# Fully Automated Fact Checking Using External Sources

Georgi Karadzhov[1], Preslav Nakov[2], Lluís Màrquez[2], Alberto Barrón-Cedeño[2], and Ivan Koychev[1]

[1]Sofia University "St. Kliment Ohridski", Bulgaria
[2]Qatar Computing Research Institute, HBKU, Qatar
*georgi.m.karadjov@gmail.com, {pnakov, lmarquez, albarron}@hbku.edu.qa*
*koychev@fmi.uni-sofia.bg*

## Abstract

Given the constantly growing proliferation of false claims online in recent years, there has been also a growing research interest in automatically distinguishing false rumors from factually true claims. Here, we propose a general-purpose framework for fully-automatic fact checking using external sources, tapping the potential of the entire Web as a knowledge source to confirm or reject a claim. Our framework uses a deep neural network with LSTM text encoding to combine semantic kernels with task-specific embeddings that encode a claim together with pieces of potentially-relevant text fragments from the Web, taking the source reliability into account. The evaluation results show good performance on two different tasks and datasets: (*i*) rumor detection and (*ii*) fact checking of the answers to a question in community question answering forums.

## 1 Introduction

Recent years have seen the proliferation of deceptive information online. With the increasing necessity to validate the information from the Internet, *automatic fact checking* has emerged as an important research topic. It is at the core of multiple applications, e.g., discovery of fake news, rumor detection in social media, information verification in question answering systems, detection of information manipulation agents, and assistive technologies for investigative journalism. At the same time, it touches many aspects, such as credibility of users and sources, information veracity, information verification, and linguistic aspects of deceptive language.

In this paper, we present an approach to fact-checking with the following design principles: (*i*) generality, (*ii*) robustness, (*iii*) simplicity, (*iv*) reusability, and (*v*) strong machine learning modeling. Indeed, the system makes very few assumptions about the task, and looks for supportive information directly on the Web. Our system works fully automatically. It does not use any heavy feature engineering and can be easily used in combination with task-specific approaches as well, as a core subsystem. Finally, it combines the representational strength of recurrent neural networks with kernel-based classification.

The system starts with a claim to verify. First, we automatically convert the claim into a query, which we execute against a search engine in order to obtain a list of potentially relevant documents. Then, we take both the snippets and the most relevant sentences in the full text of these Web documents, and we compare them to the claim. The features we use are dense representations of the claim, of the snippets and of related sentences from the Web pages, which we automatically train for the task using Long Short-Term Memory networks (LSTMs). We also use the final hidden layer of the neural network as a task-specific embedding of the claim, together with the Web evidence. We feed all these representations as features, together with pairwise similarities, into a Support Vector Machine (SVM) classifier using an RBF kernel to classify the claim as True or False.

Figure 1 presents a real example from one of the datasets we experiment with. The left-hand side of the figure contains a True example, while the right-hand side shows a False one. We show the original claims from `snopes.com`, the query generated by our system, and the information retrieved from the Web (most relevant snippet and text selection from the web page). The veracity of the claim can be inferred from the textual information.

| | Example 1 | Example 2 |
|---|---|---|
| **original claim** | Texas, teenager Ahmed Mohamed was arrested and accused of creating a "hoax bomb" after bringing a home-assembled clock to school. snopes.com[a] | Is the popular casual dining chain Chipotle closing all their locations soon? snopes.com[c] |
| **generated query** | 'Texas' ● 'Ahmed Mohamed' ● hoax ● bomb ● clock ● arrested ● accused | 'Chipotle' ● dining ● locations ● popular ● closed |
| **best snippet** | knew I'd have a blast playing with [...] Ahmed wasn't accused of making a bomb he was accused of making a look-alike, a hoax [...] it was a bomb, the kid who invented his own digital clock | Chipotle Mexican Grill, Inc is an American chain of fast casual restaurants in the United States, [...] Ellis membership. The restaurant had three locations that operated in 2011 before closing |
| **best webpage sentences** | A 14-year-old Texas student was arrested at school for building a clock. ● A ninth grader was arrested on Sept. 14 just outside Dallas, when he brought a homemade clock to school that teachers and authorities said looked like a bomb Business Insider[b] | Chipotle says it plans to open burger restaurant Tribune news services Chipotle, still struggling to win back customers after a series of food scares, plans to open its first burger restaurant this year. ● The chain known for burritos said Thursday it will open a Tasty Made location this fall in Lancaster, Ohio, which is southeast of Columbus. ● The menu will be limited. Chicago Tribune[d] |
| **label** | Factually TRUE | Factually FALSE |

[a] http://www.snopes.com/2015/09/16/ahmed-mohamed/
[b] http://www.businessinsider.com/ahmed-mohamed-arrested-irving-texas-clock-bomb-2015-9
[c] http://www.snopes.com/chipotle-closing/
[d] http://www.chicagotribune.com/business/ct-chipotle-burger-restaurant-20160728-story.html

Figure 1: Example claims and the information we use to predict whether they are factually true or false.

Our contributions can be summarized as follows:

- We propose a general-purpose light-weight framework for fully-automatic fact checking using evidence derived from the Web.

