

# Rumour verification through recurring information and an inner-attention mechanism<sup>☆</sup>

Ahmet Aker<sup>a,c,\*</sup>, Alfred Sliwa<sup>a</sup>, Fahim Dalvi<sup>b</sup>, Kalina Bontcheva<sup>c</sup>

<sup>a</sup> University of Duisburg-Essen, Germany

<sup>b</sup> QCRI, Qatar

<sup>c</sup> University of Sheffield, England

## ARTICLE INFO

### Article history:

Received 5 April 2019

Revised 22 July 2019

Accepted 22 July 2019

Available online 29 August 2019

MSC:

00-01

99-00

### Keywords:

Rumour Verification

Inner Attention Model

Recurring Terms in Rumours

## ABSTRACT

Verification of online rumours is becoming an increasingly important task with the prevalence of event discussions on social media platforms. This paper proposes an inner-attention-based neural network model that uses frequent, recurring terms from past rumours to classify a newly emerging rumour as *true*, *false* or *unverified*. Unlike other methods proposed in related work, our model uses the source rumour alone without any additional information, such as user replies to the rumour or additional feature engineering. Our method outperforms the current state-of-the-art methods on benchmark datasets (RumourEval2017) by 3% accuracy and 6% F-1 leading to 60.7% accuracy and 61.6% F-1. We also compare our attention-based method to two similar models which however do not make use of recurrent terms. The attention-based method guided by frequent recurring terms outperforms this baseline on the same dataset, indicating that the recurring terms injected by the attention mechanism have high positive impact on distinguishing between true and false rumours. Furthermore, we perform out-of-domain evaluations and show that our model is indeed highly competitive compared to the baselines on a newly released RumourEval2019 dataset and also achieves the best performance on classifying fake and legitimate news headlines.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

Social media platforms now play an important role in the lives of many. They are not only used for social interaction, entertainment and relaxation but also for information seeking and sharing [1]. However, social media platforms are rife with rumours, which following [2] we define as fast-spreading, unverified pieces of information. False rumours can cause major problems not only for everyday users, but also for professionals like journalists and institutions that rely on social information. Thus, in the past few years, automatic rumour verification has become a hot research topic. Most research on rumour analysis has focused on Twitter, as it has established itself as the go-to social platform for real-time news [3] coupled with its comprehensive API. Twitter's unmoderated nature is also the perfect ground for spreading rumours [4].

According to Zubiaga et al. [2], rumour analysis consists of a pipeline of four stages: rumour detection, rumour tracking, stance

classification and rumour verification. Our work focuses on the rumour verification stage. The aim of rumour verification is to verify a rumour and label it as *true*, *false* or *unverified*. To tackle this problem, related work has mainly performed feature engineering and investigated different state-of-the-art machine learning algorithms and learning strategies [5–13]. Unlike the related work, we focus on recurrent patterns in the messages and use these pieces of information to distinguish between the veracity levels.

According to Sen et al. [14] messages during any crisis situation are centred on a finite set of terms related to that specific event. For instance, in the event of an airplane crash, the messages are likely to repeat words related to the crash itself, number of deaths, information related to the pilots, etc. Messages about a terror act will be dominated by words related to the terrorists, number of casualties, location, etc. However, different rumour levels (true, false and unverified) will use different wordings. For instance, false rumours will aim to exaggerate the situation, whereas true rumours would aim to report an incident according to known facts and minimise sentimental overlay.

We argue that these repeating or recurrent pieces of information in the rumour classes are important for the verification task. In this research, we use recurrent words to distinguish between

<sup>☆</sup> Fully documented templates are available in the elsarticle package on CTAN.

\* Corresponding author at: University of Duisburg-Essen, Germany.

E-mail addresses: [a.aker@is.inf.uni-due.de](mailto:a.aker@is.inf.uni-due.de) (A. Aker), [a.sliwa@is.inf.uni-due.de](mailto:a.sliwa@is.inf.uni-due.de) (A. Sliwa), [dalvifahim@gmail.com](mailto:dalvifahim@gmail.com) (F. Dalvi), [k.bontcheva@sheffield.ac.uk](mailto:k.bontcheva@sheffield.ac.uk) (K. Bontcheva).

the different veracity classes. To capture the recurrent words we use a simple but effective method of uni-gram counting. Frequent uni-grams are determined from the training data and are used to teach our model where in a Twitter message to focus. Our model is a combination of an LSTM layer with an attention mechanism.

The recurrent words are extracted from the source tweets only, i.e. our approach does not look at replying tweets. Several previous approaches to rumour verification use stance information of replying tweets [12,13,15]. The advantage of our approach is that the model can be applied directly on the first message reporting a rumour and does not wait until a certain number of replying tweets appear.

