

# Improving Evidence Retrieval for Automated Explainable Fact-Checking

Chris Samarinas<sup>1</sup>, Wynne Hsu<sup>2,3</sup>, and Mong Li Lee<sup>2,3</sup>

<sup>1</sup>Institute of Data Science, National University of Singapore

<sup>2</sup>NUS Centre for Trusted Internet and Community

<sup>3</sup>School of Computing, National University of Singapore

## Abstract

Automated fact-checking on a large-scale is a challenging task that has not been studied systematically until recently. Large noisy document collections like the web or news articles make the task more difficult. We describe a three-stage automated fact-checking system, named Quin+, using evidence retrieval and selection methods. We demonstrate that using dense passage representations leads to much higher evidence recall in a noisy setting. We also propose two sentence selection approaches, an embedding-based selection using a dense retrieval model, and a sequence labeling approach for context-aware selection. We deployed a domain-specific demo of our fact-checking system to verify claims related to the COVID-19 pandemic, and a generic demo for verifying open-domain claims using results from web search engines.

## 1 Introduction

With the emergence of social media and many individual news sources online, the spread of misinformation has become a major problem with potentially harmful social consequences. Fake news can manipulate the public opinion, create conflicts, elicit unreasonable fear and suspicion. The vast amount of unverified online content stimulated the appearance of the first external post-hoc fact-checking organizations since the early 2000s. Numerous organizations from different countries have dedicated resources to verify claims online, such as PolitiFact, FactCheck.org, Snopes etc. However, manual fact-checking is time consuming and intractable on a large scale. The ability to automatically perform fact-checking is critical to minimize negative social impact.

Automated fact checking is a complex task involving evidence extraction followed by evidence reasoning and entailment. For the retrieval of rel-

evant evidence from a corpus of documents, existing systems typically utilize traditional sparse retrieval which may have poor recall, especially when the relevant passages have few overlapping words with the facts to be verified. Dense retrieval models have been proven effective in question answering as these models can better capture the latent semantic content of text. The work in (Samarinas et al., 2020) is the first to use dense retrieval in fact checking. The authors constructed a new dataset called Factual-NLI comprised of claim-evidence pairs from the FEVER dataset (Thorne et al., 2018) as well as synthetic examples generated from benchmark Question Answering datasets (Kwiatkowski et al., 2019; Nguyen et al., 2016). They demonstrated that using Factual-NLI to train a dense retriever can improve evidence retrieval significantly.

While the FEVER dataset has enabled the systematic evaluation of automated fact-checking systems, it does not reflect well the noisy nature of real-world data. Motivated by this, we introduce the Factual-NLI+ dataset, an extension of the FEVER dataset with synthetic examples from question answering datasets and noise passages from web search results. We examine how dense representations can improve the first-stage retrieval recall of passages for fact-checking in a noisy setting, and make the retrieval of relevant evidences more tractable on a large scale.

However, the selection of relevant evidence sentences for more accurate fact-checking and explainability remains a challenge. Figure 1 shows an example of a claim and the retrieved passage which has three sentences, of which only the last sentence provides the critical evidence to refute the claim. We propose two ways to select the relevant sentences, an embedding-based selection using a dense retrieval model, and a sequence labeling approach for context-aware selection. We

show that the former generalizes better with a high recall, while the latter has higher precision making it suitable for the identification of relevant evidence sentences.

We deploy a fact-checking demo system for the verification of claims about the COVID-19 pandemic using news articles and research papers. Our system is also able to verify open-domain claims using web search results.

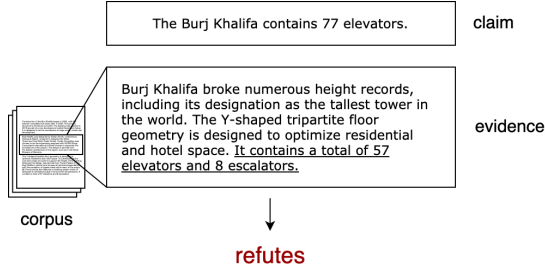


Figure 1: Sample claim and the retrieved passage evidence where only the last sentence is relevant.

