

# Rumor Detection in Social Networks via Deep Contextual Modeling

Amir Pouran Ben Veyseh

*Department of Computer and Information Science*  
*University of Oregon*  
Eugene, USA  
apouranb@cs.uoregon.edu

My T. Thai

*Department of Computer and Information Science and Engineering*  
*University of Florida*  
Gainesville, USA  
mythai@cise.ufl.edu

Thien Huu Nguyen

*Department of Computer and Information Science*  
*University of Oregon*  
Eugene, USA  
thien@cs.uoregon.edu

Dejing Dou

*Department of Computer and Information Science*  
*University of Oregon*  
Eugene, USA  
dou@cs.uoregon.edu

**Abstract**—Fake news and rumors constitute a major problem in social networks recently. Due to the fast information propagation in social networks, it is inefficient to use human labor to detect suspicious news. Automatic rumor detection is thus necessary to prevent devastating effects of rumors on the individuals and society. Previous work has shown that in addition to the content of the news/posts and their contexts (i.e., replies), the relations or connections among those components are important to boost the rumor detection performance. In order to induce such relations between posts and contexts, the prior work has mainly relied on the inherent structures of the social networks (e.g., direct replies), ignoring the potential semantic connections between those objects. In this work, we demonstrate that such semantic relations are also helpful as they can reveal the implicit structures to better capture the patterns in the contexts for rumor detection. We propose to employ the self-attention mechanism in neural text modeling to achieve the semantic structure induction for this problem. In addition, we introduce a novel method to preserve the important information of the main news/posts in the final representations of the entire threads to further improve the performance for rumor detection. Our method matches the main post representations and the thread representations by ensuring that they predict the same latent labels in a multi-task learning framework. The extensive experiments demonstrate the effectiveness of the proposed model for rumor detection, yielding the state-of-the-art performance on recent datasets for this problem.

**Index Terms**—Rumor, Fake News, Deep Learning, Structure

## I. INTRODUCTION

With the expansion of social networks, the amount of data accessible to users is ever increasing. Social Networks (e.g.,

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

ASONAM '19, August 27-30, 2019, Vancouver, Canada

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6868-1/19/08...\$15.00

<http://dx.doi.org/10.1145/3341161.3342896>

Facebook and Twitter) make it possible for their users to reach out thousands to millions of users. New contents are uploaded constantly to social networks, reflecting a near real-time view about the events in the real world. However, besides the benefit of effective information propagation, social networks present a unique challenge to verify the accuracy of the posts (i.e., fake news and rumors) (dubbed as rumor detection in this work). This problem has attracted much attention from the society recently due to its devastating implication on political, social and economical movements. For instance:

- In U.S. 2016 presidential election almost 529 different rumours about candidates were propagated via Facebook and Twitter which had influence on voters [1].
- A rumour about two explosions in White House in April 2013 in Twitter resulted in social panic and a dramatic drop in stock market [2].
- Rumors about the lost Malaysian airplane in March 2014 in Weibo, a Chinese micro blog service, made it difficult for people to follow the true news and hurt the families of the passengers [3].

It is thus the utmost interest of the social network platforms to develop effective strategies to combat against fake news and rumors. The challenges for this problem come from the fact that various sources of information might be required to recognize the misinformation, and extensive analysis and reasoning of such information sources might be needed to accurately make a decision. In order to solve this problem, the early effort from some social network platforms has relied on the reports from their users to identify the suspicious posts that are then further verified by outside fact checkers (e.g., Facebook) or automatic systems (e.g., Twitter). This approach is not very efficient as it involves substantial human effort and might need much time for the systems to recognize the fake news before they appear. Such long response time might be already sufficient for the fake news to cause serious

problems on the public, calling for more automatic methods for rumor detection to minimize the necessary human effort and expedite the recognition time. Automatic rumor detection is being studied actively in the literature and our work in this paper would contribute to this line of research by introducing a novel method to better address this problem.

There are different definitions for fake news and rumors in the literature. In this work, we will employ the definitions for these concepts as specified in [4]. In general, we can categorize rumors in two different categories:

- The rumors that contain true or false information and they would be indeed proved to be correct or incorrect (respectively) by other authorized sources later (i.e., after the time they are posted).
- The rumor that cannot be verified as presenting true or false information by the authorized sources. However, the users have identified those posts as rumors.

In other words, any piece of information whose veracity status is questionable by the time of posting would be considered as a rumor [4]. Based on this definition for rumors, rumor detection is defined as follows: Given a piece of information from social networks (e.g., user posts), we would like to predict if this piece of information is a rumor or not (i.e., cannot be verified as the true news). Due to the nature of social networks, the information that we can employ in the rumor predictions involves the conversations that the piece of information trigger (e.g., the replies from the users) and the profiles of the users who participate into such conversations. We consider the triggered conversations as the contextual information, and call a post from social networks along with its contextual information as a thread for convenience.