- We propose a deep neural network with LSTM encoding to combine semantic kernels with task-specific embeddings that encode a claim together with pieces of potentially-relevant text fragments from the Web, taking the source reliability into account.

- We further study factuality in community Question Answering (cQA), and we create a new high-quality dataset, which we release to the research community. To the best of our knowledge, we are the first to study factuality of answers in cQA forums, and our dataset is the first dataset specifically targeting factuality in a cQA setting.

- We achieve strong results on two different tasks and datasets —rumor detection and fact checking of the answers to a question in community question answering forums—, thus demonstrating the generality of the approach and its potential applicability to different fact-checking problem formulations.

The remainder of this paper is organized as follows. Section 2 introduces our method for fact checking claims using external sources. Section 3 presents our experiments and discusses the results. Section 4 describes an application of our approach to a different dataset and a slightly different task: fact checking in community question answering forums. Section 5 presents related work. Finally, Section 6 concludes and suggests some possible directions for future work.

## 2 The Fact-Checking System

Given a claim, our system searches for support information on the Web in order to verify whether the claim is likely to be true. The three steps in this process are (*i*) external support retrieval, (*ii*) text representation, and (*iii*) veracity prediction.

### 2.1 External Support Retrieval

This step consists of generating a query out of the claim and querying a search engine (here, we experiment with Google and Bing) in order to retrieve supporting documents. Rather than querying the search engine with the full claim (as on average, a claim is two sentences long), we generate a shorter query following the lessons highlighted in (Potthast et al., 2013).
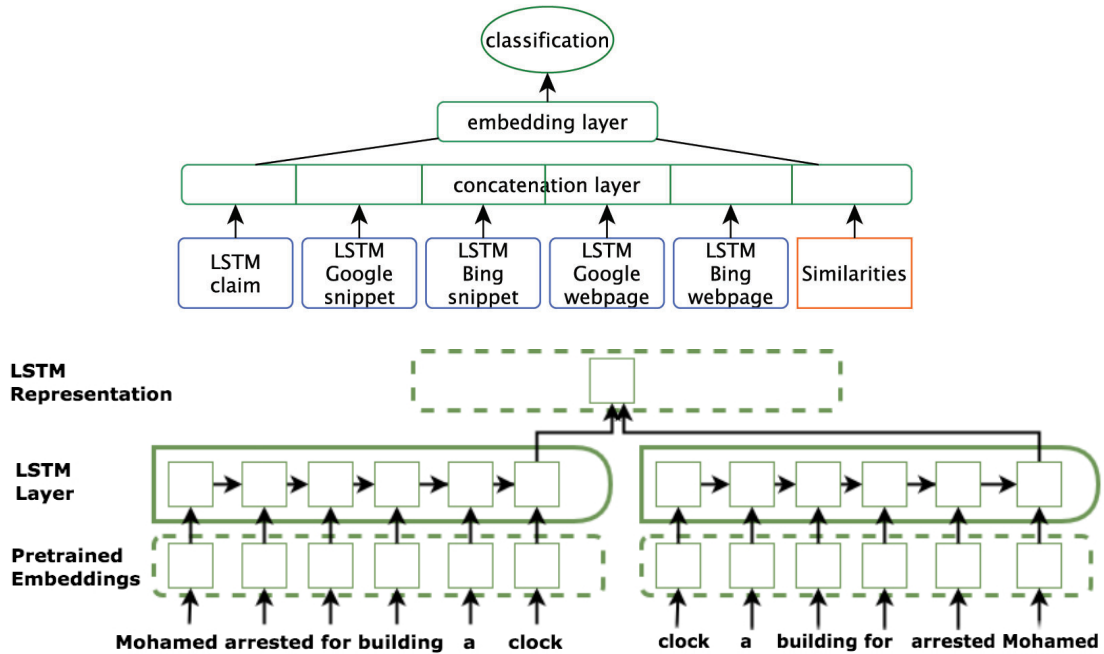
Figure 2: Our general neural network architecture (top) and detailed LSTM representation (bottom). Each blue box in the top consists of the bi-LSTM structure in the bottom.

As we aim to develop a general-purpose fact checking system, we use an approach for query generation that does not incorporate any features that are specific to claim verification (e.g., no temporal indicators).

We rank the words by means of *tf-idf*. We compute the *idf* values on a 2015 Wikipedia dump and the English Gigaword.[1] Potthast et al. (2013) suggested that a good way to perform high-quality search is to only consider the verbs, the nouns and the adjectives in the claim; thus, we exclude all words in the claim that belong to other parts of speech. Moreover, claims often contain named entities (e.g., names of persons, locations, and organizations); hence, we augment the initial query with all the named entities from the claim's text. We use IBM's AlchemyAPI[2] to identify named entities. Ultimately, we generate queries of 5–10 tokens, which we execute against a search engine. We then collect the snippets and the URLs in the results, skipping any result that points to a domain that is considered unreliable.[3] Finally, if our query has returned no results, we iteratively relax it by dropping the final tokens one at a time.