Compared to current state-of-the-art performances on known benchmark data (RumourEval2017), our results show 3% gain in accuracy and 6% in  $F-1$  – so far the best reported figures are 57.1% accuracy and 55.8%  $F-1$  [13]. To also show that the results are due to our use of the recurrent words, we compare our model to a vanilla LSTM as well as soft attention approach where recurrent words are not used. Again, the model with the recurrent words leads to best performance. We also test the model on two out-of-domain datasets. The first out-of-domain dataset is the newly released RumourEval2019 dataset that focuses on natural disasters. The second dataset entails headlines collected from fake and legitimate news. Our results show that our inner attention approach achieves satisfactory results on the RumourEval2019 dataset and reaches the best performance on the headlines dataset.

In the following section, we discuss related work. In Section 3 we describe the task as well as the data used in our experiments. Our method is presented in Section 4. Section 5 describes our experimental settings. Results are given in Section 6 and discussed in Section 7. We conclude our paper with immediate future directions in Section 8.

## 2. Related work

The veracity classification task aims to determine whether a given rumour can be confirmed as true, debunked as false, or in some cases marked as yet to be resolved (i.e. unverified). Related work has tackled the problem in a supervised fashion by applying state-of-the-art machine learning algorithms on features extracted from rumour datasets. The pioneering paper of [5] proposed message, user, topic and propagation-based features. In most subsequent studies these features have been used as baselines. Following these feature sets, Kwon et al. [6,16] proposed a new set of feature categories: temporal, structural and linguistic, and showed their importance in verifying rumours. Other than Twitter, the Chinese microblogging platform Sina Weibo also has been analysed for rumours. For this service, Yang et al. [17] proposed client and location-based features and showed that these help to increase prediction accuracy. Further studies modelled features over time [9], investigated other ways of determining the features such as from propagation trees [8,18] or applied different techniques to model them such as using Hidden Markov Models [7]. Chen et al. [19] treated rumour veracity classification as an anomaly detection problem where false rumours are regarded as anomalies. Several features relating to the content, crowd opinion and post propagation were used. Chang et al. [20] put the emphasis on the characteristics of users who post the rumours to determine the veracity. Tong et al. [21] aimed at curtailing false rumour spread rather than marking tweets as true or false. Motivated by the fact that later corrections are not as effective, the authors argued that the first post seen by a user is influential for their future opinion and thus it is important to show users rumours only once they are confirmed to be true.

Rumour veracity classification has also been studied in the RumourEval shared task at SemEval 2017 [22]. Subtask B consisted

in determining if each rumour in the dataset was *true*, *false* or remained *unverified*. Participants viewed the task either as a three-way [12,23,24] or two-way [25,26], single tweet classification task. The winning system [12] added features more specific to the distribution of stance labels in the tweets replying to the source tweet (percentage of reply tweets classified as either support, deny or query). The authors report an accuracy of 53.6% on the RumourEval dataset. Authors of [13] performed a re-implementation of this winning system and reported a slightly better performance (57.1%). In addition, the authors propose a multi-task-learning model where they learn veracity, stance and rumour detection at the same time and reported again a performance of 57.1% accuracy. Note, in both studies [12,13] stance from the responding tweets played a major role in achieving the score of 57.1%. The power of stance to verify rumours was also recently reported by Dungs et al. [15]. The survey paper of Zubiaga et al. [2] provides an extensive summary of current work on rumour verification, along with related tasks such as detection of rumours and stance classification of messages involved in rumours.

Unlike the related work we do not perform explicit feature engineering, nor do we rely on the wisdom from the crowd, e.g. stance information. Our approach works only with the source tweet and does not look at replying messages. The advantage of our method over the previous studies, especially those which rely on replying tweets and are forced to wait until a certain number of replies appear, is that it can work without any delays and therefore be applied in real time – whenever a source tweet appears, our approach can judge it in terms of veracity. In our approach, we make use of recurrent words extracted from source tweets and use them in an attention mechanism to steer the model direction to areas which are distinctive in different rumour levels. Note, attention mechanism has been applied by Chen et al. [27] to perform early rumour detection. However, unlike these authors, we apply inner-attention whereas the authors applied soft attention on top of the LSTM layers.<sup>1</sup> The idea of recurrent words or patterns have been also applied to perform rumour detection [5,28]. For instance, Zhao et al. [28] use manually constructed patterns like *<is (that / this / it) true>* to detect tweets which enquiry the truth of information conveyed by some rumour spreading tweets. The patterns are used as features in traditional machine learning settings. We extract our patterns automatically from the training data and inject them through attention directly to our model. In addition, we use them to verify rumours instead of, as done by Zhao et al. [28], detecting or tracking rumours.

