

Where Are the Facts? Searching for Fact-checked Information to Alleviate the Spread of Fake News

Nguyen Vo and Kyumin Lee

Computer Science Department, Worcester Polytechnic Institute
Worcester, Massachusetts, 01609, USA
{nkvo, kmlee}@wpi.edu

Abstract

Although many fact-checking systems have been developed in academia and industry, fake news is still proliferating on social media. These systems mostly focus on fact-checking but usually neglect *online users* who are the main drivers of the spread of misinformation. How can we use fact-checked information to improve users' consciousness of fake news to which they are exposed? How can we stop users from spreading fake news? To tackle these questions, we propose a novel framework to search for fact-checking articles, which address the content of an original tweet (that may contain misinformation) posted by online users. The search can directly warn fake news posters and online users (e.g. the posters' followers) about misinformation, discourage them from spreading fake news, and scale up verified content on social media. Our framework uses both text and images to search for fact-checking articles, and achieves promising results on real-world datasets. Our code and datasets are released at <https://github.com/nguyenvo09/EMNLP2020>.

1 Introduction

The rampant spread of biased news, partisan stories, false claims and misleading information has raised heightened societal concerns in recent years. Many reports pointed out that fabricated stories possibly caused citizens' misperception about political candidates (Allcott and Gentzkow, 2017), manipulated stock prices (Kogan et al., 2019) and threatened public health (Ashoka, 2020; Alluri, 2019).

The proliferation of misinformation has provoked the rise of fact-checking systems worldwide. Since 2014, the number of fact-checking outlets has totally increased 400% in 60 countries (Stencel, 2019). However, fabricated stories and hoaxes are still pervading our cyberspace. Fig. 1 shows an example of a fake quote related to Barack Obama.



Figure 1: An original tweet and a mock-up of how a correctly retrieved fact-checking article is presented.

The quote had been debunked by Snopes (Emery, 2016) on September 08, 2016 but two months later, it appeared again inside an original tweet posted by a Twitter user (called an original poster) and was retweeted over 28 thousand times. Perhaps, the original poster and people who shared the original tweet did not know if it was fact-checked or they might share it simply because it was suitable for their personal preferences or ideologies (Lewandowsky et al., 2012). In other words, existing fact-checking systems mainly focus on detection but neglect online users who play the critical role in spreading fake news. After detecting fake news, what are the next steps to discourage people from sharing it? Recent studies (Vo and Lee, 2018, 2019) tried to curb the above weakness. However, these approaches are not proactive since they rely on fact-checkers who may be unreliable.

Recent works showed that when seeing fact-checked information, users' likelihood to delete fake news's shares went up 400% (Friggeri et al., 2014) and 95% of the time users did not further consume or go through fake news (CNN, 2020). Observing the downside of existing methods and impacts of broadcasting verified news, our goal is to search for *fact-checking articles* (FC-articles) which address the content of original tweets (i.e. confirming, supporting, debunking or refuting). We

show a mock-up of how a relevant FC-article is linked/displayed given an original tweet in Fig. 1. By searching for FC-articles and incorporating fact-checked information into social media posts, we can warn users (e.g. followers of original posters) about fake news to which they are exposed. The search also proactively scales up volume of verified content on social media. However, achieving the goal is challenging since we need to solve two problems: **(P1)** what information in original tweets should we use to find correct FC-articles? and **(P2)** how can we design a framework to retrieve and rank FC-articles?

With the first problem **(P1)**, we can use original tweets’ text to find FC-articles. However, this approach is suboptimal since fake news can appear in many forms (e.g. text, images, videos) (Friggeri et al., 2014; O’Brien, 2018) as shown in Fig. 1. Thus, we propose to use both text and images of original tweets to search for FC-articles. Regarding the second problem **(P2)**, we propose a framework consisting of two key steps: (1) using a basic retrieval (i.e. BM25) to find initial lists of candidate FC-articles and then (2) re-ranking the initial lists by using advanced models for ranking refinement. In the first step, since original tweets’ text may be insufficient to find correct articles as shown in Fig. 1 where there is no meaningful information in the text but in the image, we propose to expand original tweets’ text by using text inside original tweets’ images. For the second step, we propose an attention mechanism to focus on key textual matching signals and jointly integrate them with visual information to boost ranking quality. By tackling these issues, our contributions are as follows:

- To the best of our knowledge, our study is the first one that searches for fact-checking articles to increase users’ awareness of fact-checked information when they are exposed to fake news.
- We propose a novel neural ranking model which jointly utilizes textual and visual matching signals. The model is also integrated with a novel attention mechanism.
- Experiments on two datasets demonstrate effectiveness and generality of our model over state-of-the-art retrieval techniques.

2 Related Work

Fake News and Fact-checking. Fake news detection methods mainly use linguistics and textual content (Zellers et al., 2019; Zhao et al., 2015; Wang,

2017; Shu et al., 2019), temporal spreading patterns (Liu and Wu, 2018; Ma et al., 2018), network structures (Wu and Liu, 2018; Liu et al., 2020) and users’ feedbacks (Vo and Lee, 2019, 2020; Shu et al., 2019). Studies about multimodal fake news detection (Gupta et al., 2013; Wang et al., 2018b) are different from ours since their inputs are text and images of tweets while our inputs are pairs of a multimodal tweet and a FC-article.

Our work is closely related to evidence-aware fact-checking. Thorne et al. (2018); Nie et al. (2019) built a pipeline to find documents and sentences to fact-check mutated claims generated from Wikipedia pages, Wang et al. (2018a) aimed to find webpages related to given FC-articles and predict their stances on claims in the FC-articles. Popat et al. (2018) only focused on fact-checking and (Shaar et al., 2020) detected previously fact-checked claims. Our paper deviates from these work since we aim to find FC-articles given multimodal fake news in social media posts. As our goal is to increase users’ awareness of verified news, studies about fact-checkers (Vo and Lee, 2018, 2019; You et al., 2019) are close to ours.

Neural Ranking Models for Text Search. Neural ranking models for text search mainly fall into two groups: semantic matching and relevance matching models. The former one seeks to learn representations of a query and a document, and measure their similarity (Huang et al., 2013; Shen et al., 2014; Severyn and Moschitti, 2015; Nie et al., 2019; Zhu et al., 2019), while the later one (Chen et al., 2018; Hui et al., 2017; Pang et al., 2016; Guo et al., 2016; Xiong et al., 2017; Hui et al., 2018; Dai et al., 2018) aims to capture relevant matching signals between a query and a document based on word interactions. There are methods unifying two categories such as Mitra et al. (2017); Rao et al. (2019a). Our model can be viewed as a relevance matching method in which a novel attention mechanism is designed to focus on crucial word interactions.

Neural Models for Multimodal Retrieval. Multimodal data (e.g. text and images) are used in cross-modal retrieval (Cao et al., 2016; Wang et al., 2017; Balaneshin-kordan and Kotov, 2018; Chen et al., 2016), visual Q&A task (Kim et al., 2018), product search (Laenen et al., 2018; Guo et al., 2018) and so on. Our work is the first using multimodal data in social media posts to search for verified information.

3 Our framework

Given an original tweet q and a FC-article d , where every original tweet q contains text and images and the article d contains text and/or images, we aim to derive function $f(q, d)$ which determines their relevancy¹. We use $f(q, d)$ to rank all FC-articles.

Following (Thorne et al., 2018), we adopt the re-ranking methodology as follows: (1) quickly retrieving candidate FC-articles/documents² for each original tweet/query³ by a basic retrieval and (2) re-ranking the candidates by our MAN (Multimodal Attention Network) as shown in Fig. 2. We describe our input representations, basic retrieval and MAN in following subsections.