## 2.2 Text Representation

Next, we build the representation of a claim and the corresponding snippets and Web pages. First, we calculate three similarities (a) between the claim and a snippet, or (b) between the claim and a Web page: (*i*) cosine with *tf-idf*, (*ii*) cosine over embeddings, and (*iii*) containment (Lyon et al., 2001). We calculate the embedding of a text as the average of the embeddings of its words; for this, we use pre-trained embeddings from GloVe (Pennington et al., 2014). Moreover, as a Web page can be long, we first split it into a set of rolling sentence triplets, then we calculate the similarities between the claim and each triplet, and we take the highest scoring triplet. Finally, as we have up to ten hits from the search engine, we take the maximum and also the average of the three similarities over the snippets and over the Web pages.

We further use as features the embeddings of the claim, of the best-scoring snippet, and of the best-scoring sentence triplet from a Web page. We calculate these embeddings (*i*) as the average of the embeddings of the words in the text, and also (*ii*) using LSTM encodings, which we train for the task as part of a deep neural network (NN). We also use a task-specific embedding of the claim together with all the above evidence about it, which comes from the last hidden layer of the NN.

## 2.3 Veracity Prediction

Next, we build classifiers: neural network (NN), support vector machines (SVM), and a combination thereof (SVM+NN).

**NN.** The architecture of our NN is shown on top of Figure 2. We have five LSTM sub-networks, one for each of the text sources from two search engines: *Claim*, *Google Web page*, *Google snippet*, *Bing Web page*, and *Bing snippet*. The claim is fed into the neural network as-is. As we can have multiple snippets, we only use the best-matching one as described above. Similarly, we only use a single best-matching triple of consecutive sentences from a Web page. We further feed the network with the similarity features described above. All these vectors are concatenated and fully connected to a much more compact hidden layer that captures the task-specific embeddings. This layer is connected to a softmax output unit to classify the claim as true or false. The bottom of Figure 2 represents the generic architecture of each of the LSTM components. The input text is transformed into a sequence of word embeddings, which is then passed to the bidirectional LSTM layer to obtain a representation for the full sequence.

**SVM.** Our second classifier is an SVM with an RBF kernel. The input is the same as for the NN: word embeddings and similarities. However, the word embeddings this time are calculated by averaging rather than using a bi-LSTM.

**SVM + NN.** Finally, we combine the SVM with the NN by augmenting the input to the SVM with the values of the units in the hidden layer. This represents a task-specific embedding of the input example, and in our experiments it turned out to be quite helpful. Unlike in the SVM only model, this time we use the bi-LSTM embeddings as an input to the SVM. Ultimately, this yields a combination of deep learning and task-specific embeddings with RBF kernels.

## 3 Experiments and Evaluation

### 3.1 Dataset

We used part of the rumor detection dataset created by Ma et al. (2016). While they analyzed a claim based on a set of potentially related tweets, we focus on the claim itself and on the use of supporting information from the Web.

The dataset consists of 992 sets of tweets, 778 of which are generated starting from a claim on `snopes.com`, which Ma et al. (2016) converted into a query. Another 214 sets of tweets are tweet clusters created by other researchers (Castillo et al., 2011; Kwon et al., 2013) with no claim behind them. Ma et al. (2016) ignored the claim and did not release it as part of their dataset. We managed to find the original claim for 761 out of the 778 `snopes.com`-based clusters.

Our final dataset consists of 761 claims from `snopes.com`, which span various domains including politics, local news, and fun facts. Each of the claims is labeled as factually *true* (34%) or as a *false* rumor (66%). We further split the data into 509 for training, 132 for development, and 120 for testing. As the original split for the dataset was not publicly available, and as we only used a subset of their data, we had to make a new training and testing split. Note that we ignored the tweets, as we wanted to focus on a complementary source of information: the Web. Moreover, Ma et al. (2016) used manual queries, while we use a fully automatic method. Finally, we augmented the dataset with Web-retrieved snippets, Web pages, and sentence triplets from Web pages.[4]

### 3.2 Experimental Setup

We tuned the architecture (i.e., the number of layers and their size) and the hyper-parameters of the neural network on the development dataset. The best configuration uses a bidirectional LSTM with 25 units. It further uses a RMSprop optimizer with 0.001 initial learning rate, L2 regularization with $\lambda$=0.1, and 0.5 dropout after the LSTM layers. The size of the hidden layer is 60 with *tanh* activations. We use a batch of 32 and we train for 400 epochs.

For the SVM model, we merged the development and the training dataset, and we then ran a 5-fold cross-validation with grid-search, looking for the best kernel and its parameters. We ended up selecting an RBF kernel with $c = 16$ and $\gamma$ =0.01.

### 3.3 Evaluation Metrics

The evaluation metrics we use are P (precision), R (recall), and $F_1$, which we calculate with respect to the false and to the true claims. We further report AvgR (macro-average recall), $AvgF_1$ (macro-average $F_1$), and Acc (accuracy).