Different approaches based on feature engineering [5] [6], propagation pattern [7], and neural networks [9] have been proposed for rumor detection. It has been demonstrated in these works that besides the content presented in the posts and their corresponding contexts (e.g., the replies), it is also crucial to model the relations between those elements to boost the performance for rumor detection [8]. Unfortunately, these prior work have only focused on the structural relations inherited directly from the social networks (e.g., the reply relation between the main posts and the replies) and failed to exploit the implicit relations that can be induced via the semantics of the posts and the contexts. For instance, the recent work on rumor detection [8] employs the direct tree structures between user posts and replies in Twitter to guide the computation of recursive neural networks to perform prediction.

In this work, we argue that the implicit semantic relations between user posts and replies are also important for rumor detection and the models should capture them appropriately to boost the performance. Consequently, we propose a novel model for rumor detection that explicitly learns the semantic similarities between pairs of the main posts and the contextual replies using the self-attention mechanism. In the proposed model, such semantic similarities are used as the weights to compute representation vectors for the main post and the replies in a thread via the weighted combinations of the other

elements in the same thread. The representation vectors are eventually aggregated into a final representation vector to represent the thread and perform the rumor prediction. Finally, as the main posts in the threads are the most important pieces of information in the threads for rumor detection, we propose to augment the final representations of the threads with the representations of the main posts to avoid the confusion for the proposed model. In particular, we introduce a novel method to ensure that the information in the final representations for the threads also involves the main information in the representation for the main posts. This is achieved by enforcing that such thread and main post representations would predict the same latent labels in a multi-task learning framework for rumor detection. The extensive experiments demonstrate the effectiveness of the proposed model and lead the the state-of-the-art performance for rumor detection on two recent benchmark datasets. In summary, our contributions in this paper include:

- We introduce a novel method for rumor detection based on the explicit modeling of semantic relations between the main posts and the contextual replies in social networks.
- We propose a novel multi-task learning framework to emphasize the main posts in the threads for rumor detection.
- We conduct extensive experiments on recent rumor detection datasets and achieve the state-of-the-art performance on those datasets.

In the following, we will first present the related work and then provide a formal definition for the task of rumor detection. Afterward, we will describe the proposed model and the experiments, followed by a conclusion in the end of this paper.

## II. RELATED WORK

In general we can categorize the previous work into three categories [4]:

- Feature engineering approach: In this category, a set of features is hand-designed to transform the posts into feature representations that are then sent to some statistical model to perform classification. The typical features in this approach include the textual information, the structural evidences [5], [6], image/media content [7], and the propagation patterns of the information diffusion in social networks [10]. The success of this approach depends crucially on the quality of the hand-designed feature sets that might be suboptimal once being applied to different social networks and domains.
- Propagation based approach: In this approach, it is assumed that the propagation pattern of the rumors is different from those for non-rumor posts and such difference can be exploited to detect rumors in social networks [11]. However, One drawback of this approach is that it does not consider the textual and visual information from the post content.
- Deep learning approach: In contrast to feature engineering, this approach automatically learns effective features

for the posts from data via deep learning architectures [8], [9]. The features induced by deep learning often capture the underlying representations of data, thus improving the generalization performance as well as the adaptability to new domains/social networks for rumor detection [9]. Our work in this paper follow this approach to develop a novel deep learning model for rumor detection.

The most related work to ours is the deep learning model to capture contextual information in Twitter for rumor detection in [8]. In particular, Recursive Neural Networks (RvNN) are used to compose the tree-like structures of the posts and the corresponding replies in Twitter based on their *tf-idf* representations. However, such prior work only exploits the explicit relations between the main posts and their replies from the network structures (i.e., the direct reply relations). Our work in this paper is different from the prior work as we go beyond the explicit structural relations in social networks and model the implicit relations among the posts based on their semantics to perform rumor detection.

In addition to the main task of rumor detection, there are other approaches that attempt to detect the stance of the replies toward the main post and then detect the rumor [13], [16]. It has been shown that for this task the evolution of the people’s stance toward the main post is very helpful and considering time series is important for this problem. Motivated by such characteristics, [12] proposes a multi-task learning framework to simultaneously predict the stance and classify the main post for rumor detection.

### III. PROBLEM DEFINITION

Following the recent work on rumor detection [8], we use Twitter as the social network in this work. Formally, an input  $I$  for rumor detection consists of the main tweet  $R_0$  along with a set of reply tweets for this main tweet  $R_1, R_2, \dots, R_T$  (where  $T$  is the number of reply tweets):  $I = (R_0, R_1, R_2, \dots, R_T)$ .  $I$  is called a thread for convenience. The goal of rumor detection is to predict whether the input thread represents a rumor or not. In order to perform such prediction, following [8], we attempt to classify  $I$  into one of the following four labels: 1) Without Rumor, (2) True Rumor, (3) False Rumor and (4) Unrecognizable. If the main tweet  $R_0$  can be proved to be false, the label for  $I$  is “False Rumor” while the label for  $I$  is “True Rumor” if  $R_0$  can be shown to be true.