3.1 Input Representations

We denote text and images of an original tweet q as (q_{text}, q_{images}) where q_{text} is a sequence of N words $\{w_i^q\}_{i=1}^N$ and q_{images} is a list of X images $\{v_i^q\}_{i=1}^X$. Similarly, text and images of a fact-checking article d are denoted as (d_{text}, d_{images}) where d_{text} is a sequence of M words $\{w_j^d\}_{j=1}^M$ and d_{images} is a list of Y images $\{v_j^d\}_{j=1}^Y$.

3.2 Basic Retrieval

We use BM25 as a basic retrieval due to its good performance compared with several ranking models (McDonald et al., 2018; Pang et al., 2017). Since using tweets’ text may be insufficient to find relevant articles, we expand queries’ text by using text extracted from images. For example, in Fig. 1, text extracted from the image is *Breaking News: Obama: "I won't leave if Trump is elected"*. Following (Vosoughi et al., 2018), we use a tool (OCR Space, 2020) to extract text in images. To our knowledge, our work is the first one using text in images to find verified information.

3.3 Multimodal Attention Network (MAN)

MAN has four components: (1) projection layers, (2) textual matching layer (3) visual matching layer and (4) unifying textual and visual information.

3.3.1 Projection Layers

We use two projection layers: one for Glove embeddings and the other one for contextual word embeddings.

Projection layer for Glove embeddings. Each word w , which can be w_i^q or w_j^d , is mapped into a

vector $\mathbf{t} \in \mathbb{R}^{300}$ by a fixed word embedding layer initialized by Glove embeddings (Pennington et al., 2014). Then, the vector \mathbf{t} is projected into $\mathbf{g} \in \mathbb{R}^P$ by a trainable linear layer shown in Eq. 1.

$$\mathbf{g} = \tanh(\mathbf{W}_1 \cdot \mathbf{t} + \mathbf{b}_1) \quad (1)$$

where $\mathbf{W}_1 \in \mathbb{R}^{P \times 300}$, $\mathbf{b}_1 \in \mathbb{R}^P$. P is projection dimensions. After going through the linear layer, we denote $\mathbf{g}_i^q \in \mathbb{R}^P$ and $\mathbf{g}_j^d \in \mathbb{R}^P$ as representations of word w_i^q and word w_j^d , respectively.

Projection layer for contextual word embeddings. Since Glove embeddings do not reflect context of words in queries and articles, we integrate ELMo (Peters et al., 2018) as a static encoder to generate contextual word embeddings. ELMo maps each word w , which can be w_i^q or w_j^d , into a vector $\ell \in \mathbb{R}^{1024}$ which is then projected into $\mathbf{h} \in \mathbb{R}^P$ by a trainable linear layer shown in Eq. 2.

$$\mathbf{h} = \tanh(\mathbf{W}_2 \cdot \ell + \mathbf{b}_2) \quad (2)$$

where $\mathbf{W}_2 \in \mathbb{R}^{P \times 1024}$, $\mathbf{b}_2 \in \mathbb{R}^P$. P is projection dimensions. After going through the linear layer, we denote $\mathbf{h}_i^q \in \mathbb{R}^P$ and $\mathbf{h}_j^d \in \mathbb{R}^P$ as contextual representations of words w_i^q and w_j^d , respectively.

3.3.2 Textual Matching Layer

We derive (1) Glove embeddings interactions, (2) attended interaction matrix and (3) contextual word embedding interactions, and input them to convolution neural networks (CNNs) for feature extraction.

Glove Embeddings Interactions. An article may be relevant to an original tweet if they have overlapping words or similar words. To capture such signals, we use cosine similarity to derive matrix $\mathbf{S} \in \mathbb{R}^{N \times M}$ as shown in Eq. 3.

$$\mathbf{S}_{ij} = \frac{\mathbf{g}_i^{qT} \cdot \mathbf{g}_j^d}{\|\mathbf{g}_i^q\| \times \|\mathbf{g}_j^d\|}, i = 1..N, j = 1..M \quad (3)$$

Let’s look at an example of matrix \mathbf{S} in Fig. 4(a) where x-axis is an article and y-axis is a query. Roughly speaking, matrix \mathbf{S} looks like a *gray-scale image* in which the overlapping phrase ‘at a costume party’ is like a *segment* at the bottom of the image, suggesting the article is relevant to the query. To capture such patterns, CNNs are widely used.

Attended Interaction Matrix. Matrix \mathbf{S} captures overlapping words between a query and an article. However, when word w_i^q is same as word w_j^d , sometime they may not have the same meaning. Thus, we need an attention mechanism to avoid over-reliance on raw similarities in matrix \mathbf{S} . Inspired by Tay et al. (2019), we measure how dissimilar w_i^q

¹Relevance means that the fact-checking article fact-checks the query

²We use fact-checking articles, articles and documents interchangeably

³We use original tweets and queries interchangeably

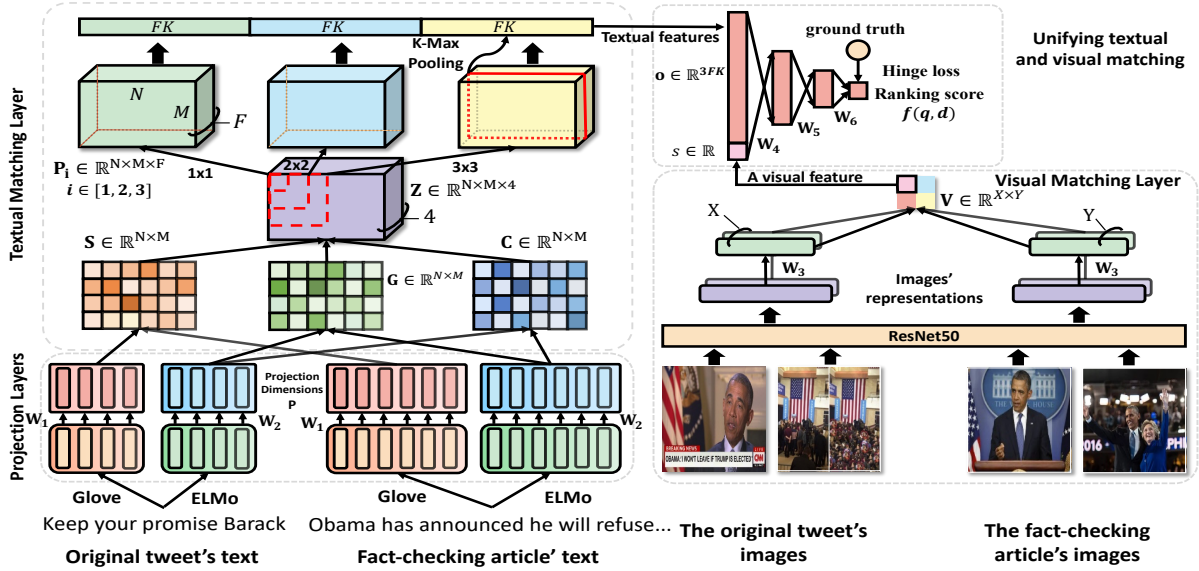


Figure 2: Our proposed model MAN

and w_j^d are based on Euclidean distance between their contextual representations as follows:

$$\mathbf{G}_{ij} = 2 \times \sigma(-\|\mathbf{h}_i^q - \mathbf{h}_j^d\|), i = 1..N, j = 1..M \quad (4)$$

where $\sigma(\cdot)$ is a sigmoid function. Since Euclidean distance is non negative, $\sigma(-\|\mathbf{h}_i^q - \mathbf{h}_j^d\|)$ will be in $(0, 0.5]$ and \mathbf{G}_{ij} will be in $(0, 1]$. Therefore, we can use \mathbf{G}_{ij} to attend to \mathbf{S}_{ij} as follows:

$$\mathbf{A}_{ij} = \mathbf{S}_{ij} \times \mathbf{G}_{ij}, i = 1..N, j = 1..M \quad (5)$$

It is clear to see that when the distance between \mathbf{h}_i^q and \mathbf{h}_j^d is large, \mathbf{G}_{ij} will be closer to 0 which helps downgrade impact of \mathbf{S}_{ij} . From Eq. 5, we can form attended interaction matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$.