---

[4] All the data, including the splits, is available at `github.com/gkaradzhov/FactcheckingRANLP`

| | False Claims | | | True Claims | | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|
| **Model** | **P** | **R** | **F1** | **P** | **R** | **F1** | **AvgR** | **AvgF$_1$** | **Acc** |
| SVM + NN | 84.1 | 86.3 | 85.2 | 71.1 | 67.5 | 69.2 | **76.9** | **77.2** | **80.0** |
| NN | 79.6 | 92.5 | 85.5 | 77.8 | 52.5 | 62.7 | **72.5** | **74.1** | **79.2** |
| SVM | 75.0 | 86.3 | 80.2 | 60.7 | 42.5 | 50.0 | **64.4** | **65.1** | **71.7** |
| all *false* | 66.7 | 100.0 | 80.0 | – | 0.0 | 0.0 | 50.0 | 40.0 | 66.7 |
| all *true* | – | 0.0 | 0.0 | 33.3 | 100.0 | 50.0 | 50.0 | 25.0 | 33.3 |

Table 1: Results on the *rumor detection dataset* using Web pages returned by the search engines.

## 3.4   Results

Table 1 shows the results on the test dataset. We can see that both the NN and the SVM models improve over the majority class baseline (all false rumors) by a sizable margin. Moreover, the NN consistently outperforms the SVM by a margin on all measures. Yet, adding the task-specific embeddings from the NN as features of the SVM yields overall improvements over both the SVM and the NN in terms of avgR, avgF$_1$, and accuracy. We can see that both the SVM and the NN overpredict the majority class (false claims); however, the combined SVM+NN model is quite balanced between the two classes.

Table 2 compares the performance of the SVM with and without task-specific embeddings from the NN, when training on Web pages vs. snippets, returned by Google vs. Bing vs. both. The NN embeddings consistently help the SVM in all cases. Moreover, while the baseline SVM using snippets is slightly better than when using Web pages, there is almost no difference between snippets vs. Web pages when NN embeddings are added to the basic SVM. Finally, gathering external support from either Google or Bing makes practically no difference, and using the results from both together does not yield much further improvement. Thus, (*i*) the search engines already do a good job at generating relevant snippets, and one does not need to go and download the full Web pages, and (*ii*) the choice of a given search engine is not an important factor. These are good news for the practicality of our approach.

Unfortunately, direct comparison with respect to (Ma et al., 2016) is not possible. First, we only use a subset of their examples: 761 out of 993 (see Section 3.1), and we also have a different class distribution. More importantly, they have a very different formulation of the task: for them, the claim is not available as input (in fact, there has never been a claim for 21% of their examples); rather an example consists of a set of tweets retrieved using *manually* written queries.

| Model | External support | AvgR | AvgF$_1$ | Acc |
|---|---|---|---|---|
| SVM + NN | Bing+Google; pages | 76.9 | 77.2 | **80.0** |
| SVM | Bing+Google; pages | 64.4 | 65.1 | **71.7** |
| SVM + NN | Bing+Google; snippets | 75.6 | 75.6 | **78.3** |
| SVM | Bing+Google; snippets | 68.1 | 69.0 | **74.2** |
| SVM + NN | Bing; pages | 77.5 | 77.0 | **79.2** |
| SVM | Bing; pages | 66.9 | 67.5 | **72.5** |
| SVM + NN | Bing; snippets | 76.3 | 76.4 | **79.2** |
| SVM | Bing; snippets | 68.8 | 69.7 | **75.0** |
| SVM + NN | Google; pages | 73.1 | 74.2 | **78.3** |
| SVM | Google; pages | 63.1 | 63.8 | **71.7** |
| SVM + NN | Google; snippets | 73.1 | 74.2 | **78.3** |
| SVM | Google; snippets | 65.6 | 66.6 | **73.3** |
| baseline (all false claims) | | 50.0 | 40.0 | **66.7** |

Table 2: Results using an SVM with and without task-specific embeddings from the NN on the *Rumor detection dataset*. Training on Web pages vs. snippets vs. both.

In contrast, our system is fully automatic and does not use tweets at all. Furthermore, their most important information source is the change in tweets volume over time, which we cannot use. Still, our results are competitive to theirs when they do not use temporal features.

To put the results in perspective, we can further try to make an indirect comparison to the very recent paper by Popat et al. (2017). They also present a model to classify true vs. false claims extracted from `snopes.com`, by using information extracted from the Web. Their formulation of the task is the same as ours, but our corpora and label distributions are not the same, which makes a direct comparison impossible. Still, we can see that regarding overall classification accuracy they improve a baseline from 73.7% to 84.02% with their best model, i.e., a 39.2% relative error reduction. In our case, we go from 66.7% to 80.0%, i.e., an almost identical 39.9% error reduction. These results are very encouraging, especially given the fact that our model is much simpler than theirs regarding the sources of information used (they model the stance of the text, the reliability of the sources, the language style of the articles, and the temporal footprint).
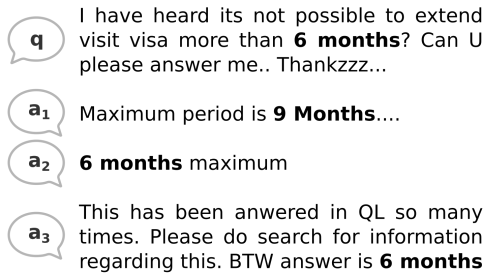
Figure 3: Example from the cQA forum dataset.

| | Label | Answers |
|---|---|---|
| + | FACTUAL - TRUE | 128 |
| − | FACTUAL - PARTIALLY TRUE | 38 |
| − | FACTUAL - CONDITIONALLY TRUE | 16 |
| − | FACTUAL - FALSE | 22 |
| − | FACTUAL - RESPONDER UNSURE | 26 |
| − | NONFACTUAL | 19 |
| | **TOTAL** | **249** |
| + | **POSITIVE** | **128** |
| − | **NEGATIVE** | **121** |

Table 3: Distribution of the answer labels.