### IV. MODEL

#### A. Word & Tweet Representation

We consider the input  $I = (R_0, R_1, R_2, \dots, R_T)$  as sequence of tweets where each tweet, in turn, is a sequence of words. The main tweet  $R_0$  is put at the beginning of the tweet sequence. As the tweets might involve different numbers of words, we pad the tweets with a special token to ensure that all the tweets have the same word length  $N$  (i.e., the maximum word length of the tweets in the dataset).

In order to represent the posts and replies, for the  $i$ -th tweet  $R_i$ , we first convert its words  $W_{i1}, W_{i2}, \dots, W_{iN}$  (i.e.,  $R_i = W_{i1}, W_{i2}, \dots, W_{iN}$ ) into their pre-trained word

embeddings  $e_{i1}, e_{i2}, \dots, e_{iN}$  respectively. Afterward, we apply the max-pooling operation over such word embeddings along each dimension to obtain the representation vector  $h_i$  for  $R_i$ :

$$h_i = \text{Elementwise\_Max}(e_{i1}, e_{i2}, \dots, e_{iN}) \quad (1)$$

This tweet-vector transformation procedure would convert the the input thread  $I = (R_0, R_1, R_2, \dots, R_T)$  into a sequence of representation vectors  $(h_0, h_1, h_2, \dots, h_T)$  (respectively) that are then fed into the following steps for further computation.

#### B. Contextualizing Tweet Representations

Due to the nature of social networks, the main tweet and the replies are not independent and the content of the main tweet or replies has substantial influence on other tweets in the same thread. In the previous work, it has been shown that capturing such relations among the main tweets and replies can help to boost the performance of rumor detection [8]. However, such previous work has only considered the explicit relation between the main tweet and its replies (i.e., the reply trees in Twitter from the network structures), neglecting the implicit relations among the tweets based on their semantic similarities. In this work, we propose to exploit such implicit semantic relations to further improve the performance for rumor detection.

In order to capture the semantic relations between the tweets, we learn the pairwise similarities among them based on the self-attention mechanism. In particular, inspired by the transformer architecture in [24], we first compute the key and query vectors for each tweet based on its representation  $h_i$ :

$$\begin{aligned} k_i &= W_k * h_i + b_k \\ q_i &= W_q * h_i + b_q \end{aligned} \quad (2)$$

Given these key and query vectors for the tweets, we obtain the similarity  $a_{ij}$  between the  $i$ -th tweet and the  $j$ -th tweet in the input thread  $I$  via the dot product:

$$a_{i,j} = k_i \cdot q_j / \gamma \quad (3)$$

where  $\gamma$  is a normalization factor. Once all the similarities  $a_{ij}$  for all the tweet pairs in a thread have been computed, we exploit these similarities as the weights to compute more abstract representations for the tweets based on the weighted sums:

$$h'_i = \sum_j a_{i,j} * h_j \quad (4)$$

In the next step, we obtain the overall representation vector  $h'$  for the input thread  $I$  by applying the max-pooling operation one more time over all the tweet representations  $h'_i$ :

$$h' = \text{Elementwise\_Max}(h'_0, h'_1, h'_2, \dots, h'_T) \quad (5)$$

Afterward, a 2-layer feed-forward neural network followed by a softmax layer is employed to produce a probability distribution  $P(y|R_0, R_1, R_2, \dots, R_T; \theta)$  over the possible labels

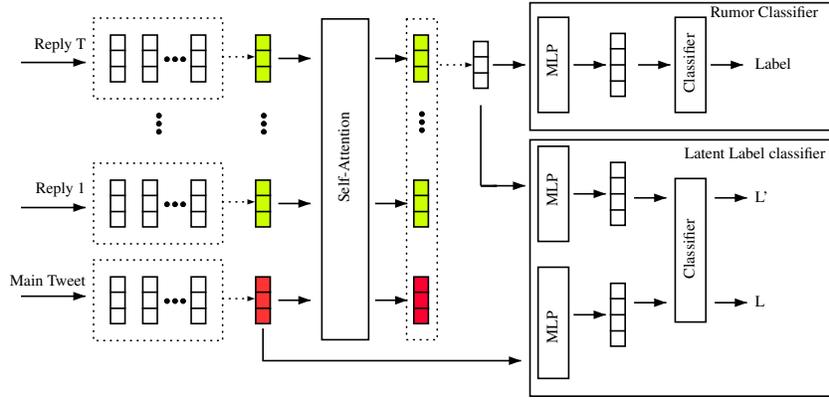


Fig. 1: The Proposed Model. Inputs to the model are word embeddings. Green vectors represent replies and red vectors represent the main tweet. The self-attention component incorporates contextual information into the representation of each tweet. Max-pooled representations are fed into the rumor and latent label classifiers.

for rumor prediction for  $I$  ( $\theta$  is the model parameter). Finally, we optimize the negative log-likelihood function to train the proposed model for rumor detection:

$$L_{label} = -\log P(y^* | R_0, R_1, R_2, \dots, R_T; \theta) \quad (6)$$

where  $y^*$  is the correct label for  $I$ .