To our knowledge, our work is the first one using dissimilarity between contextual word embeddings to attend to interactions of GloVe embeddings.

Contextual Word Embeddings Interactions. In our case studies in Section 6.5, we find that contextual word embeddings are able to capture high similarity between a typo and a normal word (e.g. hillar vs. hillary) while GloVe embeddings fail to do so. To further exploit contextual embeddings, we derive matrix $\mathbf{C} \in \mathbb{R}^{N \times M}$ as follows:

$$\mathbf{C}_{ij} = \frac{\mathbf{h}_i^{qT} \cdot \mathbf{h}_j^d}{\|\mathbf{h}_i^q\| \times \|\mathbf{h}_j^d\|}, i = 1..N, j = 1..M \quad (6)$$

Again, we can view matrix \mathbf{C} as a greyscale image as shown in Fig. 4(d). In addition to cosine similarities, we found that using bilinear function (Rao et al., 2019a) works pretty well as well.

Textual Feature Extraction. We stack matrices \mathbf{S} (Eq. 3), \mathbf{A} (Eq. 5) and \mathbf{C} (Eq. 6) and $\mathbf{S} - \mathbf{C}$ to generate a tensor $\mathbf{Z} \in \mathbb{R}^{N \times M \times 4}$ shown in Eq. 7. The matrix $\mathbf{S} - \mathbf{C}$ is used to make our model aware

of differences between interaction matrices.

$$\mathbf{Z} = [\mathbf{S} \oplus \mathbf{A} \oplus \mathbf{C} \oplus (\mathbf{S} - \mathbf{C})] \quad (7)$$

' \oplus ' denotes matrix stacking. We apply n CNNs on tensor \mathbf{Z} to extract features. The i^{th} CNN is performed with kernel size, stride and the number of filters equal to $i \times i \times 4, 1$ and F , respectively. The output feature map of the i^{th} CNN layer is $\mathbf{P}_i \in \mathbb{R}^{N \times M \times F}, i \in \{1..n\}$. Note, padding zeros are used to ensure \mathbf{P}_i has size of $N \times M \times F$.

Next, we apply k-max pooling on each j^{th} output channel of \mathbf{P}_i denoted as $\mathbf{P}_i[:, :, j] \in \mathbb{R}^{N \times M}$ to generate vector $\mathbf{o}_{i,j} \in \mathbb{R}^K$ as shown in Eq. 8.

$$\mathbf{o}_{i,j} = \text{kmax}(\mathbf{P}_i[:, :, j]), i = 1..n, j = 1..F \quad (8)$$

Finally, nF vectors $\mathbf{o}_{i,j}$ are concatenated to create textual features vector $\mathbf{o} \in \mathbb{R}^{nFK}$ shown in Eq. 9.

$$\mathbf{o} = [\mathbf{o}_{1,1}; \dots; \mathbf{o}_{i,j}; \dots; \mathbf{o}_{n,F}] \in \mathbb{R}^{nFK} \quad (9)$$

3.4 Visual Matching Layer

A fixed pretrained ResNet50 (He et al., 2016) maps an image v , which is either an image of an original tweet v_i^q or an image of a FC-article v_j^d , into vector $\mathbf{v} \in \mathbb{R}^H$ which is then projected into vector $\mathbf{m} \in \mathbb{R}^T$ by a trainable linear layer: $\mathbf{m} = \mathbf{W}_3 \cdot \mathbf{v} + \mathbf{b}_3$, where $\mathbf{W}_3 \in \mathbb{R}^{T \times H}$ and $\mathbf{b}_3 \in \mathbb{R}^T$. H and T are set to 2048 and 300, respectively. After the linear layer, we denote $\mathbf{m}_i^q \in \mathbb{R}^T$ and $\mathbf{m}_j^d \in \mathbb{R}^T$ as representations of v_i^q and v_j^d , respectively.

Intuitively, an article is relevant to a query if the article has images similar to the query's images. Thus, we derive matrix $\mathbf{V} \in \mathbb{R}^{X \times Y}$ of pairwise

similarities of images in Eq. 10.

$$\mathbf{V}_{ij} = \frac{\mathbf{m}_i^{qT} \cdot \mathbf{m}_j^d}{\|\mathbf{m}_i^q\| \times \|\mathbf{m}_j^d\|}, i = 1..X, j = 1..Y \quad (10)$$

Similar to (Rao et al., 2019a,b), we pool the largest pairwise similarity s as a visual feature as follows:

$$s = \max(\mathbf{V}), \text{ where } s \in \mathbb{R} \quad (11)$$

When the article has no images, s is set to -1 .

3.5 Unifying Textual and Visual Information

We unify textual and visual information by appending scalar s (Eq. 11) to vector \mathbf{o} (Eq. 9), denoted as $[\mathbf{o}; s]$, and derive $f(q, d)$ as shown in Eq. 12.

$$f(q, d) = \mathbf{W}_6 \cdot \text{relu}(\mathbf{W}_5 \cdot \text{relu}(\mathbf{W}_4 \cdot [\mathbf{o}; s])) \quad (12)$$

where $\mathbf{W}_4 \in \mathbb{R}^{128 \times (nFK+1)}$, $\mathbf{W}_5 \in \mathbb{R}^{64 \times 128}$ and $\mathbf{W}_6 \in \mathbb{R}^{1 \times 64}$. We remove biases to avoid clutter. Our model is trained on triples consisting of a query q , relevant document d^+ and non-relevant document d^- , minimizing hinge loss in Eq. 13.

$$\mathcal{L}(q, d^+, d^-) = \max(0, 1 - f(q, d^+) + f(q, d^-)) \quad (13)$$

4 Data Collection

Finding FC-articles, which address an original tweet, is laborious since we have to read many FC-articles even when using search engines (Popat et al., 2017, 2018). To reduce labeling efforts, we looked at existing datasets (Jiang and Wilson, 2018; Vosoughi et al., 2018; Vo and Lee, 2019) and found that a dataset in Vo and Lee (2019) met our need. The dataset provides non-anonymized pairs of an original tweet and its reply in which FC-articles, from two major fact-checking sites *snopes.com* and *politifact.com*, are embedded. *Fact-checkers* in Vo and Lee (2019) replied to the original tweet posters with FC-articles as evidence. From the original tweets’ replies, we generate pairs of an original tweet q and a FC-article d . We also only kept original tweets where text and images are both available.

After preprocessing, we have 19,341 original tweet in English and FC-article pairs (q, d) in which there are 18,961 unique original tweets and 2,845 FC-articles. Following Vosoughi et al. (2018), a labeling step is conducted to ensure that in each pair, the article fact-checks the original tweet. We hired native U.S. English speakers since they were more likely to be familiar with topics in the tweets and FC-articles. The labelers labeled each pair (q, d) as 1 if the article d fact checked the tweet q . Otherwise, they labeled it as 0. They were trained directly by the authors and were asked to

label several examples as exercises to ensure that they fully understood the task. We required labelers to read the original tweet’s text, the article’s text and images, and developed a labeling UI to help labelers to quickly explore the linked FC-articles shown in Fig. 5 in our appendix. For each pair, three different labelers labeled it. The final label is based on the majority vote. The Kappa value is 0.56, suggesting moderate agreement among the labelers (Viera et al., 2005).