## 2 Related Work

Automated claim verification using a large corpus has not been studied systematically until the availability of the Fact Extraction and VERification dataset (FEVER) (Thorne et al., 2018). This dataset contains claims that are supported or refuted by specific evidence from Wikipedia articles. Prior to the work in (Samarinas et al., 2020), fact-checking solutions relied on sparse passage retrieval, followed by a claim verification (entailment classification) model (Nie et al., 2019). Other approaches used the mentions of entities in a claim and/or very basic entity linking to retrieve documents and a machine learning model such as logistic regression or an enhanced sequential inference model to decide whether an article most likely contains the evidence (Yoneda et al.; Chen et al., 2017; Hanselowski et al., 2018).

However, retrieval based on sparse representations and exact keyword matching can be rather restrictive for various queries. This restriction can be mitigated by dense representations using BERT-based language models (Devlin et al., 2019). The works in (Lee et al., 2019; Karpukhin et al., 2020; Xiong et al., 2020; Chang et al., 2020) have successfully used such models and its variants for passage retrieval in open-domain question answering. The results can be further improved using passage re-ranking with cross-attention BERT-based models (Nogueira et al., 2019). The work in

(Samarinas et al., 2020) is the first to propose a model called QR-BERT to retrieve passages for fact-checking.

Apart from passage retrieval, sentence (evidence) selection is also a critical task in fact-checking. These evidence sentences provide an explanation why a claim has been assessed to be credible or not. Recent works have proposed a BERT-based model for extracting relevant evidence sentences from multi-sentence passages (Atanasova et al., 2020). The authors observe that joint training on veracity prediction and explanation generation performs better than training separate models. The work in (Stammbach and Ash) investigates how the few-shot learning capabilities of the GPT-3 model (Brown et al., 2020) can be used for generating fact-checking explanations.

## 3 The Quin+ System

The automated claim verification task can be defined as follows: given a textual claim  $c_m$  and a corpus  $D = \{d_1, d_2, \dots, d_n\}$ , where every passage  $d_i = \{s_1^i, s_2^i, \dots, s_k^i\}$  is comprised of sentences  $s_j^i$ , a system must return a set of evidence sentences  $E_m = \{s_1, s_2, \dots, s_k\} \subset \bigcup d_i$  and a label  $\hat{y}_m \in \{\text{supports}, \text{refutes}, \text{not enough info}\}$ . We develop an automated fact-checking system, called Quin+, to perform claim verification in three stages: *passage retrieval* from a corpus, *sentence selection* and *entailment classification* as shown in Figure 2. The label of a claim can be decided with two different aggregation approaches depending on the nature of the corpus; (a) by doing entailment classification on each selected passage and assigning the most frequent class or (b) by doing entailment classification on the concatenation of selected evidence sentences. In our evaluation, where each claim has only one source of relevant evidence, we used the second approach, while in our demos, where in the general case we have multiple sources of relevant evidence, we used the first.

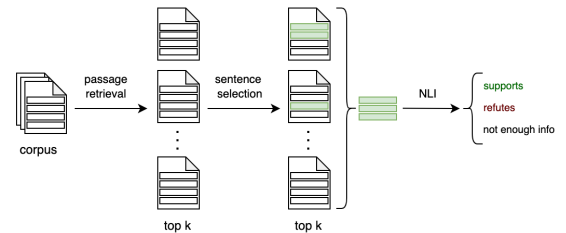


Figure 2: Three stages of claim verification in Quin+.

### 3.1 Passage Retrieval.

Quin+’s passage retrieval model is based on QR-BERT, a *dense* retrieval model (Samarinas et al., 2020). QR-BERT creates dense vectors for passages by calculating their average token embedding. The relevance of a passage  $d$  to a claim  $c$  is then given by their dot product:

$$r(c, d) = \phi(c)^T \phi(d) \quad (1)$$

Dot product search can run efficiently using an approximate nearest neighbors index implemented using the FAISS library (Johnson et al., 2017). QR-BERT maximizes the sampled softmax loss:

$$L_\theta = \sum_{(c,d) \in D_B^+} r_\theta(c, d) - \log \left( \sum_{d_i \in D_B} e^{r_\theta(c, d_i)} \right) \quad (2)$$

where  $D_B$  is the set of passages in a training batch  $B$ , and  $D_B^+$  is the set of positive claim-passage pairs in  $B$ .

