decoder is responsible for producing a sequence of evidence sentence identifiers $[e_1, e_2, \dots, e_{|\hat{E}|}]$ with respect to $|\hat{E}|$ relevant evidences given the claim and relevant documents. In this work, we define the evidence sentence identifier e_g as the unique number denoting the semantic representation of each relevant sentence. Specifically, we first identify the relevant documents from \mathcal{K} based on the titles generated by the title decoder. Then, we collect all the sentences in relevant documents as the candidate sentence set. Finally, we use another BART encoder to encode these sentences to obtain their semantic representations.

Therefore, the evidence decoder predicts the g -th sentence identifier e_g as follows,

$$p(e_g | e_{<g}, c) = \text{Decoder}(e_{<g}, H_{enc}),$$

where we also adopt the BART decoder as the evidence decoder. Note the embedding vocabulary for the evidence decoder is dynamic and it is composed of the semantic representations of the candidate sentences with respect to different claims. Similar to the title decoder, we can not only model the dependency information among evidences, but also achieve a dynamic set of evidences. The evidence generation objective over the training corpus \mathcal{D} can be further formalized as,

$$\mathcal{L}_{evidence} = \arg \max_{\beta} \sum_{c \in \mathcal{D}} \sum_g \log p(e_g | e_{<g}, c; \beta),$$

where β denotes the model parameters. At the inference time, we adopt the greedy decoding strategy due to the small search space.

3.5 Model Training

We take advantage of fine-tuning the pre-trained encoder-decoder architecture (i.e., we use BART weights [12]) and train GERE by

maximizing the title generation objective jointly with evidence generation objective, as follow,

$$\mathcal{L}_{total} = \mathcal{L}_{title} + \mathcal{L}_{evidence}.$$

In this way, the document retrieval and sentence retrieval component can be integrated into a unified framework and trained simultaneously. Given a new claim, we can easily employ the learned GERE framework to provide relevant documents as well as evidences to existing claim verification models.

4 EXPERIMENTS

In this section, we conduct experiments to verify the effectiveness of our proposed method.

4.1 Experimental Settings

We conduct experiments on the benchmark dataset FEVER [30], which consists of 185, 455 annotated claims with 5, 416, 537 Wikipedia documents. All claims are classified as SUPPORTS, REFUTES or NOT ENOUGH INFO by annotators. For each claim, the average number of relevant documents and evidences is 1.28 and 1.86, respectively. The dataset partition is kept the same with [31].

We implement our model in PyTorch based on fairseq library¹. We use the BART weights in the large version and optimize the model using Adam [11] with the warmup technique, where the learning rate increases over the first 10% of batches, and then decays linearly to zero. In the training process, we make use of the given order of documents and evidences in FEVER as the ground-truth sequential document titles and evidence identifiers. The learning rate is $3e^{-5}$, the label smoothing rate is 0.1, and the batch size is dynamic with a requirement that the max number of tokens is 4096 for each update. The size of the beams is 5. All hyper-parameters are tuned using the development set.

4.2 Evaluation Metrics

Following previous works [20, 30, 33], for the document retrieval and sentence retrieval, Precision (P), Recall (R) and F1 are used in our experiments. Note F1 is our main metric because it can directly show the performance of methods on retrieval of precise documents and evidences. For the claim verification, we adopt the official evaluation metrics, i.e., FEVER and Label Accuracy (LA). FEVER measures the accuracy with a requirement that the predicted evidences fully covers the ground-true evidences. LA measures the accuracy without considering the validity of the retrieved evidences.

4.3 Results on Document Retrieval

Following the previous works [7, 20], we choose three types of document retrieval baselines that are widely adopted in FV, including feature-based (i.e., BM25 [26], TF-IDF [20]), mention-based (i.e., UKP-Athene [7]), keyword-based (i.e., Keyword Matching [20], NSMN [20]) methods. We also select two dense retrieval methods (i.e., DPR [9] and RAG [13]) for comparison. Specifically, documents with top-5 relevance scores are selected in the baselines following the common setting [7, 20]. We do not truncate the generated sequence of documents in GERE, where 91.24% claims are provided with less than 5 documents.

¹The data and code can be found at <https://github.com/ChrisKuei/GERE>

Table 1: Comparisons of the document retrieval performance achieved by GERE and baselines; ‡ indicates statistically significant improvements over all the baselines (p-value < 0.05).