## 3. Task and data

### 3.1. Task: rumour verification

Rumour verification refers to the task of determining a rumour source post being either *true*, *false* or *unverified* and thus is treated as a 3-label classification task [2]. The label *unverified* is used when there is not enough information to verify the message. The input to the task is a source tweet reporting a rumour. Optionally, the source tweet may have replying tweets. Based on the source tweet and the optional replying tweets the task is to label the rumour with a class label.

### 3.2. Data

We use two rumour datasets for the rumour verification task. Both datasets include conversation threads in Twitter about several

<sup>1</sup> We implemented a soft attention method as baseline to compare it with inner-attention approach.



Fig. 1. RumourEval2017: example conversation structure.

**Table 1**  
Distribution of RumourEval2017 dataset.

Split set	# Rumours	# True	# False	# Unverified
Train	272	127	50	95
Dev	25	10	12	3
Test	28	8	12	8
<b>Total</b>	<b>325</b>	<b>145</b>	<b>74</b>	<b>106</b>

events. Each conversation thread is represented in a tree structure with the source tweet as a root node and several branches of replying tweets. A reply chain is a sequence of tweets from root node to a leaf node.

### 3.2.1. RumourEval2017

Our first data resource is the RumourEval2017 dataset [29] which is derived from the PHEME dataset [30]. The RumourEval2017 dataset contains 325 Twitter conversation threads discussing rumours with respect to eight different events like Germanwings Air Crash, Charlie Hebdo, Ottawa Shootings, etc. All the events are man made. Each thread in the dataset is annotated as *true*, *false* or *unverified*. Also each reply to a source tweet is annotated with one of the labels: *supporting*, *denying*, *questioning* and *commenting*. An example conversation structure is shown in Fig. 1. As it was used for the RumourEval2017 challenge it is split into training, development and held-out test set. Table 1 shows the class distribution for the three sets.

Note the split of the dataset into training, development and testing sets is performed randomly. Also note that the rumours in the dataset are independent of each other so that there is no timely development between the rumours. This is also the case for the RumourEval2019 dataset (see next section).

### 3.2.2. RumourEval2019

Our second data contains rumours about natural disasters such as hurricanes, tornados, floods, etc. and thus has a completely different domain focus as the RumourEval2017 data. The structure

**Table 2**  
Distribution of RumourEval2019 dataset.

Split set	# Rumours	# True	# False	# Unverified
Test	56	22	30	4

of the data is similar to the RumourEval2017 and contains conversations with a source tweet that starts a rumour and replies reacting to the source tweet. Also similar to the RumourEval2017 data, the veracity levels are *true*, *false* and *unverified* and have been also annotated for stance (*supporting*, *denying*, *questioning* and *commenting*). This data was adopted by the RumourEval2019 challenge as testing data. In this work we refer to this data as the RumourEval2019 dataset. Statistical details about it are shown in Table 2.

### 3.2.3. Headlines

In addition to the rumour datasets described in the previous sections, we also used news headlines to validate our model. More precisely we have news headlines extracted from manually determined fake news articles and legitimate articles. Here, the idea is to validate our model whether it is able to distinguish legitimate headlines which follow rather a reporting style from fake news headlines that use catchy and exaggerating terms.

We use the data repository FakeNewsSet<sup>2</sup> provided by Shu et al. [31] for the purpose of fake news headline instances. For the legitimate news we took random 100 articles from the Guardian and the independent and used their headlines as legitimate instances. The domains of the articles – fake and non-fake – are from the political and economic areas. In total we have 100 headlines from the fake articles and 100 headlines from the true news. We label the fake news related headlines as *false* and the headlines from the legitimate articles with *true*. Example of true and false news headlines are shown in Table 3.

<sup>2</sup> <https://github.com/KaiDMML/FakeNewsNet>.

**Table 3**  
Example fake and non-fake headlines.

Veracity class	Example headlines
False	WHOA! NEW DISTURBING VIDEO Shows HILLARY'S Campaign Likely FAKED Her Audience At NC Rally
True	New car sales plunge 20% in September

**Table 4**  
Example recurring terms extracted for each veracity class from RumourEval 2017 dataset.

Veracity class	Example terms
False	Rumours, terrorist, #muslim, internet, horrifying, false, cold, fake
True	Live, died, freed, free, follow, identified
Unverified	Alleged, disappeared, military, give