### C. Information Preservation

The current model treats the main tweets as equally important as the replies. This is undesirable as the main tweets involve the most important content in the threads over which the model should emphasize to produce good performance. Such emphasis can be done by ensuring that the overall final representations of the threads also cover the core information represented in the main tweet representations (i.e., preserving the information in the main tweets). In this work, we propose to achieve this goal by enforcing the same latent labels induced by the thread representations and the main tweet representations. In particular, we first transform the thread representation  $h'$  and the main tweet representation  $h_0$  into two probability distributions over the same number  $K$  of possible latent labels  $P'(L' | R_0, R_1, R_2, \dots, R_T) = F'(h')$  and  $P_0(L | R_0) = F(h_0)$  (respectively) using two different 2-layer feed-forward neural networks with softmax layers in the end (i.e.,  $F$  and  $F'$  are different 2-layer feed-forward networks). In order to ensure  $h'$  covers the information in  $h_0$ , the latent label  $L'$  that  $P'(L' | R_0, R_1, R_2, \dots, R_T)$  predicts should be the same as the latent label  $L$  predicted by  $P_0(L | R_0)$ . To this end, we additionally optimize the following negative log-likelihood function for the latent labels when we train the proposed model:

$$L = \operatorname{argmax}_L P_0(L | R_0) \quad (7)$$

$$L_{info} = -\log P'(L | R_0, R_1, R_2, \dots, R_T) \quad (8)$$

Consequently, the loss function to train the model in this work is the weighted sum of the rumor label loss and the information preservation loss:

$$Loss = L_{Label} + \alpha L_{info} \quad (9)$$

where  $\alpha$  is the hyper-parameter controlling the contribution of the information preservation loss to the final loss function. Figure 1 shows the main building blocks of the proposed model.

### D. Baseline Models

The goal of the information preservation component in the proposed model is to emphasize main tweet  $R_0$  in the final representation of the input thread  $h'$ . In order to demonstrate the benefits of the latent label prediction mechanism, we explore two other methods to achieve this goal, serving as the baselines for the proposed model.

In the first baseline, we emphasize the main tweet  $R_0$  in  $h'$  by directly imposing that the thread representations  $h'$  and the main tweet representations  $h_0$  are similar. In particular, we first transform  $h'$  and  $h_0$  into more abstract representation vectors  $l'$  and  $l$  (respectively) with the same dimension using different 2-layer feed-forward networks. We would then enforce that  $h'$  and  $h_0$  are similar by replacing  $L_{info}$  in the proposed model with the squared difference loss between  $l'$  and  $l$ :

$$L_{info}^1 = \|l - l'\|^2 \quad (10)$$

where  $\|\cdot\|^2$  is the  $L_2$  norm of the vector. We denote this model by *Diff* in our experiments.

In the second baseline for the information preservation component, we observe that the thread representation  $h'$  is computed based on both the main tweet representation  $h_0$  and the reply tweet representations  $h_1, h_2, \dots, h_T$ . In order to establish a stronger influence of the main tweet on the thread representation, we can ensure that the thread representation  $h'$  is more similar to the main tweet representation  $h_0$  than the reply tweet representations  $h_1, h_2, \dots, h_T$ . In particular, we first aggregate the reply tweet representations  $h_1, h_2, \dots, h_T$  into a single vector  $h_{rep}$  to facilitate the comparison with  $h'$  via the max-pooling operation:

$$h_{rep} = \text{Elementwise\_Max}(h_1, h_2, \dots, h_T) \quad (11)$$

Afterward, we estimate the similarity  $s_0$  between the thread representation  $h'$  and the main tweet representation  $h_0$  as well as the similarity  $s_{rep}$  between the thread representation  $h'$  and the reply tweet representation  $h_{rep}$  by applying a feed-forward neural network  $FF$  with a single output unit (i.e., for the similarity) on the concatenation of the representations:  $s_0 = FF(h', h_0)$ ,  $s_{rep} = FF(h', h_{rep})$ . Note that we ensure the similarities  $s_0$  and  $s_{rep}$  are between 0 and 1 by introducing the sigmoid function in the end of  $FF$ . Consequently, we replace the loss function  $L_{info}$  in Equation 7 of the proposed model with the following margin loss function to push  $h'$  closer to  $h_0$  than  $h_{rep}$ :

$$L_{info}^2 = 1 - s_0 + s_{rep} \quad (12)$$

We name this model as *Discriminator* in the experiments. Note that all the three variants of the information preservation component (i.e., latent label prediction, *Diff* and *Discriminator*) aim to retain more information of the main tweet representation  $h_0$  in the thread representation  $h'$ . However, the label prediction mechanism approaches this via the implicit constraint of the same predicted latent label from the representations while *Diff* targets a more extreme and direct method of similar representations. *Discriminator*, on the other hand, considers the information preservation for the main tweet in the context with the reply tweet representations and the margin loss function.