Model	Dev		
	P	R	F1
BM25	14.42	66.22	23.68
TF-IDF [20]	42.83	87.45	57.50
UKP-Athene [7]	35.33	92.51	51.13
KM [20]	44.90	83.30	58.35
NSMN [20]	52.73	88.63	66.12
DPR	55.42	89.35	68.41
RAG	62.17	91.63	74.08
GERE	84.43 ‡	78.01	81.10 ‡

Table 2: Comparisons of the sentence retrieval performance achieved by GERE and baselines; ‡ indicates statistically significant improvements over all the baselines (p-value < 0.05).

Model	Dev			Test		
	P	R	F1	P	R	F1
TF-IDF [30]	-	-	17.20	11.28	47.87	18.26
ColumbiaNLP [3]	-	78.04	-	23.02	75.89	35.33
UKP-Athene [7]	-	87.10	-	23.61	85.19	36.97
GEAR [38]	24.08	86.72	37.69	23.51	84.66	36.80
NSMN [20]	36.49	86.79	51.38	42.27	70.91	52.96
KGAT [15]	27.29	94.37	42.34	25.21	87.47	39.14
DREAM [37]	26.67	87.64	40.90	25.63	85.57	39.45
DQN [33]	54.75	79.92	64.98	52.24	77.93	62.55
GERE	58.43 ‡	79.61	67.40 ‡	54.30 ‡	77.16	63.74 ‡

As shown in Table 1, we can find that: (1) GERE significantly outperforms the state-of-the-art methods in terms of F1, demonstrating the superiority of GERE in retrieval of relevant documents. (2) The best performance of precision GERE brings comes at the cost of reduced recall compared to baselines. The reason might be that GERE generates a preciser but more compact set of documents with respect to each claim. That is, GERE leads to a lower number of retrieved documents than baselines, with the average number as 1.91. Therefore, the baselines have a higher probability to recall the ground-truth evidences than our GERE.

4.4 Results on Sentence Retrieval

For sentence retrieval, we adopt several representative models for comparison, including traditional probabilistic ranking models (i.e., TF-IDF [30], ColumbiaNLP [3]), neural ranking models (i.e., UKP-Athene [7], GEAR [38], NSMN [20]) and pre-trained models (i.e., KGAT [15], DREAM [37], DQN [33]). For the online evidence evaluation, only the first 5 sentences of predicted evidences that the

Table 3: Comparisons of different claim verification models using the evidences obtained from the original paper and that achieved by our GERE.

Model	Dev		Test	
	LA	FEVER	LA	FEVER
BERT Concat [15]	73.67	68.89	71.01	65.64
BERT Concat+GERE	74.41	70.25	71.83	66.40
BERT Pair [15]	73.30	68.90	69.75	65.18
BERT Pair+GERE	74.59	69.92	70.33	66.51
GEAR [38]	74.84	70.69	71.60	67.10
GEAR+GERE	75.96	71.88	72.52	68.34
GAT [15]	76.13	71.04	72.03	67.56
GAT+GERE	77.09	72.36	72.81	69.40
KGAT [15]	78.29	76.11	74.07	70.38
KGAT+GERE	79.44	77.38	75.24	71.17

candidate system provides are used for scoring. To meet the requirements of the online evaluation, top-5 ranked sentences are selected in the baselines and the first 5 generated sentences are kept in GERE as the evidence set. Note that 83.57% claims are provided with less than 5 sentences in GERE.

As shown in Table 2, we can see that: (1) Similar to document retrieval, GERE gives the best performance in terms of precision and F1, and performs worse than the baselines in terms of recall for sentence retrieval. This is due to that the average number of generated sentences is 2.42 in GERE, i.e., GERE provides more compact but preciser evidence set compared with baselines. (2) DQN leverages a post-processing strategy that can also find precise evidences. Although it achieves slightly better performance on recall than ours, its precision and F1 are much worse, indicating that post processing is not a optimal way to find precise evidences.

4.5 Results on Claim Verification

To verify the effectiveness of the evidences obtained by GERE, we choose several advanced claim verification models for comparison, including BERT-based models (i.e., BERT-concat [15] and BERT-pair [15]), and graph-based models (i.e., GEAR [38], GAT [15] and KGAT [15]). Specifically, these models are provided by the evidences obtained from the original paper and that achieved by GERE, respectively. Note we select the first 5 sentences generated by GERE as the evidences for these models to ensure a fair comparison.