The moderate agreement between labelers was because there were many pairs of an original tweet and a FC-article where the tweet and the article are topically similar but the article does not fact-check the tweet. For example, the tweet is about Hillary Clinton’s mishandled classified emails while the article fact-checks if she gave uranium to Russia. Both the tweet and the article were about Hillary Clinton but the article did not precisely fact-check the tweet’s content. As we utilized the dataset in (Vo and Lee, 2019) which was collected during the 2016 U.S. presidential election, many tweets and FC-articles were about misinformation related to Hillary Clinton and Donald Trump, leading to topically similar pairs which might confuse labelers. After labeling, we have a full dataset of 13,239 positive pairs made by 13,091 original tweets and 2,170 FC-articles.

We observe that there may be *false negatives* in the full dataset, meaning that a FC-article actually fact-checks an original tweet but the article is viewed as an irrelevant one (i.e., 100% precision but less than 100% recall) because the FC-article was not embedded in a fact-checker’s reply. For example, an original tweet is fact-checked by both a Snopes article and a Politifact one but only the Snopes article was embedded in the fact-checker’s reply to the original tweet while the Politifact one was not included in the reply. If we build a model on the full dataset, the false negatives may mislead our model. To mitigate impact of this problem, we split the full dataset into two sub datasets called Snopes and Politifact datasets. The former one contains pairs where FC-articles are from *snopes.com* and the later one contains pairs where FC-articles are from *politifact.com*. Note, there still may be false negatives in each sub dataset since an original tweet may have multiple fake news stories fact-checked by different articles from the same fact-checking website but a fact-checker did not embed all of the articles in the reply. But the number of

false negatives under this case would be smaller than those in the full dataset. In Snopes dataset, we have 11,202 positive pairs made by 11,167 tweets and 1,703 FC-articles. In PolitiFact dataset, we have 2,037 positive pairs made by 2,026 tweets and 467 FC-articles. There are 102 overlapping tweets between the two datasets. The number of unique original posters is 8,277 and 1,482 in Snopes and PolitiFact respectively. On average, each original poster posted ~ 1.35 tweets in our datasets.

5 Data Analysis

Topics of original tweets/queries. Since the topic of an original tweet is related to the topic of a corresponding FC-article, we extracted topics of relevant FC-articles to understand the topical distribution of tweets. By analyzing each FC-article, top 5 topics of tweets in Snopes are as follows: Politics (42.3%), Fauxtography (22.7%), Junk News (8.1%), Uncategorized (6.8%), Quotes (4.8%). For PolitiFact, tweets’ topics are mostly about politics due to its political mission. In conclusion, our datasets captured various topics.

Similarity of text in tweets and text in images.

As we utilize text in images to enhance ranking performance, we seek to understand how similar text in tweets and text in images. For each query/tweet having text in its images, we transformed its text in tweet and its text in images into two vectors of TF-IDF values, and computed their cosine similarity. From all queries of a dataset, we computed mean cosine similarity. The mean similarity is 0.083 and 0.102 for Snopes and PolitiFact respectively, indicating that text in tweets is less similar to text in images. The number of tweets/queries containing text in images is 8,494 (76%) and 1,742 (86%) for Snopes and PolitiFact respectively.

6 Experiments

6.1 Neural Ranking Baselines

We compare with 9 state-of-the-art neural ranking baselines, divided into 3 groups as follows: (1) multimodal retrieval methods including DVSH (Cao et al., 2016) and TranSearch (Guo et al., 2018), (2) semantic matching models including ESIM (Chen et al., 2017) and NSMN (Nie et al., 2019), and (3) relevance matching methods including Match-Pyramid (Pang et al., 2016), KNRM (Xiong et al., 2017), ConvKNRM (Dai et al., 2018), CoPACRR (Hui et al., 2018) and DUET (Mittra et al., 2017). Please see the appendix for details of the baselines.

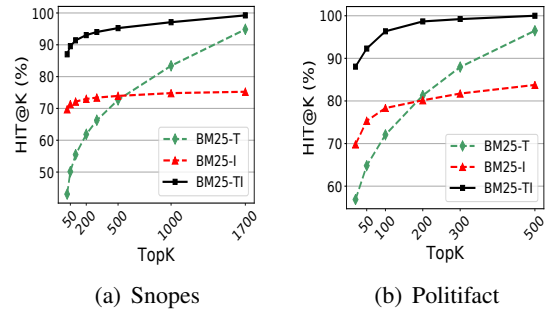


Figure 3: Performance of basic retrieval methods

6.2 Experimental Design

Evaluation Metrics. We adopt NDCG@K (Xiong et al., 2017) and HIT@K (He et al., 2017) as evaluation metrics. We report mean HIT@K and NDCG@K where $K \in [1, 3, 5]$ based on all queries. Since over 99.5% queries have only one relevant document, HIT@K is almost equal to Recall@K. Note, HIT@1 is equal to NDCG@1.

Performance of the Basic Retrieval. We test BM25 in three cases shown in Fig. 3: (1) queries are tweets’ text (BM25-T), (2) queries are text in tweets’ images (BM25-I) and (3) queries are tweets’ text + text in tweets’ images (BM25-TI).

In Fig. 3(a), HIT@50 of BM25-T is only 50% while BM25-I’s HIT@50 is 70%, suggesting that a lot of fake news appear in images. This is because tweets’ text has at most 280 characters. Images are more attractive to online users and easier to convey fake news to them. When K is larger, BM25-I’s HIT@K saturates quickly since only 76% queries have text inside their images. Finally, BM25-TI is the best. Its HIT@50 is 89.6%. Similar patterns appear in PolitiFact in Fig. 3(b). With BM25-TI, HIT@50 is 94%. From these results, we choose BM25-TI as the basic retrieval of our framework.

Split Datasets. We need to choose value of K - the number of initial candidates for each query. If K is too small, initial candidates may not have relevant articles, leading to a meaningless re-ranking step. If K is too large, rerankers’ running time may be high for online apps. We set K to 50 for both datasets. The number of queries, each of which has at least one relevant article in top 50 candidates, is 10,003 out of 11,167 for Snopes and 1,870 out of 2,026 for PolitiFact. Similar to Thorne et al. (2018), from these queries of each dataset, we randomly split them into train, validation and testing sets with ratio 80%/10%/10% as shown in Table 1. There are 1,164 and 156 leftover queries in Snopes and PolitiFact, respectively. Note, having leftover queries is a common issue for re-ranking based systems

Table 1: Split datasets

Datasets	Snopes			Politifact		
	Train	Valid	Test	Train	Valid	Test
Original Tweets	8,002	1,000	1,001	1,496	187	187
FC-Articles	1,703	1,697	1,697	467	467	467

(Thorne et al., 2018). Initial candidates output by BM25-TI are used by all neural ranking models.

Testing Scenarios. All models are tested in the re-ranking step with 2 scenarios (**SC1** and **SC2**). The main difference between them is whether to extract text from images of original tweets and FC-articles, and to incorporate the text with the other information (i.e., text and images of both the tweets and FC-articles). **SC1** and **SC2** are without text from images and with text from images, respectively.

Experimental Settings. For all baselines and our model, we use Adam optimizer (Kingma and Ba, 2014) with learning rate set to 0.001, and early stop training based on HIT@3 and NDCG@3 on a validation set with patience set to 10 epochs. Weight decay of L2-norm is 0.001. Batch size is 16. The number of negative documents sampled for each positive document during training is 3. Maximum |words| in queries is 50 for **SC1** and 100 for **SC2** respectively. Maximum |words| in documents is 1,000. Vocab size V in Snopes dataset is 25,932 in **SC1** and 40,670 in **SC2** respectively. Vocab size V in Politifact dataset is 10,957 in **SC1** and 15,747 in **SC2** respectively. Maximum |images| in queries is 4. Maximum |images| in docs is 17. Images’ shape is (224, 224, 3).

For our model, the number of projection dimensions P is chosen from {64, 128, 256, 512}. The number of output channels F is chosen from {16, 24}. The value of k in kmax pooling is chosen from {16, 32, 48}. The number of CNNs n is chosen from {1, 2, 3}. Our model performs best on Snopes with P , F , k and n equal to 256, 16, 32 and 2, respectively. It performs best on PolitiFact with P , F , k and n equal to 256, 16, 48 and 3, respectively. We implement our model with PyTorch 0.4.1 and test it on a NVIDIA 1080 GTX GPU.