The work in (Samarinas et al., 2020) introduce **Factual-NLI** dataset that extends the FEVER dataset (Thorne et al., 2018) with more diverse synthetic examples derived from question answering datasets. There are 359,190 new entailed claims with evidence and additional contradicted claims from a rule-based approach.

To ensure robustness, we compiled a new large-scale *noisy* version of Factual-NLI called Factual-NLI+. This dataset includes all the 5 million Wikipedia passages in the FEVER dataset. We add ’noise’ passages as follows: For every claim  $c$  in the FEVER dataset, we retrieve the top 30 web results from the Bing search engine and keep passages with the highest BM25 score that are classified as neutral by the entailment model. For claims generated from MSMARCO queries (Nguyen et al., 2016), we included the irrelevant passages that are found in the MSMARCO dataset for those queries. This results in 418,650 additional passages. The new dataset reflects better the nature of a large-scale corpus that would be used by real-world fact-checking system. We train a dense retrieval model using this extended dataset.

The Quin+ system utilizes a hybrid model that combines the results from the dense retrieval model described above and BM25 sparse retrieval to obtain the final list of retrieved passages. For efficient sparse retrieval, we used the Rust-based Tantivy full text search engine <sup>1</sup>.

<sup>1</sup><https://github.com/tantivy-search/tantivy>

### 3.2 Sentence Selection.

We propose and experiment with two sentence selection methods: an embedding-based selection and context-aware sentence selection method.

The embedding-based selection method relies on the dense representations learned by the dense passage retrieval model QR-BERT. For a given claim  $c$ , we select the sentences  $s_j$  from a given passage  $d = \{s_1, s_2, \dots, s_k\}$  such that the relevance score  $r(c, d)$  from the QR-BERT model is greater than some threshold  $\lambda$  which is set experimentally.

The context-aware sentence selection method uses a sequence labelling approach. We adopt the BIO tagging format so that all the irrelevant tokens are classified as  $O$ , the first token of an evidence sentence classified as  $B$  evidence and the rest tokens of an evidence sentence as  $I$  evidence (see Figure 3). We train a model based on RoBERTa-large (Liu et al., 2019), minimizing the cross-entropy loss:

$$L_\theta = - \sum_{i=1}^N \sum_{j=1}^{l_i} \log(p_\theta(y_j^i)) \quad (3)$$

where  $N$  is the number of examples in the training batch,  $l_i$  the number of non-padding tokens of the  $i^{th}$  example and  $p_\theta(y_j^i)$  is the estimated softmax probability of the correct label for the  $j^{th}$  token of the  $i^{th}$  example.

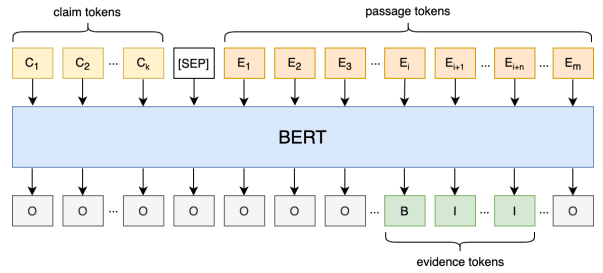


Figure 3: Sequence labeling model for evidence selection from a passage for a given claim.

We trained the sentence selection model in Quin+ on Factual-NLI with batch size 64, Adam optimizer and initial learning rate  $5 \times 10^{-5}$  until convergence.