As the input for rumor detection is a sequence of tweets, a common approach in natural language processing is to handle such sequential data with Recurrent Neural Networks (RNN) (e.g., Long-short Term Memory Networks - LSTM). Consequently, in addition to the baseline models for information preservation, we investigate another baseline model to fine tune the tweet representations based on LSTM. In particular, for this baseline model, we employ a LSTM layer over the sequence of tweet representations  $h_0, h_1, h_2, \dots, h_T$  before the self-attention component. The hidden states of the LSTM layer  $\hat{h}_0, \hat{h}_1, \hat{h}_2, \dots, \hat{h}_T$  would then replace the tweet representations  $h_1, h_2, \dots, h_T$  respectively in the proposed model. This baseline is called *RNN* in the experiments.

## V. EXPERIMENTS

### A. Datasets, Resources & Parameters

Following the previous work on rumor detection, we use the Twitter datasets (described in [12]) for the evaluations in this paper (i.e., the Twitter 15 and Twitter 16 datasets). There are 1381 and 1118 main tweets in these datasets respectively for which a main tweet would correspond to one tree of replies.

Regarding the parameters and resources for the proposed model, we use the Glove [21] embedding (of size 300) to initialize the word vectors. 300 hidden units are employed for the key and query vectors in Equations 2. The feed-forward layer for the rumor classifier has two layers with 200 hidden units. The feed-forward layer for the information preservation

component, on the other hand, has two layers with 100 hidden units and employ three latent labels for prediction. We use the Adagrad optimizer with initial learning rate of 0.3 with the trade-off parameter of  $\alpha = 1$  for the loss function. Following the previous work [8], we use the 5-fold cross validation procedure to tune the parameters and obtain the performance for the models in this work.

### B. Comparing to the state-of-the-art

This section compares the proposed model (called *Semantic Graph*) with the state-of-the-art model on our datasets. For performance measure we use the accuracy on all classes and F1 score per class for each dataset. We compare with two types of models: **1) Feature based models:** These models rely on feature engineering to extract features for such statistical models as Decision Tree, SVM and Random Forest [5], [10], [11], [14], [15], [17]. **2) Deep learning models:** These models use deep learning models to learn features for rumor detection (i.e., Recurrent Neural Networks or Recursive Neural Networks with GRU-RNN, BU-RvNN and TD-RvNN in [8], [9]). Tables II and III show the results on Twitter 15 and Twitter 16 respectively.

These tables show that deep learning models outperform feature based models for rumor detection due to the capacity of the deep learning models to automatically learn effective features from data. In addition, comparing *Semantic Graph* and RvNN models with GRU-RNN, we see that the structural information (e.g., the reply or semantic relations) helps to improve the performance for rumor detection. Finally, by incorporating implicit semantic relations among all the tweets in a thread, *Semantic Graph* achieves the state-of-the-art performance on both datasets in terms of accuracy, and outperforms all the other models in three out of four classes in terms of F1 Score. This clearly demonstrates the effectiveness of the proposed method for rumor detection.

### C. Word Embedding

The input to the proposed model is the embeddings of the words in the tweets. This section compares the performance of the proposed model when different word embeddings are employed. In particular, we investigate two types of word embeddings in this work:

**1) Contextualized Embeddings:** Contextualized word embeddings involve pre-trained models that can compute the embedding for a word based on its context. The following contextualized word embeddings are compared in this work:

- **ELMo:** This is a bidirectional language model with multiple layers of LSTMs [20]. We use the ELMo embeddings of dimension 1024 in the experiments.
- **BERT:** This is a bidirectional language model trained with a transformer [19]. We employ the vectors in the last layer of the BERT model (with dimension 768) for the word embeddings in the experiments.
- **GPT:** Similar to BERT, GPT uses a transformer to train a language model. However, this is an unidirectional model

TABLE I: Heatmap of attention. Numbers in front of each tweet show the index of the tweet in the heatmap. Numbers in heatmap show the attention weights between the corresponding tweets indexed at the columns and rows.

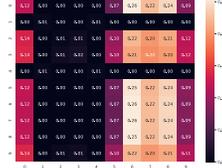
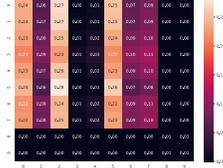
False Rumor	True Rumor
	
<p><b>Main:</b> really? amber alert website goes dark under government shutdown (6)  <b>Reply:</b> How? They only need webmaster to run it (7)  <b>Reply:</b> barry proves it is not (4)  <b>Reply:</b> What's up with this. It's up nicely (13)</p>	<p><b>Main:</b> lego letter from the 1970s still offers a powerful message to parents 40 years later (6)  <b>Reply:</b> lego lives on my brother. my daughter and my grandson all enjoyed it thro the generations (8)  <b>Reply:</b> forty years later still screams from parents as we step on the bloody stuff (7)  <b>Reply:</b> times may change but the truth of this lego letter from the 1970s never has in 40 years via (0)</p>
Unverified	Non-Rumor
	