## 4 Application to cQA

Next, we tested the generality of our approach by applying it to a different setup: fact-checking the answers in community question answering (cQA) forums. As this is a new problem, for which no dataset exists, we created one. We augmented with factuality annotations the cQA dataset from SemEval-2016 Task 3 (CQA-QA-2016) (Nakov et al., 2016). Overall, we annotated 249 question–answer, or $q$-$a$, pairs (from 71 threads): 128 factually true and 121 factually false answers.

Each question in CQA-QA-2016 has a subject, a body, and meta information: ID, category (e.g., *Education*, and *Moving to Qatar*), date and time of posting, user name and ID. We selected only the factual questions such as "*What is Ooredoo customer service number?*", thus filtering out all (*i*) socializing, e.g., "*What was your first car?*", (*ii*) requests for opinion/advice/guidance, e.g., "*Which is the best bank around??*", and (*iii*) questions containing multiple sub-questions, e.g., "*Is there a land route from Doha to Abudhabi. If yes; how is the road and how long is the journey?*"

Next, we annotated for veracity the answers to the retained questions. Note that in CQA-QA-2016, each answer has a subject, a body, meta information (answer ID, user name and ID), and a judgment about how well it addresses the question of its thread: GOOD vs. POTENTIALLY USEFUL vs. BAD . We only annotated the GOOD answers.[5] We further discarded answers whose factuality was very time-sensitive (e.g., "*It is Friday tomorrow.*", "*It was raining last week.*")[6], or for which the annotators were unsure.

We targeted very high quality, and thus we did not use crowdsourcing for the annotation, as pilot annotations showed that the task was very difficult and that it was not possible to guarantee that *Turkers* would do all the necessary verification, e.g., gathering evidence from trusted sources. Instead, all examples were first annotated independently by four annotators, and then *each example* was discussed in detail to come up with a final label. We ended up with 249 GOOD answers to 71 different questions, which we annotated for factuality: 128 POSITIVE and 121 NEGATIVE examples. See Table 3 for details.

We further split our dataset into 185 $q$–$a$ pairs for training, 31 for development, and 32 for testing, preserving the general positive:negative ratio, and making sure that the questions for the $q$–$a$ pairs did not overlap between the splits.

Figure 3 presents an excerpt of an example from the dataset, with one question and three answers selected from a longer thread. Answer $a_1$ contains false information, while $a_2$ and $a_3$ are true, as can be checked on an official governmental website.[7]

We had to fit our system for this problem, as here we do not have claims, but a question and an answer. So, we constructed the query from the concatenation of $q$ and $a$. Moreover, as Google and Bing performed similarly, we only report results using Google. We limited our run to snippets only, as we have found them rich enough above (see Section 3). Also, we had a list of reputed and Qatar-related sources for the domain, and we limited our results to these sources only. This time, we had more options to calculate similarities compared to the rumors dataset: we can compare against $q$, $a$, and $q$–$a$; we chose to go with the latter. For the LSTM representations, we use both the question and the answer.

---

[5]See (Nakov et al., 2017a) for an overview of recent approaches to finding GOOD answers for cQA.

[6]Arguably, many answers are somewhat time-sensitive, e.g., "*There is an IKEA in Doha.*" is true only after IKEA opened, but not before that. In such cases, we just used the present situation as a point of reference.

[7]https://www.moi.gov.qa/site/english/departments/PassportDept/news/2011/01/03/23385.html

| | False Claims | | | True Claims | | | Overall | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | P | R | F1 | P | R | F1 | AvgR | AvgF1 | Acc |
| SVM + NN | 72.2 | 76.5 | 74.3 | 73.3 | 68.8 | 71.0 | **72.7** | **72.7** | **72.7** |
| SVM | 70.6 | 70.6 | 70.6 | 68.8 | 68.8 | 68.8 | **69.7** | **69.7** | **69.7** |
| NN | 61.1 | 64.7 | 62.9 | 60.0 | 56.3 | 58.1 | **60.5** | **60.5** | **60.6** |
| all *false* | 51.5 | 100.0 | 68.0 | – | 0 | 0 | **50.0** | **34.0** | **51.5** |
| all *true* | – | 0.0 | 0.0 | 48.5 | 100.0 | 65.3 | **50.0** | **32.7** | **48.5** |

Table 4: Results on the cQA answer fact-checking problem.

Table 4 shows the results on the cQA dataset. Once again, our models outperformed all baselines by a margin. This time, the predictions of all models are balanced between the two classes, which is probably due to the dataset being well balanced in general. The SVM model performs better than the NN by itself. This is due to the fact that the cQA dataset is significantly smaller than the *rumor detection* one. Thus, the neural network could not be trained effectively by itself. Nevertheless, the task-specific representations were useful and combining them with the SVM model yielded consistent improvements on all the measures once again.