### 3.3 Entailment Classification.

Natural Language Inference (NLI), also known as textual entailment classification, is the task of

detecting whether a hypothesis statement is entailed by a premise passage. It is essentially a text classification problem, where the input is a pair of premise-hypothesis  $(P, H)$  and the output a label  $y \in \{\text{entailment, contradiction, neutral}\}$ . An NLI model is often a core component of many automated fact-checking systems. Datasets like the Stanford Natural Language Inference corpus (SNLI) (Bowman et al., 2015), Multi-Genre Natural Language Inference corpus (Multi-NLI) (Williams et al., 2018) and Adversarial-NLI (Nie et al., 2020) have facilitated the development of models for this task.

In our system, we used an NLI model based on RoBERTa-large with a linear transformation of the [CLS] token embedding:

$$o = \text{softmax}(W \cdot \text{BERT}_{[\text{CLS}]}([P; H]) + b) \quad (4)$$

where  $P; H$  is the concatenation of the premise with the hypothesis,  $W_{3 \times 1024}$  a linear transformation matrix and  $b_{3 \times 1}$  the added bias. We trained the entailment model by minimizing the cross-entropy loss on the concatenation of the three popular NLI datasets (SNLI, Multi-NLI and Adversarial-NLI) with batch size 64, Adam optimizer and initial learning rate  $5 \times 10^{-5}$  until convergence.

#### 4 Performance of Quin+

We evaluate the three individual components of Quin+ (retrieval, sentence selection and entailment classification) and finally perform an end-to-end evaluation using various configurations.

Table 1 gives the recall@k and Mean Reciprocal Rank (MRR@100) of the passage retrieval models on the FEVER, the noisy extension of FEVER (subset of Factual-NLI+) and Factual-NLI+. We see that the hybrid passage retrieval model in Quin+ gives the best performance compared to BM25 and the dense retrieval model. When adding noise to the FEVER dataset, the gap between Quin’s hybrid passage retrieval model and sparse retrieval widens. This demonstrates the limitations of using sparse retrieval, and why it is crucial to have a dense retrieval model to surface relevant passages from a noisy corpus.

Table 2 shows the token-level precision, recall and F1 score of the proposed sentence selection methods on the Factual-NLI dataset and a domain-specific (medical) claim verification dataset, SciFact (Wadden et al., 2020). We also compare the

(a) FEVER Dataset					
Model	R@5	R@10	R@20	R@100	MRR
BM25	50.53	58.92	67.93	82.93	0.381
Dense	65.47	69.61	72.51	75.71	0.535
Hybrid	<b>71.71</b>	<b>78.60</b>	<b>83.65</b>	<b>91.09</b>	<b>0.556</b>

(b) FEVER + noise					
Model	R@5	R@10	R@20	R@100	MRR
BM25	35.17	44.18	53.89	73.95	0.2649
Dense	54.10	62.13	68.09	75.24	0.4053
Hybrid	<b>54.89</b>	<b>64.61</b>	<b>73.33</b>	<b>86.11</b>	<b>0.4074</b>

(c) Factual-NLI+ Dataset					
Model	R@5	R@10	R@20	R@100	MRR
BM25	45.02	53.20	61.56	77.96	0.347
Dense	59.66	67.09	72.23	78.52	0.461
Hybrid	<b>61.29</b>	<b>70.03</b>	<b>77.51</b>	<b>87.90</b>	<b>0.465</b>

Table 1: Performance of passage retrieval models.

performance to a baseline sentence level NLI approach.

The context-aware sequence labeling (SL) model gives the highest precision, recall and F1 score when tested on the Factual-NLI dataset. Further, the precision is significantly higher than the other methods. On the other hand, for the SciFact dataset, we see that sequence labeling method remains the top performer in terms of precision and F1 score after fine-tuning, but its recall is lower than the embedding-based method. This shows that sequence labeling model is able to mitigate the high false positive rate observed with the embedding-based selection method by taking into account the surrounding context.

(a) Factual-NLI Dataset			
Model	P	R	F1
Sequence labeling	<b>94.78</b>	<b>92.11</b>	<b>93.43</b>
Embedding-based	66.12	90.29	76.34
Sentence Entailment	67.74	91.87	77.98

(b) SciFact Dataset			
Model	P	R	F1
Sequence labeling (F-NLI)	66.25	55.12	60.17
Sequence labeling (fine-tuned)	<b>69.38</b>	68.45	<b>68.91</b>
Embedding-based	43.30	<b>92.36</b>	58.96
Sentence Entailment	62.21	71.54	66.55

Table 2: Performance of Sentence selection.