<p><b>Main:</b> donald trump: " ... laziness is a trait in blacks. it really is, i believe that. it s not anything they can control. (2)  <b>Reply:</b> that is not from trump. it is found in a book by a former trump employee, he alleges he heard trump say that. (0)  <b>Reply:</b> now he is trying to win the african american vote and wonders why it is not working. the gop wed itself to terminal stupidity. (5)  <b>Reply:</b> it is appalling that this nonsensical revisionism is being forcefully persued. (7)</p>	<p><b>Main:</b> colombian gov. and left-wing farc rebel movement agree ceasefire to end decades-long conflict (11)  <b>Reply:</b> hmmm, not exactly as promising as first reported, but a step in the right direction at least. (10)  <b>Reply:</b> I really hope that it will last. (9)  <b>Reply:</b> funny how the farc still stealing children and bombing towns as i write this. (0)</p>

TABLE II: Model Performance on Twitter 15.

Model	Accuracy	F1 NR	F1 FR	F1 TR	F1 UR
DTR [14]	0.409	0.501	0.311	0.364	0.473
DTC [5]	0.454	0.733	0.355	0.317	0.415
RFC [10]	0.565	0.810	0.422	0.401	0.543
SVM-TS [15]	0.544	0.796	0.472	0.404	0.483
SVM-BOW [8]	0.548	0.564	0.524	0.582	0.512
SVM-HK [17]	0.493	0.650	0.439	0.342	0.336
SVM-TK [11]	0.667	0.619	0.669	0.772	0.645
GRU-RNN [9]	0.641	0.684	0.634	0.688	0.571
BU-RvNN [8]	0.708	0.695	0.728	0.759	0.653
TD-RvNN [8]	0.723	0.682	0.758	<b>0.821</b>	0.654
Semantic Graph	<b>0.770</b>	<b>0.814</b>	<b>0.764</b>	0.775	<b>0.743</b>

TABLE III: Model Performance on Twitter 16.

Model	Accuracy	F1 NR	F1 FR	F1 TR	F1 UR
DTR [14]	0.414	0.394	0.273	0.630	0.344
DTC [5]	0.465	0.643	0.393	0.419	0.403
RFC [10]	0.585	0.752	0.415	0.547	0.563
SVM-TS [15]	0.574	0.755	0.420	0.571	0.526
SVM-BOW [8]	0.585	0.553	0.655	0.582	0.578
SVM-HK [17]	0.511	0.648	0.434	0.473	0.451
SVM-TK [11]	0.662	0.643	0.623	0.783	0.655
GRU-RNN [9]	0.633	0.617	0.715	0.577	0.527
BU-RvNN [8]	0.718	0.723	0.712	0.779	0.659
TD-RvNN [8]	0.737	0.662	0.743	<b>0.835</b>	0.708
Semantic Graph	<b>0.768</b>	<b>0.825</b>	<b>0.751</b>	0.768	<b>0.789</b>

[18]. In the experiments, we use the last layer of GPT with dimension of 768 for the word embeddings.

2) **Non-Contextualized Embeddings:** In contrast to the contextualized embeddings, the same embedding vector is

produced for each word in the vocabulary regardless of its context. The following non-contextualized embeddings are considered in our experiments:

- Word2Vec: This is a word embedding method based the skip gram model in [22]. The Word2Vec embeddings in our experiments have 300 dimensions.
- fastText: This is a library for efficient text representation learning. We use the model trained on Wikipedia, UMBC webbase corpus and statmt.org news dataset [23]. The embedding dimension is 300.

TABLE IV: Model Performance on Twitter 15

Model	Accuracy	F1 NR	F1 FR	F1 TR	F1 UR
Word2Vec	0.635	0.648	0.650	0.663	0.641
fastText	0.628	0.621	0.663	0.645	0.632
ELMO	0.698	0.679	0.695	0.712	0.640
BERT	0.729	0.701	0.720	0.756	0.661
GPT	0.767	0.810	<b>0.765</b>	0.771	<b>0.751</b>
Glove	<b>0.770</b>	<b>0.814</b>	0.764	<b>0.775</b>	0.743

TABLE V: Model Performance on Twitter 16

Model	Accuracy	F1 NR	F1 FR	F1 TR	F1 UR
Word2Vec	0.643	0.630	0.629	0.642	0.663
FastText	0.651	0.639	0.626	0.630	0.671
ELMO	0.711	0.683	0.703	0.699	0.715
BERT	0.742	0.801	<b>0.755</b>	<b>0.779</b>	0.770
GPT	0.761	0.821	0.742	0.768	0.783
Glove	<b>0.768</b>	<b>0.825</b>	0.751	0.768	<b>0.789</b>

Table IV and V show the performance of the proposed model with different word embeddings on Twitter 15 and Twitter 16 respectively. These tables show that the contextualized word embeddings in general outperform the non-contextualized word embeddings for rumor detection. This result is expected since the contextualized word embedding capture the appropriate meanings of the words based on their context while non-contextualized models use the same representations for all meanings of the words. However, among all embedding models, Glove actually achieves the best performance. Our hypothesis is that the language in our rumor detection datasets is very similar to those in the corpus used to train Glove (i.e., Wikipedia). This leads to better overlapping between the vocabularies of our datasets and Glove embeddings to provide more prior knowledge for the rumor detection models.