6.3 Performance of Multimodal Attention Network and Variants

We also show results of MAN’s variants as follows: (1) only using text (Eq. 9) and (2) only using images (Eq. 11). We call the former *Contextual Text Matching* (CTM) and the later *Visual Matching Network* (VMN). We show MAN’s improvements wrt. the best result of baselines in each metric.

SC1: Re-ranking using images and text in

tweets. In Table 2, our CTM outperforms the best baselines, achieving maximum improvements of 4.7% on NDCG@1. Our VMN amazingly outperforms text-based ranking baselines in Snopes perhaps because fauxtography is one of the most popular categories on Snopes (Friggeri et al., 2014) while Politifact mainly fact-checks political claims. By using both text and images, our MAN shows an average increase of 17.2% over the best baselines with the maximum improvement of 39.6%.

SC2: Re-ranking using images, tweets’ text and images’ text. We omit VMN from Table 3 since its results are same as Table 2. In Table 3, both our MAN and CTM outperform baselines on two datasets. Interestingly, MAN has lower performance than CTM on Snopes while it has higher performance than CTM on Politifact. We suspect that the abundance of textual signals between original tweets and FC-articles in **SC2** unintentionally makes MAN tend to favor textual signals and neglect visual signals. To remedy this issue, we propose to augment training data in **SC2** with training data in **SC1** while keeping the same validation and testing set from **SC2**. Intuitively, the augmented training data may regularize MAN better (Yu et al., 2018) by letting it observe both rich textual overlapping pairs in **SC2** and pairs with sparse textual signals in **SC1**. We name our model trained under the augmented training data as MAN-A. In Table 3, MAN-A mitigates the above issue with an average increase of 4.8% over the best baselines with the maximum improvement of 11.2%. Text in images has a high impact on performance of CTM and MAN. In Table 3, when using text in images to expand textual content of queries, performance of CTM and MAN increased by 17~34% compared with their performances in Table 2.

From Tables 2 and 3, semantic matching models and multimodal baselines perform worse than relevance matching methods because the first two groups’ goal is to compress whole queries and articles into dense vectors and measure their similarities. However, when compressing textual contents, some irrelevant information may be captured, leading to poor representations (Rao et al., 2019a).

In conclusion, our model MAN outperforms all baselines in both two testing scenarios .

Experiments on the leftover original tweets (i.e., 1,164 tweets in Snopes and 156 tweets in Politifact). We further test benefits of using text and images on each leftover query where we rank its x

Table 2: Performance of our models and baselines when using images and text in tweets

Ranking Models Types	Ranking Models	Snopes					Politifact				
		NDCG@1	NDCG@3	HIT@3	NDCG@5	HIT@5	NDCG@1	NDCG@3	HIT@3	NDCG@5	HIT@5
Exact Matching	BM25-T	0.20579	0.27642	0.32867	0.30420	0.39461	0.18182	0.29162	0.37968	0.31348	0.43316
Multimodal Retrieval (Group 1)	DVSH-B	0.38661	0.51091	0.60040	0.54084	0.67333	0.26203	0.33333	0.38503	0.36003	0.44920
	TransSearch	0.31668	0.46081	0.56444	0.50062	0.66034	0.28342	0.37925	0.44920	0.40040	0.50267
Semantic Matching (Group 2)	ESIM	0.33367	0.46608	0.56444	0.50372	0.65534	0.14973	0.28722	0.39037	0.34871	0.53476
	NSMN	0.45754	0.60097	0.70330	0.63220	0.77822	0.37968	0.47718	0.55080	0.53128	0.67914
Relevance Matching (Group 3)	DUET	0.36863	0.48875	0.57842	0.52628	0.66833	0.29412	0.41009	0.49733	0.43505	0.55615
	MatchPyramid	0.48052	0.58523	0.66034	0.61565	0.73327	0.29412	0.38903	0.45455	0.40812	0.50267
	KNRM	0.48951	0.61081	0.69730	0.63686	0.76124	0.42246	0.54935	0.63636	0.58456	0.72193
	ConvKNRM	0.52148	0.63168	0.70929	0.65942	0.77522	0.45989	0.57229	0.65241	0.62117	0.77005
Ours	CoPACRR	0.53247	0.64469	0.72328	0.67208	0.78921	0.45455	0.59344	0.69519	0.62761	0.77540
	CTM	0.55744	0.67555	0.75624	0.70156	0.81918	0.47059	0.61669	0.71658	0.64292	0.78075
	VMN	0.68931	0.73540	0.76723	0.75019	0.80320	0.24599	0.26821	0.31551	0.28363	0.35829
	MAN	0.74326	0.82197	0.87712	0.83447	0.90609	0.55080	0.65435	0.73262	0.67644	0.78610
MAN vs. the best result of baselines		39.59%	27.50%	21.27%	24.16%	14.81%	19.77%	10.26%	5.38%	7.78%	1.38%

Table 3: Performance of our models and baselines when using images, text in tweets and text in images

Ranking Models Types	Ranking Models	Snopes					Politifact				
		NDCG@1	NDCG@3	HIT@3	NDCG@5	HIT@5	NDCG@1	NDCG@3	HIT@3	NDCG@5	HIT@5
Exact Matching	BM25-TI	0.63736	0.69650	0.73826	0.71058	0.77223	0.27807	0.34928	0.40642	0.38909	0.50267
Multimodal Retrieval (Group 1)	DVSH-B	0.32667	0.46849	0.56843	0.49640	0.63636	0.21925	0.29335	0.34759	0.32626	0.42246
	TransSearch	0.45854	0.58410	0.67433	0.61832	0.75724	0.39572	0.50878	0.58824	0.52397	0.62567
Semantic Matching (Group 2)	ESIM	0.61139	0.70660	0.77323	0.72999	0.83117	0.33155	0.44658	0.52941	0.48617	0.62567
	NSMN	0.78821	0.85732	0.90809	0.87148	0.94106	0.58824	0.70002	0.77540	0.73500	0.86096
Relevance Matching (Group 3)	DUET	0.51848	0.63605	0.71928	0.67075	0.80220	0.41711	0.53087	0.60963	0.55757	0.67380
	MatchPyramid	0.86513	0.91150	0.94406	0.91791	0.95904	0.64171	0.74872	0.82353	0.77702	0.89305
	KNRM	0.84815	0.89118	0.92008	0.90271	0.94805	0.65775	0.75464	0.82353	0.77237	0.86631
	ConvKNRM	0.85914	0.90829	0.94306	0.91401	0.95704	0.66310	0.79163	0.88235	0.80705	0.91979
Ours	CoPACRR	0.86913	0.91166	0.94006	0.91851	0.95604	0.66845	0.77419	0.84492	0.79191	0.88770
	CTM	0.89910	0.93191	0.95504	0.94008	0.97502	0.71123	0.82512	0.89840	0.84331	0.94118
	MAN	0.88412	0.92563	0.95604	0.93238	0.97203	0.72193	0.83104	0.90374	0.85313	0.95722
	MAN-A	0.90909	0.94204	0.96503	0.94892	0.98202	0.74332	0.84905	0.91979	0.85987	0.94652
MAN-A vs. best result of baselines		4.60%	3.33%	2.22%	3.31%	2.40%	11.20%	7.25%	4.24%	6.54%	2.91%

Table 4: Ranking performances on leftover queries when using images, text in tweets and text in images