Table 3 shows the performance of the entailment classification model (claim verification performance) on the ground-truth (oracle) passages of the Factual-NLI+ dataset. The neutral claims do not have any ground-truth passages. In order

Evidence	Supports			Refutes			Neutral			Macro Avg		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Oracle sentences	82.15	60.05	69.38	77.54	89.32	83.02	84.62	91.28	87.83	81.44	80.22	80.07
Oracle passages	63.40	53.93	58.28	33.95	40.65	37.00	70.33	71.37	70.85	55.90	55.32	55.38
Selected sentences	74.40	56.68	64.34	75.27	81.96	78.47	81.67	88.32	84.87	77.11	75.66	75.89

Table 3: Performance of NLI / claim verification with oracle passage retrieval on Factual-NLI+

to examine the performance for this class as well, we use the top passage returned by QR-BERT as ground-truth for every neutral claim. An interesting observation is that when we pass the whole ground-truth passages to the entailment classification model, the model tends to predict the neutral class most of the time. Selecting the correct evidence sentences using the proposed sequence labeling model alleviates this issue significantly (see figure 4). The reason behind this is that NLI models identify the correct class much easier when the premise is short, because they are trained mostly on short premises from the currently available datasets.

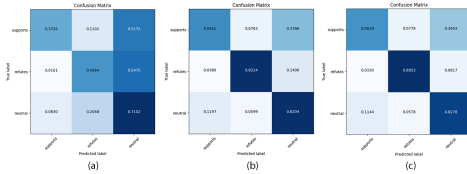


Figure 4: The confusion matrices for the entailment classification of the claim-passage pairs in Factual-NLI development set using: the whole passage (a), only evidence sentences from the sequence labeling model (b), and only golden evidence sentences (c).

Finally, we do end-to-end evaluation of our system on Factual-NLI+ using various configurations of passage retrieval and evidence selection approaches (table 4). We report the macro-average F1 score for the three classes (supporting, refuting, neutral). From our experiments, dense or hybrid retrieval with evidence selection using the proposed sequence labeling model gives the best results. Even though hybrid retrieval seems to lead to slightly worse performance, it requires much less passages (6 instead of 50) and makes the system more efficient.

## 5 System Demonstration

We have deployed the Quin+ system<sup>2</sup> for verifying claims about the COVID-19 pandemic from

<sup>2</sup><https://quin.algoprog.com>

Retrieval	Evidence	F1
Oracle	Oracle	80.07
BM25, k=5	Embedding-based	52.76
BM25, k=20	Embedding-based	47.65
BM25, k=5	Sequence labeling	49.65
Dense, k=5	Embedding-based	49.03
Dense, k=5	Sequence labeling	52.83
Dense, k=50	Sequence labeling	<u>58.22</u>
Hybrid, k=6	Embedding-based	50.29
Hybrid, k=6	Sequence labeling	<u>57.24</u>
Hybrid, k=50	Sequence labeling	52.60

Table 4: End-to-end claim verification evaluation on Factual-NLI+ for different setups of passage retrieval and evidence selection

over 400K news articles and 100K research papers from the CORD-19 dataset (Wang et al., 2020). In addition, we created a demo for verifying any open-domain claim using the top 20 results from the Bing search engine. For a given claim, Quin+ returns relevant text snippets with highlighted evidence. The snippets are grouped into two sets, one of supporting and another of refuting snippets. Based on the number of supporting and refuting snippets, it displays a veracity rating. If there are significantly more supporting, it returns "Probably True", in the opposite case "Probably False", and when there's very small number of relevant evidence it returns "Not Enough Evidence".