## 5 Related Work

Journalists, online users, and researchers are well aware of the proliferation of false information on the Web, and topics such as information credibility and fact checking are becoming increasingly important as research directions. For example, there was a recent 2016 special issue of the ACM Transactions on Information Systems journal on Trust and Veracity of Information in Social Media (Papadopoulos et al., 2016), there was a SemEval-2017 shared task on Rumor Detection (Derczynski et al., 2017), and there is an upcoming lab at CLEF-2018 on Automatic Identification and Verification of Claims in Political Debates (Gencheva et al., 2017).

The credibility of contents on the Web has been questioned by researches for a long time. While in the early days the main research focus was on online news portals (Brill, 2001; Finberg et al., 2002; Hardalov et al., 2016), the interest has eventually shifted towards social media (Castillo et al., 2011; Zubiaga et al., 2016; Popat et al., 2017; Karadzhov et al., 2017), which are abundant in sophisticated malicious users such as opinion manipulation *trolls*, paid (Mihaylov et al., 2015b) or just perceived (Mihaylov et al., 2015a; Mihaylov and Nakov, 2016), *sockpuppets* (Maity et al., 2017), *Internet water army* (Chen et al., 2013), and *seminar users* (Darwish et al., 2017).

For instance, Canini et al. (2011) studied the credibility of Twitter accounts (as opposed to tweet posts), and found that both the topical content of information sources and social network structure affect source credibility. Other work, closer to ours, aims at addressing credibility assessment of rumors on Twitter as a problem of finding false information about a newsworthy event (Castillo et al., 2011). This model considers user reputation, writing style, and various time-based features, among others.

Other efforts have focused on news communities. For example, several truth discovery algorithms are combined in an ensemble method for veracity estimation in the VERA system (Ba et al., 2016). They proposed a platform for end-to-end truth discovery from the Web: extracting unstructured information from multiple sources, combining information about single claims, running an ensemble of algorithms, and visualizing and explaining the results. They also explore two different real-world application scenarios for their system: fact checking for crisis situations and evaluation of trustworthiness of a rumor. However, the input to their model is structured data, while here we are interested in unstructured text as input.

Similarly, the task defined by Mukherjee and Weikum (2015) combines three objectives: assessing the credibility of a set of posted articles, estimating the trustworthiness of sources, and predicting user's expertise. They considered a manifold of features characterizing language, topics and Web-specific statistics (e.g., review ratings) on top of a continuous conditional random fields model. In follow-up work, Popat et al. (2016) proposed a model to support or refute claims from `snopes.com` and Wikipedia by considering supporting information gathered from the Web. They used the same task formulation for claims as we do, but different datasets. In yet another follow-up work, Popat et al. (2017) proposed a complex model that considers stance, source reliability, language style, and temporal information.

Our approach to fact checking is related: we verify facts on the Web. However, we use a much simpler and feature-light system, and a different machine learning model. Yet, our model performs very similarly to this latter work (even though a direct comparison is not possible as the datasets differ), which is a remarkable achievement given the fact that we consider less knowledge sources, we have a conceptually simpler model, and we have six times less training data than Popat et al. (2017).

Another important research direction is on using tweets and temporal information for checking the factuality of rumors. For example, Ma et al. (2015) used temporal patterns of rumor dynamics to detect false rumors and to predict their frequency. Ma et al. (2015) focused on detecting false rumors in Twitter using time series. They used the change of social context features over a rumor's life cycle in order to detect rumors at an early stage after they were broadcast.

A more general approach for detecting rumors is explored by Ma et al. (2016), who used recurrent neural networks to learn hidden representations that capture the variation of contextual information of relevant posts over time. Unlike this work, we do not use microblogs, but we query the Web directly in search for evidence. Again, while direct comparison to the work of Ma et al. (2016) is not possible, due to differences in dataset and task formulation, we can say that our framework is competitive when temporal information is not used. More importantly, our approach is orthogonal to theirs in terms of information sources used, and thus, we believe there is potential in combining the two approaches.

In the context of question answering, there has been work on assessing the credibility of an answer, e.g., based on intrinsic information (Banerjee and Han, 2009), i.e., without any external resources. In this case, the reliability of an answer is measured by computing the divergence between language models of the question and of the answer. The spawn of community-based question answering Websites also allowed for the use of other kinds of information. Click counts, link analysis (e.g., PageRank), and user votes have been used to assess the quality of a posted answer (Agichtein et al., 2008; Jeon et al., 2006; Jurczyk and Agichtein, 2007). Nevertheless, these studies address the answers' credibility level just marginally.

Efforts to determine the credibility of an answer in order to assess its overall quality required the inclusion of content-based information (Su et al., 2010), e.g., verbs and adjectives such as *suppose* and *probably*, which cast doubt on the answer. Similarly, Lita et al. (2005) used source credibility (e.g., does the document come from a government Website?), sentiment analysis, and answer contradiction compared to other related answers.