As shown in Table 3, we can observe that: (1) All the claim verification models using the evidences obtained by our GERE significantly outperform the corresponding original versions. This result demonstrates the superiority of our method in retrieval of high-quality evidences. Modeling the dependency among documents and dependency among sentences does help improve the quality of evidences and contribute to claim verification. (2) By conducting further analysis, we find that the provided evidences in original papers generally contain conflicting pieces, some of which support the claim while the other refute. In this way, the verification process will be misled. For GERE, it can provide preciser evidences

Table 4: Comparisons on the memory footprint, the number of model parameters and inference time.

Model	Memory	Parameter	Time
NSMN	19GB	502M	28.51ms
DPR	70.9GB	220M	13.89ms
RAG	40.4GB	626M	9.46ms
GERE	2.1GB	581M	5.35ms

and the consistency between evidences are improved for different claims. The analysis again indicates that the generative formulation is helpful in improving the quality of evidences.

4.6 Memory and Inference Efficiency

In general, the document retrieval step constitutes the major part of the memory and computation cost for FV. Therefore, we evaluate the inference time of document retrieval, by GERE and three baselines, i.e., NSMN, DPR and RAG. Besides, we compare the memory footprint (disk space) needed by the overall GERE framework (including document and sentence retrieval) and the document retrieval baselines. The results are shown in Table 4. GERE has a significant reduction of memory footprint and inference time of document retrieval. The major memory computation of GERE is a prefix tree of the document titles and the number of model parameters as opposed to a large document index and a dense vector for each document in existing works. These results suggest that our model is suitable for deployment in resource-limited platforms, e.g., the online verification device.

5 CONCLUSION

In this work, we proposed *GERE* (for *Generative Evidence REtrieval*), a novel generation-based framework to address the document retrieval and sentence retrieval jointly for FV. The generative formulation leads to several advantages with respect to current solutions, including improved time-efficiency and memory-efficiency, the ability to model document/sentence dependency, and a dynamic number of evidences. GERE can be easily adapted to existing claim verification models for better claim assessment. The experimental results on the FEVER dataset demonstrates the effectiveness of GERE. In the future work, we would like to achieve the document retrieval and the sentence retrieval in a single decoder. Besides, it is valuable to design an end-to-end FV system in a fully generative way.

ACKNOWLEDGMENTS

This work was funded by the National Natural Science Foundation of China (NSFC) under Grants No. 62006218, 61902381, and 61872338, the Youth Innovation Promotion Association CAS under Grants No. 20144310, and 2021100, the Lenovo-CAS Joint Lab Youth Scientist Project, and the Foundation and Frontier Research Key Program of Chongqing Science and Technology Commission (No. cstc2017jcyjBX0059).

REFERENCES

- [1] Wasi Uddin Ahmad, Kai-Wei Chang, and Hongning Wang. 2019. Context attentive document ranking and query suggestion. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 385–394.
- [2] Nicola De Cao, Gautier Izacard, Sebastian Riedel, and Fabio Petroni. 2021. Autoregressive Entity Retrieval. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*.
- [3] Tuhin Chakrabarty, Tariq Alhindi, and Smaranda Muresan. 2018. Robust Document Retrieval and Individual Evidence Modeling for Fact Extraction and Verification. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*. 127–131.
- [4] Anton Chernyavskiy and Dmitry Ilvovsky. 2019. Extract and aggregate: A novel domain-independent approach to factual data verification. In *Proceedings of the Second Workshop on Fact Extraction and VERification (FEVER)*. 69–78.
- [5] Norbert Fuhr. 2008. A probability ranking principle for interactive information retrieval. *Information Retrieval* 11, 3 (2008), 251–265.
- [6] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W Bruce Croft. 2016. A deep relevance matching model for ad-hoc retrieval. In *Proceedings of the 25th ACM international conference on information and knowledge management*. 55–64.
- [7] Andreas Hanselowski, Hao Zhang, Zile Li, Daniil Sorokin, Benjamin Schiller, Claudia Schulz, and Iryna Gurevych. 2018. UKP-Athens: Multi-Sentence Textual Entailment for Claim Verification. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*. Association for Computational Linguistics, 103–108. <https://doi.org/10.18653/v1/w18-5516>
- [8] Christopher Hidey and Mona Diab. 2018. Team SWEEPer: Joint sentence extraction and fact checking with pointer networks. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*. 150–155.
- [9] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6769–6781.
- [10] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL-HLT*. 4171–4186.
- [11] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [12] 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 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7871–7880.
- [13] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 9459–9474. <https://proceedings.neurips.cc/paper/2020/file/6b493230205f780e1bc26945df7481e5-Paper.pdf>
- [14] Tianda Li, Xiaodan Zhu, Quan Liu, Qian Chen, Zhigang Chen, and Si Wei. 2019. Several experiments on investigating pretraining and knowledge-enhanced models for natural language inference. *arXiv preprint arXiv:1904.12104* (2019).
- [15] Zhenghao Liu, Chenyan Xiong, Maosong Sun, and Zhiyuan Liu. 2020. Fine-grained Fact Verification with Kernel Graph Attention Network. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 7342–7351.
- [16] Jackson Luken, Nanjiang Jiang, and Marie-Catherine de Marneffe. 2018. QED: A fact verification system for the FEVER shared task. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*. 156–160.
- [17] Jing Ma, Wei Gao, Shafiq Joty, and Kam-Fai Wong. 2019. Sentence-level evidence embedding for claim verification with hierarchical attention networks. Association for Computational Linguistics.
- [18] Sean MacAvaney, Andrew Yates, Arman Cohan, and Nazli Goharian. 2019. CEDR: Contextualized embeddings for document ranking. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1101–1104.
- [19] Donald Metzler, Yi Tay, Dara Bahri, and Marc Najork. 2021. Rethinking search: making domain experts out of dilettantes. In *ACM SIGIR Forum*, Vol. 55. ACM New York, NY, USA, 1–27.
- [20] 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*, Vol. 33. 6859–6866.
- [21] Yixin Nie, Songhe Wang, and Mohit Bansal. 2019. Revealing the Importance of Semantic Retrieval for Machine Reading at Scale. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 2553–2566.