#### D. Ablation Study

There are two major components in the proposed model for rumor detection, i.e., self-attention (called SA) and information preservation with latent label prediction (called *Prediction*). In order to see the contributions of these components for the proposed model, we take turns to exclude these components from the models. The first section in Tables VI and VII shows the performance of the proposed model when the SA and/or *Prediction* component are removed from the model on the Twitter 15 and 16 datasets respectively.

As we can see from the tables, when SA is excluded, the model performance drops dramatically. This shows that contextual information is important for representing each tweet. In addition, the information preservation component with latent label prediction is also necessary for the model to achieve the best performance. It suggests that the multi-task learning setting proposed by this work is effective to preserve important information during the model computation for rumor detection.

In order to assess effectiveness of the proposed latent label prediction for preserving information about the main tweet, we conduct experiments with the *Diff* and *Discriminator* baselines. The results are presented in the second sections of Tables VI and VII, showing that the proposed latent label prediction mechanism outperforms the baselines *Diff* and *Discriminator* for information preservation. We attribute the poorer performance of *Diff* and *Discriminator* to the fact that their constraint mechanisms for the similarity of the main tweet and the thread representations have an effect to eliminate the information about the reply tweets on the final thread representation. In particular, *Diff* enforces the main tweet representation and the thread representation to be similar over all the possible dimensions, leaving no space for the reply tweet representations in the thread representation. *Discriminator*, on the other hand, explicitly separates the reply tweet representations from the overall thread representation. The limited reply tweet information in the thread representation hinders an important source of information and hurts the performance for rumor detection. The latent label prediction mechanism avoids this problem as it only imposes the similarity of the main tweet representation and the thread representation over the most important dimensions via the same latent label, reserving some space for the reply tweet information in the final thread representation for better performance.

Finally the last rows in Tables VI and VII report the performance of the *RNN* baseline mentioned in Section IV.D. As we can see in the tables, although RNN is a common component in many NLP tasks, it hurts the performance of the proposed model for rumor detection. This is due to our consideration of the input thread  $I = (R_0, R_1, R_2, \dots, R_T)$  as a sequence of tweets that do not reflect its original tree structures in social networks. Better modeling approach that inherits such tree structures might help to improve the representation of the tweets and further advance the proposed model.

TABLE VI: Ablation study on Twitter 15

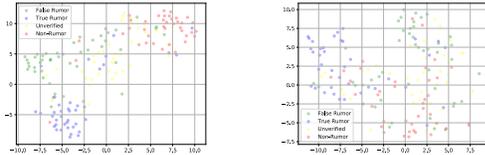
Model	Accuracy	F1 NR	F1 FR	F1 TR	F1 UR
Full Model	<b>0.770</b>	<b>0.814</b>	<b>0.764</b>	<b>0.775</b>	<b>0.743</b>
-SA	0.652	0.643	0.629	0.675	0.660
-Prediction	0.751	0.789	0.756	0.742	0.729
-SA -Prediction	0.601	0.612	0.603	0.663	0.549
Diff	0.755	0.740	0.751	0.772	0.731
Discriminator	0.763	0.791	0.758	0.772	0.740
RNN	0.753	0.746	0.755	0.769	0.740

#### E. Case study on Self-Attention

This section provides a case study on the impact of the self-attention component. Table I shows the attention heatmap for

TABLE VII: Ablation study on Twitter 16

Model	Accuracy	F1 NR	F1 FR	F1 TR	F1 UR
Full Model	<b>0.768</b>	<b>0.825</b>	0.751	<b>0.768</b>	<b>0.789</b>
-SA	0.665	0.643	0.651	0.608	0.689
-Prediction	0.749	0.795	0.723	0.741	0.770
-SA -Prediction	0.592	0.604	0.624	0.598	0.651
Diff	0.756	0.812	0.742	0.751	0.761
Discriminator	0.761	0.819	0.745	0.758	0.769
RNN	0.755	0.809	<b>0.756</b>	0.759	0.772



(a) With Self-Attention (b) Without Self-Attention

Fig. 2: t-SNE diagrams for thread representations.

four input threads in Twitter15 from different classes. These heatmaps illustrate the attention weights between the pairs of tweets in the input threads. One notable point is that there are no direct relations between the reply tweets in such threads based on the inherent structures of the networks. However, from the heatmaps, we see that the reply tweets can be related to each other via the semantic modeling. This helps to enrich the representations for the individual tweets in the threads (based on the self-attention mechanism) to improve the overall performance for rumor detection.

Finally, to further understand the effect of the self-attention component on the proposed model, we draw the t-SNE diagrams for the final representations of the threads in the test datasets when the self-attention component is included and excluded from the model. Figure 2 shows the diagrams. These diagrams suggest that combining the contextual information from the other tweets of the same threads helps to promote the separability of the classes, thus improving the prediction performance for rumor detection.