Overall, *credibility* assessment for question answering has been mostly modeled at the feature level, with the goal of assessing the quality of the answers. A notable exception is the work of (Nakov et al., 2017b), where credibility is treated as a task of its own right. Yet, note that *credibility* is different from *factuality* (our focus here) as the former is a subjective perception about whether a statement is credible, rather than verifying it as true or false as a matter of fact; still, these notions are often wrongly mixed in the literature. To the best of our knowledge, no previous work has targeted fact-checking of answers in the context of community Question Answering by gathering external support.

## 6 Conclusions and Future Work

We have presented and evaluated a general-purpose method for fact checking that relies on retrieving supporting information from the Web and comparing it to the claim using machine learning. Our method is lightweight in terms of features and can be very efficient because it shows good performance by only using the snippets provided by the search engines. The combination of the representational power of neural networks with the classification of kernel-based methods has proven to be crucial for making balanced predictions and obtaining good results. Overall, the strong performance of our model across two different fact-checking tasks confirms its generality and potential applicability for different domains and for different fact-checking task formulations.

In future work, we plan to test the generality of our approach by applying it to these and other datasets in combination with complementary methods, e.g., those focusing on microblogs and temporal information in Twitter to make predictions about rumors (Ma et al., 2015, 2016). We also want to explore the possibility of providing justifications for our predictions, and we plan to integrate our method into a real-world application.

## Acknowledgments

## References

Eugene Agichtein, Carlos Castillo, Debora Donato, Aristides Gionis, and Gilad Mishne. 2008. Finding high-quality content in social media. In *Proceedings of the International Conference on Web Search and Data Mining*. Palo Alto, California, USA, WSDM '08, pages 183–194.

Mouhamadou Lamine Ba, Laure Berti-Equille, Kushal Shah, and Hossam M. Hammady. 2016. VERA: A platform for veracity estimation over web data. In *Proceedings of the 25th International Conference Companion on World Wide Web*. Montréal, Québec, Canada, WWW '16, pages 159–162.

Protima Banerjee and Hyoil Han. 2009. Answer credibility: A language modeling approach to answer validation. In *Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Boulder, Colorado, USA, NAACL-HLT '09, pages 157–160.

Ann M Brill. 2001. Online journalists embrace new marketing function. *Newspaper Research Journal* 22(2):28.

Kevin R. Canini, Bongwon Suh, and Peter L. Pirolli. 2011. Finding credible information sources in social networks based on content and social structure. In *Proceedings of the IEEE International Conference on Privacy, Security, Risk, and Trust, and the IEEE International Conference on Social Computing*. Boston, Massachusetts, USA, SocialCom/PASSAT '11, pages 1–8.

Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. 2011. Information credibility on Twitter. In *Proceedings of the 20th International Conference on World Wide Web*. Hyderabad, India, WWW '11, pages 675–684.

Cheng Chen, Kui Wu, Venkatesh Srinivasan, and Xudong Zhang. 2013. Battling the Internet Water Army: detection of hidden paid posters. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. Niagara, Ontario, Canada, ASONAM '13, pages 116–120.

Kareem Darwish, Dimitar Alexandrov, Preslav Nakov, and Yelena Mejova. 2017. Seminar users in the Arabic Twitter sphere. In *Proceedings of the 9th International Conference on Social Informatics*. Oxford, UK, SocInfo '17, pages 91–108.

Leon Derczynski, Kalina Bontcheva, Maria Liakata, Rob Procter, Geraldine Wong Sak Hoi, and Arkaitz Zubiaga. 2017. SemEval-2017 Task 8: RumourEval: Determining rumour veracity and support for rumours. In *Proceedings of the 11th International Workshop on Semantic Evaluation*. Vancouver, Canada, SemEval '17, pages 60–67.

Howard Finberg, Martha L Stone, and Diane Lynch. 2002. Digital journalism credibility study. *Online News Association. Retrieved November* 3:2003.

Pepa Gencheva, Preslav Nakov, Lluís Màrquez, Alberto Barrón-Cedeño, and Ivan Koychev. 2017. A context-aware approach for detecting worth-checking claims in political debates. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*. Varna, Bulgaria, RANLP '17.

Momchil Hardalov, Ivan Koychev, and Preslav Nakov. 2016. In search of credible news. In *Proceedings of the 17th International Conference on Artificial Intelligence: Methodology, Systems, and Applications*. Varna, Bulgaria, AIMSA '16, pages 172–180.

Jiwoon Jeon, W. Bruce Croft, Joon Ho Lee, and Soyeon Park. 2006. A framework to predict the quality of answers with non-textual features. In *Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. Seattle, Washington, USA, SIGIR '06, pages 228–235.

Pawel Jurczyk and Eugene Agichtein. 2007. Discovering authorities in question answer communities by using link analysis. In *Proceedings of the 16th ACM Conference on Conference on Information and Knowledge Management*. Lisbon, Portugal, CIKM '07, pages 919–922.

Georgi Karadzhov, Pepa Gencheva, Preslav Nakov, and Ivan Koychev. 2017. We built a fake news & clickbait filter: What happened next will blow your mind! In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*. Varna, Bulgaria, RANLP '17.