- [22] Rodrigo Nogueira and Kyunghyun Cho. 2019. Passage Re-ranking with BERT. *arXiv preprint arXiv:1901.04085* (2019).
- [23] Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. 2020. Document Ranking with a Pretrained Sequence-to-Sequence Model. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*. 708–718.
- [24] Rodrigo Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document expansion by query prediction. *arXiv preprint arXiv:1904.08375* (2019).
- [25] Beatrice Portelli, Jason Zhao, Tal Schuster, Giuseppe Serra, and Enrico Santus. 2020. Distilling the evidence to augment fact verification models. In *Proceedings of the Third Workshop on Fact Extraction and VERification (FEVER)*. 47–51.
- [26] Stephen Robertson and Hugo Zaragoza. 2009. *The probabilistic relevance framework: BM25 and beyond*. Now Publishers Inc.
- [27] Amir Soleimani, Christof Monz, and Marcel Worring. 2020. BERT for evidence retrieval and claim verification. *Advances in Information Retrieval* 12036 (2020), 359–366.
- [28] Shyam Subramanian and Kyumin Lee. 2020. Hierarchical Evidence Set Modeling for Automated Fact Extraction and Verification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 7798–7809.
- [29] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. *Advances in neural information processing systems* 27 (2014).
- [30] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a Large-scale Dataset for Fact Extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. 809–819.
- [31] James Thorne, Andreas Vlachos, Oana Cocarascu, Christos Christodoulopoulos, and Arpit Mittal. 2018. The Fact Extraction and VERification (FEVER) Shared Task. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*. 1–9.
- [32] Santosh Tokala, G Vishal, Avirup Saha, and Niloy Ganguly. 2019. AttentiveChecker: A bi-directional attention flow mechanism for fact verification. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 2218–2222.
- [33] Hai Wan, Haicheng Chen, Jianfeng Du, Weilin Luo, and Rongzhen Ye. 2021. A DQN-based Approach to Finding Precise Evidences for Fact Verification. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 1030–1039.
- [34] Wenpeng Yin and Dan Roth. 2018. TwoWingOS: A Two-Wing Optimization Strategy for Evidential Claim Verification. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. 105–114.
- [35] Takuma Yoneda, Jeff Mitchell, Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Ucl machine reading group: Four factor framework for fact finding (hexaf). In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*. 97–102.
- [36] Chen Zhao, Chenyan Xiong, Corby Rosset, Xia Song, Paul Bennett, and Saurabh Tiwary. 2020. Transformer-XH: Multi-evidence Reasoning with Extra Hop Attention. In *The Eighth International Conference on Learning Representations (ICLR 2020)*. <https://www.microsoft.com/en-us/research/publication/transformer-xh-multi-evidence-reasoning-with-extra-hop-attention/>
- [37] 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 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 6170–6180.
- [38] 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 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 892–901.