The system utilizes hybrid retrieval and the embedding-based evidence selection model. Even though the sequence labeling model gave better results on Factual-NLI+, in our demos the embedding-based selection model seemed to work much better, possibly because it maintains much higher recall across different domains (as seen in table 2). When we have multiple sources of relevant evidence and not just one or two as in the FEVER and Factual-NLI+ datasets, the recall is more important as it leads to more relevant evidence selected and more accurate veracity rating based on multiple sources. Figure 5 shows the system in action.



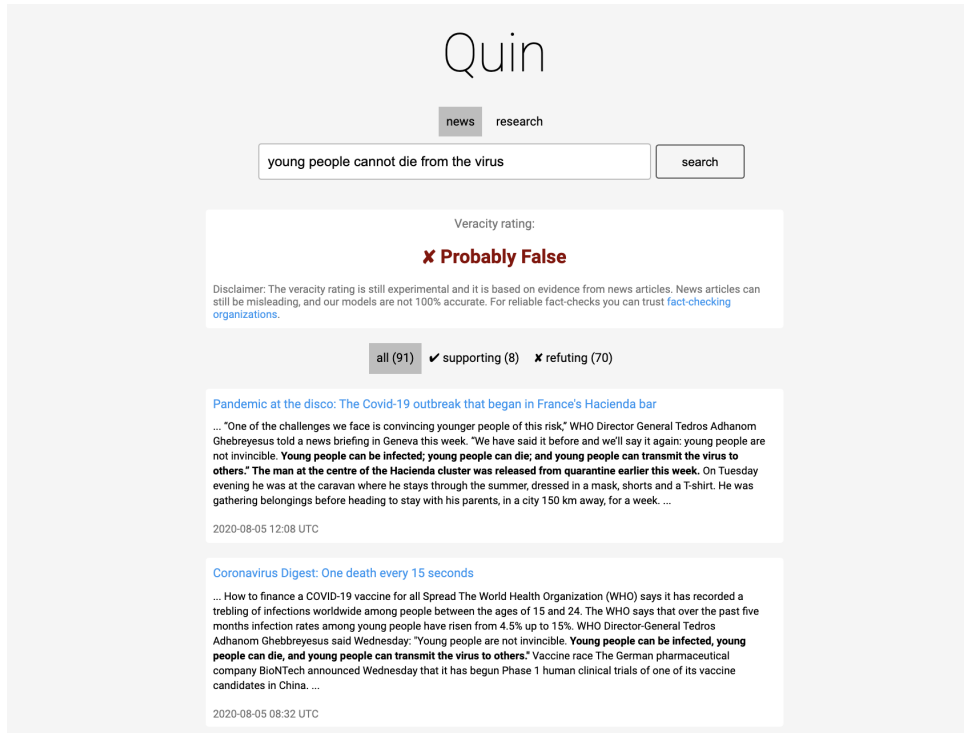


Figure 5: The Quin system returning relevant evidence and a veracity rating for a claim about COVID-19.

## 6 Conclusion & Future Work

In this paper we examined how the first stages of an automated fact-checking system can be improved. More specifically, we investigated how a dense passage retrieval model can lead to higher recall when retrieving relevant evidence, and proposed a sequence labeling model for selecting relevant evidence sentences for improving the final verification accuracy. For our experiments we also introduced an extension of the Factual-NLI dataset, that constitutes a more challenging benchmark for passage retrieval in fact-checking.

This work focused on improving the first two steps of a fact-checking system, without attempting to improve the last step, the entailment classification. Even though recent pre-trained language models seem to perform very well on the two popular NLI datasets (SNLI and Multi-NLI), they are not equally effective in a real-world setting. Recent work has identified bias in these two datasets (Poliak et al., 2018), which has a negative effect in the generalization ability of the trained models. Various approaches have been proposed to unlearn or avoid this bias (He et al., 2019; Belinkov et al., 2019). Another weakness of these datasets, is that they are comprised of short, usually single-sentence, premises. As a consequence, models trained on these datasets usually do not

perform well on noisy real-world data with multiple sentences. These issues led to the development of additional more challenging datasets, such as Adversarial NLI (Nie et al., 2020), which we also used for training our NLI model. Including examples from Adversarial-NLI improved the entailment classification F1 score from 74.14% to 80.07% (table 3), but still leaves a lot of room for improvement. Generalizable Natural Language Inference remains an open problem, and requires understanding of the reasons behind the misclassifications, the introduction of more diverse training examples without annotation artifacts and better language models.