Ranking Models	Snopes			Politifact		
	NDCG@1	NDCG@3	HIT@3	NDCG@1	NDCG@3	HIT@3
TransSearch	0.20361	0.31856	0.40292	0.12821	0.23542	0.31410
NSMN	0.32646	0.41123	0.47595	0.34615	0.46871	0.55769
MatchPyramid	0.26031	0.33194	0.38488	0.28846	0.34257	0.37821
ConvKNRM	0.29124	0.40280	0.48282	0.36538	0.53479	0.65385
CoPACRR	0.30928	0.40748	0.48024	0.33333	0.43887	0.51923
MAN-A	0.58591	0.68348	0.75258	0.51282	0.64598	0.73718
Impr. MAN-A	79.47%	66.20%	55.87%	40.35%	20.79%	12.74%

Table 5: Effects of contextual word embeddings

Ranking Models	Snopes			Politifact		
	NDCG@1	NDCG@3	HIT@3	NDCG@1	NDCG@3	HIT@3
Glove	0.84216	0.90017	0.94106	0.60428	0.75713	0.86096
ELMo	0.88511	0.92865	0.95804	0.70588	0.80080	0.86631
Glove+ELMo	0.89910	0.93191	0.95504	0.71123	0.82512	0.89840

relevant articles against 50 $-x$ negative documents randomly sampled by following Wan et al. (2016); Wu et al. (2017). It means there are 50 FC-articles per query/tweet. Table 4 shows results of our best model MAN-A and best baselines in each group. As expected, MAN-A outperforms all the baselines due to sparse textual content in leftover queries.

6.4 Effect of Contextual Word Embeddings

To understand effects of word embeddings on our model, we remove visual information and study re-ranking results of our model when (1) using only Glove embeddings, (2) using only contextual word embeddings from ELMo and (3) Glove+ELMo. In

Table 5, when combining Glove and ELMo, we consistently achieve best NDCG in both datasets.

6.5 Case Studies

Qualitative comparison with the best baseline.

An example tweet is ‘You won’t have to wait long ⁴’ embedded with a picture of an Antifa member beating a police officer. Clearly, the tweet’s text does not have any meaningful information while the image contains useful information. Given the tweet, the best baseline CoPACRR failed to find relevant FC-articles, whereas MAN ranked the correct FC-article (Evon, 2017) in top-3 results.

Visualization of interaction matrices and attended matrix Fig. 4 visualizes matrices \mathbf{S} , \mathbf{G} , \mathbf{A} and \mathbf{C} in Eq. 3, 4, 5, 6 respectively of an original tweet and its FC-article from a testing set. Note, these matrices are learned by our model. In Fig. 4(a), Glove embeddings help reveal overlapping phrases (e.g. *at a costume party*, *clinton*) but closeness of *hillar* and *hillary* is not well captured (i.e. $\text{sim}(\text{hillar}, \text{hillary}) = 0.3$). In contrast, $\text{sim}(\text{hillar}, \text{hillary})$ is 0.86 in Fig. 4(d), indicating quality of contextual word embeddings. The matrix \mathbf{G} in Fig. 4(b) has high values for key interactions (e.g. a list of values for *at a costume party* is [0.56,

⁴<https://bit.ly/3ngtsBK>

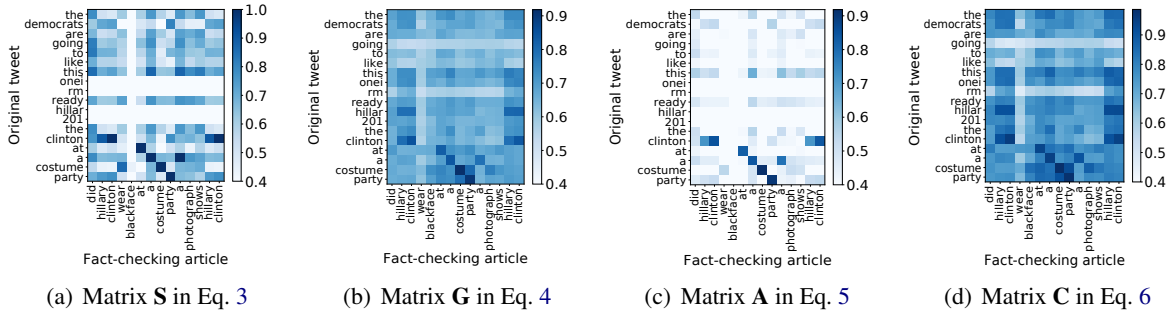


Figure 4: Visualization of matrix S , matrix G , attended matrix A and matrix C (Best viewed in color)

0.66, 0.68, 0.69]) in lieu of uniform values [1.0, 1.0, 1.0, 1.0] in matrix S). When combining matrix S and G , we have a sparse matrix A in Fig. 4(c) which pays more attention to key interactions (e.g. *costume* and *party*). In conclusion, the attention mechanism helps us capture key matching signals. **Impact of Searching for FC-Articles.** We measure how much impact we can make on online users when correct FC-articles are retrieved (i.e. HIT@1 = 1). Totally, our best model, MAN-A, accurately finds FC-articles for 910 original tweets in test set of Snopes dataset. From these tweets, the total number of their retweets is 527,299 and total number of followers of the original posters who posted 910 original tweets is 233M. Roughly speaking, we can inform fact-checked information to millions of users. Security systems can prevent half million shares of fake news in those original tweets.

7 Discussion

Since Snopes and Politifact are the most popular fact-checking sites, building two models for them is an acceptable cost. When facing a real-life social media post, we run two trained models sequentially. If there is no found FC-article, we can inform users that the post is unverified and suggest related pages from verified sites (e.g. governments' sites). When tweets do not have any images, we can use CTM which may find less relevant articles compared with MAN. However, CTM still performed better than the baselines as shown in Tables 2 and 3. We also built our best model (MAN-A) on the full dataset but observed some reduction in NDCG@1 and NDCG@3, but not HIT@3 compared with results of **SC2** on separate datasets maybe because of the *false negatives* described in Section 4. However, our model still outperformed the baselines.

There are a few things our work could be improved. First, our basic retrieval BM25-TI does

not consider images' similarities. To improve BM25-TI, we may combine images' similarities and BM25's score. We leave it as future work. Second, we create train/test data based on unique original tweets. Though there are no retweets and quotes, it is hard to completely ensure all queries's content are unique. However, our settings are applied to all models for fair comparisons. In addition, as shown in Fig. 1, online users tend to re-post fake news. Therefore, it may be reasonable to have similar original tweets' content. Third, we tried to fine tune BERT but did not achieve good results perhaps because we did not have enough data. Interestingly, prior work (Shaar et al., 2020) also had a similar observation when fine-tuning BERT.

8 Conclusions

In this paper, we propose a novel method to alleviate the spread of fake news. By searching for FC-articles and incorporating fact-checked information into social media posts, we can warn users about fake news and discourage them from spreading misinformation. Our framework uses text and images to search for FC-articles, achieving an average increase of 4.8% over best baselines with the maximum improvement of 11.2%. Complementary to fake news detection methods, our method proactively scales up verified content on social media.

Our framework can be used for other multimodal retrieval tasks (e.g. searching for verified sites as we suggested in the previous section).

Acknowledgment

This work was supported in part by NSF grant CNS-1755536, AWS Cloud Credits for Research, and Google Cloud. Any opinions, findings and conclusions or recommendations expressed in this material are the author(s) and do not necessarily reflect those of the sponsors.