Sejeong Kwon, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. 2013. Prominent features of rumor propagation in online social media. In *Proceedings of the 13th IEEE International Conference on Data Mining*. Dallas, Texas, USA, ICDM '13, pages 1103–1108.

Lucian Vlad Lita, Andrew Hazen Schlaikjer, Wei-Chang Hong, and Eric Nyberg. 2005. Qualitative dimensions in question answering: Extending the definitional QA task. In *Proceedings of the 20th National Conference on Artificial Intelligence*. Pittsburgh, Pennsylvania, USA, AAAI '05, pages 1616–1617.

Caroline Lyon, James Malcolm, and Bob Dickerson. 2001. Detecting short passages of similar text in

large document collections. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Pittsburgh, Pennsylvania, USA, EMNLP '01, pages 118–125.

Jing Ma, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J. Jansen, Kam-Fai Wong, and Meeyoung Cha. 2016. Detecting rumors from microblogs with recurrent neural networks. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence*. New York, New York, USA, IJCAI '16, pages 3818–3824.

Jing Ma, Wei Gao, Zhongyu Wei, Yueming Lu, and Kam-Fai Wong. 2015. Detect rumors using time series of social context information on microblogging websites. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. Melbourne, Australia, CIKM '15, pages 1751–1754.

Suman Kalyan Maity, Aishik Chakraborty, Pawan Goyal, and Animesh Mukherjee. 2017. Detection of sockpuppets in social media. In *Proceedings of the ACM Conference on Computer Supported Cooperative Work and Social Computing*. Portland, Oregon, USA, CSCW '17, pages 243–246.

Todor Mihaylov, Georgi Georgiev, and Preslav Nakov. 2015a. Finding opinion manipulation trolls in news community forums. In *Proceedings of the Nineteenth Conference on Computational Natural Language Learning*. Beijing, China, CoNLL '15, pages 310–314.

Todor Mihaylov, Ivan Koychev, Georgi Georgiev, and Preslav Nakov. 2015b. Exposing paid opinion manipulation trolls. In *Proceedings of the International Conference Recent Advances in Natural Language Processing*. Hissar, Bulgaria, RANLP '15, pages 443–450.

Todor Mihaylov and Preslav Nakov. 2016. Hunting for troll comments in news community forums. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*. Berlin, Germany, ACL '16, pages 399–405.

Subhabrata Mukherjee and Gerhard Weikum. 2015. Leveraging joint interactions for credibility analysis in news communities. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. Melbourne, Australia, CIKM '15, pages 353–362.

Preslav Nakov, Doris Hoogeveen, Lluís Màrquez, Alessandro Moschitti, Hamdy Mubarak, Timothy Baldwin, and Karin Verspoor. 2017a. SemEval-2017 task 3: Community question answering. In *Proceedings of the 11th International Workshop on Semantic Evaluation*. Vancouver, Canada, SemEval '17, pages 27–48.

Preslav Nakov, Lluís Màrquez, Alessandro Moschitti, Walid Magdy, Hamdy Mubarak, abed Alhakim Freihat, Jim Glass, and Bilal Randeree. 2016. SemEval-2016 task 3: Community question answering. In *Proceedings of the 10th International Workshop on Semantic Evaluation*. San Diego, California, USA, SemEval '16, pages 525–545.

Preslav Nakov, Tsvetomila Mihaylova, Lluís Màrquez, Yashkumar Shiroya, and Ivan Koychev. 2017b. Do not trust the trolls: Predicting credibility in community question answering forums. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing*. Varna, Bulgaria, RANLP '17.

Symeon Papadopoulos, Kalina Bontcheva, Eva Jaho, Mihai Lupu, and Carlos Castillo. 2016. Overview of the special issue on trust and veracity of information in social media. *ACM Trans. Inf. Syst.* 34(3):14:1–14:5.

Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. GloVe: Global vectors for word representation. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Doha, Qatar, EMNLP '14, pages 1532–1543.

Kashyap Popat, Subhabrata Mukherjee, Jannik Strötgen, and Gerhard Weikum. 2016. Credibility assessment of textual claims on the web. In *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*. Indianapolis, Indiana, USA, CIKM '16, pages 2173–2178.

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*. Perth, Australia, WWW '17, pages 1003–1012.

Martin Potthast, Matthias Hagen, Tim Gollub, Martin Tippmann, Johannes Kiesel, Paolo Rosso, Efstathios Stamatatos, and Benno Stein. 2013. Overview of the 5th international competition on plagiarism detection. In *Proceedings of the Conference on Multilingual and Multimodal Information Access Evaluation*. Valencia, Spain, CLEF '13, pages 301–331.

Qi Su, Helen Kai-Yun Chen, and Chu-Ren Huang. 2010. Incorporate credibility into context for the best social media answers. In *Proceedings of the 24th Pacific Asia Conference on Language, Information and Computation*. Sendai, Japan, PACLIC '10, pages 535–541.

Arkaitz Zubiaga, Maria Liakata, Rob Procter, Geraldine Wong Sak Hoi, and Peter Tolmie. 2016. Analysing how people orient to and spread rumours in social media by looking at conversational threads. *PLoS ONE* 11(3):1–29.