Last but not least, we believe that for the future development of large-scale fact-checking systems, a new benchmark needs to be introduced. The currently available datasets (including Factual-NLI+) are not suitable for evaluating the performance of systems that have a large number of evidence sources (such as news publishers or research papers). Even though we found a best-performing system configuration on Factual-NLI+, different configurations seemed to work better in our demos. This shows the need for a new benchmark for fact-checking that takes into account the existence of many individual evidence sources.

## References

- Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. [Generating fact checking explanations](#).
- Yonatan Belinkov, Adam Poliak, Stuart M. Shieber, Benjamin Van Durme, and Alexander M. Rush. 2019. [Don't take the premise for granted: Mitigating artifacts in natural language inference](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 877–891. Association for Computational Linguistics.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
- Wei-Cheng Chang, Felix X. Yu, Yin-Wen Chang, Yiming Yang, and Sanjiv Kumar. 2020. [Pre-training tasks for embedding-based large-scale retrieval](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2017. Enhanced lstm for natural language inference. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1657–1668.
- J. Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.
- Andreas Hanselowski, Hao Zhang, Zile Li, Daniil Sorokin, Benjamin Schiller, Claudia Schulz, and Iryna Gurevych. 2018. Ukp-athene: Multi-sentence textual entailment for claim verification. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 103–108.
- He He, Sheng Zha, and Haohan Wang. 2019. Unlearn dataset bias in natural language inference by fitting the residual. In *DeepLo@EMNLP-IJCNLP*.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2017. Billion-scale similarity search with gpus. *arXiv preprint arXiv:1702.08734*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Ledell Wu, Sergey Edunov, Danqi Chen, and Wenteau Yih. 2020. Dense passage retrieval for open-domain question answering. *arXiv preprint arXiv:2004.04906*.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. [Latent retrieval for weakly supervised open domain question answering](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 6086–6096. Association for Computational Linguistics.
- 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*.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. [MS MARCO: A human generated machine reading comprehension dataset](#). *CoRR*, abs/1611.09268.
- 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.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. 2020. [Adversarial NLI: A new benchmark for natural language understanding](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4885–4901, Online. Association for Computational Linguistics.
- Rodrigo Nogueira, W. Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-stage document ranking with bert. *ArXiv*, abs/1910.14424.
- Adam Poliak, Jason Naradowsky, Aparajita Haldar, Rachel Rudinger, and Benjamin Van Durme. 2018. [Hypothesis only baselines in natural language inference](#). In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pages 180–191, New Orleans, Louisiana. Association for Computational Linguistics.
- Chris Samarin, Wynne Hsu, and Mong Li Lee. 2020. Latent retrieval for large-scale fact-checking and question answering with nli training.

Dominik Stammach and Elliott Ash. e-fever: Explanations and summaries for automated fact checking.

James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and verification. In *NAACL-HLT*.

David Wadden, Kyle Lo, Lucy Lu Wang, Shanchuan Lin, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. *ArXiv*, abs/2004.14974.

Lucy Lu Wang, Kyle Lo, Yoganand Chandrasekhar, Russell Reas, Jiangjiang Yang, Darrin Eide, K. Funk, Rodney Michael Kinney, Ziyang Liu, W. Merrill, P. Mooney, D. Murdick, Devvret Rishi, Jerry Sheehan, Zhihong Shen, B. Stilson, A. Wade, K. Wang, Christopher Wilhelm, Boya Xie, D. Raymond, Daniel S. Weld, Oren Etzioni, and Sebastian Kohlmeier. 2020. Cord-19: The covid-19 open research dataset. *ArXiv*.

Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1112–1122. Association for Computational Linguistics.

Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2020. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *arXiv preprint arXiv:2007.00808*.

Takuma Yoneda, Jeff Mitchell, Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. Ucl machine reading group: Four factor framework for fact finding (hexaf). In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 97–102.