References

- Hunt Allcott and Matthew Gentzkow. 2017. Social media and fake news in the 2016 election. *Journal of economic perspectives*, 31(2):211–36.
- Aparna Alluri. 2019. Whatsapp: The 'black hole' of fake news in india's election. <https://www.bbc.com/news/world-asia-india-47797151>.
- Ashoka. 2020. Misinformation spreads faster than coronavirus: How a social organization in turkey is fighting fake news. <https://bit.ly/36qqmmH>.
- Saeid Balaneshin-kordan and Alexander Kotov. 2018. Deep neural architecture for multi-modal retrieval based on joint embedding space for text and images. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pages 28–36.
- Yue Cao, Mingsheng Long, Jianmin Wang, Qiang Yang, and Philip S Yu. 2016. Deep visual-semantic hashing for cross-modal retrieval. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1445–1454.
- Haolan Chen, Fred X Han, Di Niu, Dong Liu, Kunfeng Lai, Chenglin Wu, and Yu Xu. 2018. Mix: Multi-channel information crossing for text matching. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 110–119.
- 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.
- Tao Chen, Xiangnan He, and Min-Yen Kan. 2016. Context-aware image tweet modelling and recommendation. In *Proceedings of the 24th ACM international conference on Multimedia*, pages 1018–1027.
- CNN. 2020. How facebook is combating spread of covid-19 misinformation. <https://cnn.it/3gjtBkg>.
- Zhuyun Dai, Chenyan Xiong, Jamie Callan, and Zhiyuan Liu. 2018. Convolutional neural networks for soft-matching n-grams in ad-hoc search. In *Proceedings of the eleventh ACM international conference on web search and data mining*, pages 126–134.
- David Emery. 2016. President obama confirms he will refuse to leave office if trump is elected. <https://bit.ly/36tbWCp>.
- Dan Evon. 2017. Antifa member photographed beating police officer? <https://bit.ly/31cj2RL>.
- Adrien Friggeri, Lada Adamic, Dean Eckles, and Justin Cheng. 2014. Rumor cascades. In *Eighth International AAI Conference on Weblogs and Social Media*.
- 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 on Conference on Information and Knowledge Management*, pages 55–64.
- Jiafeng Guo, Yixing Fan, Xiang Ji, and Xueqi Cheng. 2019. Matchzoo: A learning, practicing, and developing system for neural text matching. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1297–1300.
- Yangyang Guo, Zhiyong Cheng, Liqiang Nie, Xin-Shun Xu, and Mohan Kankanhalli. 2018. Multi-modal preference modeling for product search. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 1865–1873.
- Aditi Gupta, Hemank Lamba, Ponnurangam Kumaraguru, and Anupam Joshi. 2013. Faking sandy: characterizing and identifying fake images on twitter during hurricane sandy. In *Proceedings of the 22nd international conference on World Wide Web*, pages 729–736.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*, pages 173–182.
- Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Information & Knowledge Management*, pages 2333–2338.
- Kai Hui, Andrew Yates, Klaus Berberich, and Gerard De Melo. 2018. Co-pacrr: A context-aware neural ir model for ad-hoc retrieval. In *Proceedings of the eleventh ACM international conference on web search and data mining*, pages 279–287.
- Kai Hui, Andrew Yates, Klaus Berberich, and Gerard de Melo. 2017. Pacrr: A position-aware neural ir model for relevance matching. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1049–1058.
- Shan Jiang and Christo Wilson. 2018. Linguistic signals under misinformation and fact-checking: Evidence from user comments on social media. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW):1–23.

- Jin-Hwa Kim, Jaehyun Jun, and Byoung-Tak Zhang. 2018. Bilinear attention networks. In *Advances in Neural Information Processing Systems*, pages 1564–1574.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*.
- Shimon Kogan, Tobias J Moskowitz, and Marina Niessner. 2019. Fake news: Evidence from financial markets. Available at SSRN 3237763.
- Katrien Laenen, Susana Zoghbi, and Marie-Francine Moens. 2018. Web search of fashion items with multimodal querying. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, pages 342–350.
- Stephan Lewandowsky, Ullrich KH Ecker, Colleen M Seifert, Norbert Schwarz, and John Cook. 2012. Misinformation and its correction: Continued influence and successful debiasing. *Psychological Science in the Public Interest*, 13(3):106–131.
- Yang Liu and Yi-Fang Brook Wu. 2018. Early detection of fake news on social media through propagation path classification with recurrent and convolutional networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Zhenghao Liu, Chenyan Xiong, and Maosong Sun. 2020. Fine-grained fact verification with kernel graph attention network.
- Jing Ma, Wei Gao, and Kam-Fai Wong. 2018. Rumor detection on twitter with tree-structured recursive neural networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1980–1989.
- Ryan McDonald, George Brokos, and Ion Androutsopoulos. 2018. Deep relevance ranking using enhanced document-query interactions. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1849–1860.
- Bhaskar Mitra, Fernando Diaz, and Nick Craswell. 2017. Learning to match using local and distributed representations of text for web search. In *Proceedings of the 26th International Conference on World Wide Web*, pages 1291–1299.
- 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.
- Sara Ashley O’Brien. 2018. Deepfakes are coming, is big tech ready? <https://cnnmon.ie/2AWCCix>.
- OCR Space. 2020. Free ocr api and online ocr. <https://ocr.space/>.
- Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Shengxian Wan, and Xueqi Cheng. 2016. Text matching as image recognition. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- Liang Pang, Yanyan Lan, Jiafeng Guo, Jun Xu, Jingfang Xu, and Xueqi Cheng. 2017. Deeprank: A new deep architecture for relevance ranking in information retrieval. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 257–266.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of NAACL-HLT*, pages 2227–2237.
- Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2017. Where the truth lies: Explaining the credibility of emerging claims on the web and social media. In *Proceedings of the 26th International Conference on World Wide Web Companion*, pages 1003–1012. International World Wide Web Conferences Steering Committee.
- Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2018. Declare: Debunking fake news and false claims using evidence-aware deep learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 22–32.
- Jinfeng Rao, Linqing Liu, Yi Tay, Wei Yang, Peng Shi, and Jimmy Lin. 2019a. Bridging the gap between relevance matching and semantic matching for short text similarity modeling. 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)*, pages 5373–5384.
- Jinfeng Rao, Wei Yang, Yuhao Zhang, Ferhan Ture, and Jimmy Lin. 2019b. Multi-perspective relevance matching with hierarchical convnets for social media search. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 232–240.
- Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to rank short text pairs with convolutional deep neural networks. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*, pages 373–382.
- Shaden Shaar, Giovanni Da San Martino, Nikolay Babulov, and Preslav Nakov. 2020. That is a known lie: Detecting previously fact-checked claims. In

- Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3607–3618.
- Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Grégoire Mesnil. 2014. A latent semantic model with convolutional-pooling structure for information retrieval. In *Proceedings of the 23rd ACM international conference on conference on information and knowledge management*, pages 101–110.
- Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. 2019. defend: Explainable fake news detection. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 395–405.
- Mark Stencel. 2019. Number of fact-checking outlets surges to 188 in more than 60 countries. <https://bit.ly/36y3S31>.
- Yi Tay, Anh Tuan Luu, Aston Zhang, Shuohang Wang, and Siu Cheung Hui. 2019. Compositional de-attention networks. In *Advances in Neural Information Processing Systems*, pages 6132–6142.
- 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)*, pages 809–819.
- Anthony J Viera, Joanne M Garrett, et al. 2005. Understanding interobserver agreement: the kappa statistic. *Fam med*, 37(5):360–363.
- Nguyen Vo and Kyumin Lee. 2018. The rise of guardians: Fact-checking url recommendation to combat fake news. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 275–284.
- Nguyen Vo and Kyumin Lee. 2019. Learning from fact-checkers: Analysis and generation of fact-checking language. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 335–344.
- Nguyen Vo and Kyumin Lee. 2020. Standing on the shoulders of guardians: Novel methodologies to combat fake news. In *Disinformation, Misinformation, and Fake News in Social Media*, pages 183–210. Springer.
- Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. *Science*, 359(6380):1146–1151.
- Shengxian Wan, Yanyan Lan, Jiafeng Guo, Jun Xu, Liang Pang, and Xueqi Cheng. 2016. A deep architecture for semantic matching with multiple positional sentence representations. In *Thirtieth AAAI Conference on Artificial Intelligence*.
- Bokun Wang, Yang Yang, Xing Xu, Alan Hanjalic, and Heng Tao Shen. 2017. Adversarial cross-modal retrieval. In *Proceedings of the 25th ACM international conference on Multimedia*, pages 154–162.
- William Yang Wang. 2017. “liar, liar pants on fire”: A new benchmark dataset for fake news detection. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 422–426.
- Xuezhi Wang, Cong Yu, Simon Baumgartner, and Flip Korn. 2018a. Relevant document discovery for fact-checking articles. In *Companion Proceedings of the The Web Conference 2018*, pages 525–533.
- Yaqing Wang, Fenglong Ma, Zhiwei Jin, Ye Yuan, Guangxu Xun, Kishlay Jha, Lu Su, and Jing Gao. 2018b. Eann: Event adversarial neural networks for multi-modal fake news detection. In *Proceedings of the 24th acm sigkdd international conference on knowledge discovery & data mining*, pages 849–857.
- Liang Wu and Huan Liu. 2018. Tracing fake-news footprints: Characterizing social media messages by how they propagate. In *Proceedings of the eleventh ACM international conference on Web Search and Data Mining*, pages 637–645.
- Yu Wu, Wei Wu, Chen Xing, Ming Zhou, and Zhoujun Li. 2017. Sequential matching network: A new architecture for multi-turn response selection in retrieval-based chatbots. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 496–505.
- Chenyang Xiong, Zhuyun Dai, Jamie Callan, Zhiyuan Liu, and Russell Power. 2017. End-to-end neural ad-hoc ranking with kernel pooling. In *Proceedings of the 40th International ACM SIGIR conference on research and development in information retrieval*, pages 55–64.
- Di You, Nguyen Vo, Kyumin Lee, and Qiang Liu. 2019. Attributed multi-relational attention network for fact-checking url recommendation. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, pages 1471–1480.
- Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V Le. 2018. Qanet: Combining local convolution with global self-attention for reading comprehension. In *Proceedings of the 6th International Conference on Learning Representations*.
- Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. In *Advances in Neural Information Processing Systems*, pages 9051–9062.

Zhe Zhao, Paul Resnick, and Qiaozhu Mei. 2015. Enquiring minds: Early detection of rumors in social media from enquiry posts. In *Proceedings of the 24th international conference on world wide web*, pages 1395–1405.

Ming Zhu, Aman Ahuja, Wei Wei, and Chandan K Reddy. 2019. A hierarchical attention retrieval model for healthcare question answering. In *The World Wide Web Conference*, pages 2472–2482.

A Appendix for the Reproducibility

A.1 Labeling UI

We developed a labeling UI as shown in Fig. 5 to support labelers to quickly explore linked articles. It includes text and images of original tweets as well as text and images of a FC-article.

A.2 Descriptions and Hyperparameters of the Baselines

Multimodal Retrieval Models. DVSH (Cao et al., 2016) accepts a pair of a multimodal query and a multimodal article, and outputs similarity score. It uses cosine max-margin loss. We also tried to compare with DVSH by using hashcode of queries’ text to match articles’ images and vice versa. However, DVSH did not perform well perhaps because queries’ text and documents’ images may be not semantically similar. We implemented DVSH by ourselves because there is no publicly downloadable code. We set its hidden size to 300 and used AlexNet to extract visual features by following (Cao et al., 2016).

TranSearch (Guo et al., 2018) learns representations of queries by using queries’ text and representations of documents by using text and images of the documents. For TranSearch, we omitted the pretraining step because our datasets do not have *also_viewed* or *buy_after_viewing* information. VGG19 was used to extract visual features by following (Guo et al., 2018). We used the publicly accessible TranSearch implementation.

Semantic Matching Models. We compare with ESIM (Chen et al., 2017) and NSMN (Nie et al., 2019). Both models utilize BiLSTM encoders to learn contextual representations and measure similarity between queries and documents. NSMN also uses contextual word embeddings from ELMo with skip connections for better performance. The pre-trained ELMo⁵ with 93.6M parameters was used for NSMN and our proposed models. Its hidden size was 4,096, and output size was 512 with using 2 highway layers.

In ESIM, a hidden size was set to 300. In NSMN, we set its hidden size to 100 to all BiLSTM layers. We also tried to set its hidden size to {200, 300} but we got out-of-memory error on our GPU because NSMN is memory-intensive due to concatenation of word embeddings, contextual embeddings from

ELMo and multiple BiLSTM layers on our documents with 1,000 tokens.

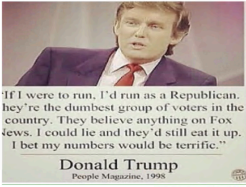
Relevance Matching Models. We compare with several state-of-the-art models in this category. MatchPyramid (Pang et al., 2016) uses CNN to capture spatial patterns. KNRM (Xiong et al., 2017) and ConvKNRM (Dai et al., 2018) use RBF kernel to pool n-gram matching signals. CoPACRR (Hui et al., 2018) uses similarities between queries’ representations and context-aware representations of words in documents to attend to matching signals. DUET (Mitra et al., 2017) unifies semantic and relevance matching signals into one model.

Implementation of MatchPyramid, KNRM, ConvKNRM, and DUET was obtained from MatchZoo (Guo et al., 2019). In MatchPyramid, we used default setting of MatchZoo. The number of kernels of KNRM and ConvKNRM was chosen from {7, 9, 11}. In ConvKNRM, we set |filters| to 300 like the word embeddings’ dimension size. n-gram was chosen from {1, 2, 3}. In DUET, we followed the same architecture proposed in (Mitra et al., 2017).

In CoPACRR, the number of CNN layers was chosen from {2, 3, 4}, the number of filters was chosen from {5, 10, 15}, top k largest matching signals ns was chosen from {5, 10, 15}. The number of segments to conduct k-max pooling $cpos$ was chosen from {4, 5, 6} and context windows were chosen from {5, 9, 11}. We used the publicly accessible CoPACRR implementation.

⁵<https://allennlp.org/elmo>

The sad part is he's right! #NeverTrump #neverhilliary #GOP #RNC <https://t.co/IQkRjOrmpW>



The sad part is he's right! #NeverTrump #neverhilliary #GOP #RNC <https://t.co/IQkRjOrmpW>

If I were to run, I'd run as a Republican. They're the dumbest group of voters in the country. They believe anything on Fox News. I could lie and they'd still eat it up. I bet my numbers would be terrific."

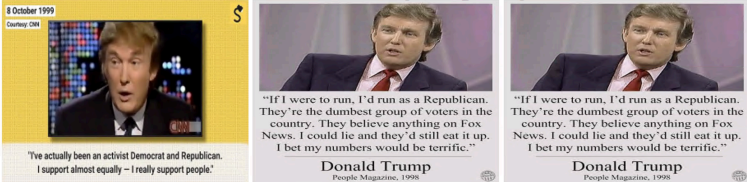
Donald Trump
People Magazine, 1998

<http://pbs.twimg.com/media/CoBFft3XYAAsCeS.jpg>
images_otweet_DataC/31676.png

Text & Images of an original tweet

Fact-checking URL: <http://www.snopes.com/1998-trump-people-quote/#>

Article's Title: Did Donald Trump Say Republicans Are the "Dumbest Group of Voters"?



8 October 1999
Courtney CN

"I've actually been an activist Democrat and Republican. I support almost equally -- I really support people."

"If I were to run, I'd run as a Republican. They're the dumbest group of voters in the country. They believe anything on Fox News. I could lie and they'd still eat it up. I bet my numbers would be terrific."

Donald Trump
People Magazine, 1998

"If I were to run, I'd run as a Republican. They're the dumbest group of voters in the country. They believe anything on Fox News. I could lie and they'd still eat it up. I bet my numbers would be terrific."

Donald Trump
People Magazine, 1998

The fact-checking article's URL, title and images

Figure 5: Labeling UI