## VI. CONCLUSION

Fake news and rumors are one of the main problems in social networks. Due to the nature of these networks, rumors can reach out thousands of people quickly and have devastating effects on different aspects of the society. In this paper we introduce a novel method for rumor detection in Twitter. Compared to the previous work, instead of using the explicit structures in the social network structure (i.e., the direct reply relations), our model learns the implicit relations among the main tweet and its replies based on their content. In addition, we present a novel method for preserving important information about the main tweet across the model computations. The extensive experiments prove the effectiveness of the proposed model for rumor detection in Twitter.

## ACKNOWLEDGMENT

This research is partially supported by the NSF grant CNS-1747798 to the IUCRC Center for Big Learning and the HDTRA1-14-1-0055 DTRA grant.

## REFERENCES

- [1] Jin, Zhiwei, Juan Cao, Han Guo, Yongdong Zhang, Yu Wang, and Jiebo Luo. *Detection and analysis of 2016 us presidential election related rumors on twitter*. In SBP-BRiMS, 2017.
- [2] Domm, Patti. *False rumor of explosion at White House causes stocks to briefly plunge; AP confirms its Twitter feed was hacked*. In CNBC.COM, 2013.
- [3] Jin, Zhiwei, Juan Cao, Yu-Gang Jiang, and Yongdong Zhang. *News credibility evaluation on microblog with a hierarchical propagation model*. In ICDM, 2014.
- [4] Cao, Juan, Junbo Guo, Xirong Li, Zhiwei Jin, Han Guo, and Jintao Li. *Automatic rumor detection on microblogs: A survey*. In arXiv preprint arXiv:1807.03505, 2018.
- [5] Castillo, Carlos, Marcelo Mendoza, and Barbara Poblete. *Information credibility on twitter*. In WWW, 2011.
- [6] Yang, Fan, Yang Liu, Xiaohui Yu, and Min Yang. *Automatic detection of rumor on Sina Weibo*. In Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics, 2012.
- [7] Gupta, Manish, Peixiang Zhao, and Jiawei Han. *Evaluating event credibility on twitter*. In SDM, 2012.
- [8] Ma, Jing, Wei Gao, and Kam-Fai Wong. *Rumor detection on twitter with tree-structured recursive neural networks*. In ACL, 2018.
- [9] Ma, Jing, Wei Gao, Prasenjit Mitra, Sejeong Kwon, Bernard J. Jansen, Kam-Fai Wong, and Meeyoung Cha. *Detecting rumors from microblogs with recurrent neural networks*. In IJCAI, 2016.
- [10] Kwon, Sejeong, Meeyoung Cha, Kyomin Jung, Wei Chen, and Yajun Wang. *Prominent features of rumor propagation in online social media*. In ICDM, 2013.
- [11] Ma, Jing, Wei Gao, and Kam-Fai Wong. *Detect rumors in microblog posts using propagation structure via kernel learning*. In ACL, 2017.
- [12] Ma, Jing, Wei Gao, and Kam-Fai Wong. *Detect rumor and stance jointly by neural multi-task learning*. In WWW, 2018.
- [13] Zubiaga, Arkaitz, Elena Kochkina, Maria Liakata, Rob Procter, and Michal Lukasik. *Stance classification in rumours as a sequential task exploiting the tree structure of social media conversations*. In COLING, 2016.
- [14] Zhao, Zhe, Paul Resnick, and Qiaozhu Mei. *Enquiring minds: Early detection of rumors in social media from enquiry posts*. In WWW, 2015.
- [15] Ma, Jing, Wei Gao, Zhongyu Wei, Yueming Lu, and Kam-Fai Wong. *Detect rumors using time series of social context information on microblogging websites*. In CIKM, 2015.
- [16] Amir Veyseh, Javid Ebrahimi, Dejng Dou and Daniel Lowd, *A temporal attentional model for rumor stance classification*. In CIKM, 2017.
- [17] Wu, Ke, Song Yang, and Kenny Q. Zhu. *False rumors detection on sina weibo by propagation structures*. In ICDE, 2015.
- [18] Radford, Alec, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. *Language models are unsupervised multitask learners*. Technical report, OpenAI, 2018.
- [19] Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. *Bert: Pre-training of deep bidirectional transformers for language understanding*. In arXiv preprint arXiv:1810.04805, 2018.
- [20] Peters, Matthew E., Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. *Deep contextualized word representations*. In NAACL-HLT, 2018.
- [21] Pennington, Jeffrey, Richard Socher, and Christopher Manning. *Glove: Global vectors for word representation*. In EMNLP, 2014.
- [22] Mikolov, Tomas, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean. *Distributed representations of words and phrases and their compositionality*. In NIPS, 2013.
- [23] Bojanowski, Piotr, Edouard Grave, Armand Joulin, and Tomas Mikolov. *Enriching word vectors with subword information*. In TACL, (2017)
- [24] Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, ukasz Kaiser, and Illia Polosukhin. *Attention is all you need*. In NIPS, 2017.