## 4. Method

### 4.1. Core idea

An event causes many social media users to form groups around related issues and hashtags [32]. These groups use language that is narrow and is centred on a finite set of terms related to the event [14], i.e. there are terms that repeat or recur among the messages. We argue that these recurring terms carry useful information to distinguish between rumour classes. In Table 4 we list some frequent uni-grams for different rumour classes. In the *true* rumour case we can see that the terms indicate some fact reporting behaviour such as “live, died, freed, identified”, etc. However, unlike *true* rumours the *false* ones have a different style. They tend to contain words commonly used to exaggerate like “horrifying, terrorist”, etc. This indicates that they rely on emotionally coloured language aiming to fuel fear in readers. Interestingly, false rumours even contain terms like “rumours, false, fake”, which are used to attack other pieces of information reporting about the event and claim that they are all false. This has been also reported by earlier work such as Starbird [33] who reports that false messages accuse true ones of reporting fake information. For the *unverified* messages the language is vague and thus difficult to classify as one of the previous classes – probably this is also why the messages in this class were marked as unverified. However, as in the true rumours the terms in the unverified messages seem to follow a reporting behaviour such as “alleged, disappeared” but the message is not as clear as in the case of true rumour. Our aim is to capture these recurring terms and use them to guide our classification model. We capture the recurring terms using the frequency distribution of uni-grams observed in the training data.

### 4.2. Inner-attention LSTM

In our work, we tackle the veracity classification using an inner-attention based LSTM neural network [34]. Our model resembles the ones used in earlier work for question answering [35] and topic specific argument retrieval [36]. Fig. 2 shows the architecture of our inner-attention based LSTM network.

This neural network is able to learn the importance weighting of source tweet words with respect to a given set of recurring terms. Both the source words as well as the recurring words are represented using word embeddings (Word2Vec [37] or Glove [38] trained on the Twitter dataset<sup>3</sup>). For each source word, an individual weighting factor is learned denoting the importance of that word with respect to the recurring terms. The following

formula is used to obtain this weighting factor  $\alpha_i$  for each word embedding  $x_i$ :

$$\alpha_i = \sigma(u^T W_s x_i) \quad (1)$$

where  $u$  is the averaged recurring terms embedding, and  $W_s$  are the parameters of the embedding mechanism. Each word embedding  $x_i$  is then multiplied by its respective relevance weight factor  $\alpha_i$  using the below equation.

$$\bar{x}_i = \alpha_i x_i \quad (2)$$

In this way, non-relevant source words are filtered out or weighted less. This leads to weighted source tweet word embeddings, on which a Long Short Term Memory (LSTM) [39] layer is applied in the next step. The LSTM reads each weighted tweet word at a time and produces an output when the end of the tweet is reached (last word). The output depends on the last word but also on what the LSTM has read till the last word. This history of information starts with the first word and is built up with every new input. However, instead of letting every new word contribute to the historical knowledge LSTM unit applies, similar to humans, some kind of information focus. It uses only those words in the historical knowledge which are also worth to consider. Words that are not worth to consider are not used in construction of the history knowledge. Once the tweet ends LSTM constructs the output. In our model this output is passed to a dense layer (feed forward neural network) which is followed by softmax. Note the LSTM units also used in the baseline systems work on the same principle as described here.

### 4.3. Real time usage

As noted earlier, in our method we tackle the rumour verification with the source tweet alone, i.e. the tweet that is the origin of the rumour. This means that reply chains with respect to the source tweet are not required. This makes our approach applicable for any source tweet with or without replies. Furthermore, since we do not use replying tweets, we also do not require knowledge such as stance information from them. This makes our approach independent from any additional pipelines, which might introduce further noise for the subsequent parts. Due to these facts – focus only on the source tweet and no use of additional pipelines – we argue that our approach can be applied in real time situations and on any short message reporting a rumour.

Note, what we refer with real time usage is that since our approach works only on the source tweet that it can be applied once such a source tweet about an event such as the German Aircraft crash is posted. The source would be carrying a rumour that requires verification. In this current study our aim is to investigate the contribution of recurring terms for this task. Also as noted in Section 3.2.1 the rumours in both RumourEval datasets are independent of each other. This setting makes our approach suitable for distinguishing between rumours that carry exaggeration and false information and those that aim to report proper pieces of information.

## 5. Experiments

### 5.1. Baselines

#### 5.1.1. NileTMRG\*

The best system on the RumourEval2017 shared task achieved 53.6% accuracy on the test dataset [12]. The winning system relied on a traditional feature engineering approach by proposing a linear SVC with a concatenated feature set consisting of a Bag of Word (BoW) source tweet representation, the existence of URL and hashtags as well as the percentage of stance information

<sup>3</sup> <https://nlp.stanford.edu/projects/glove/>.



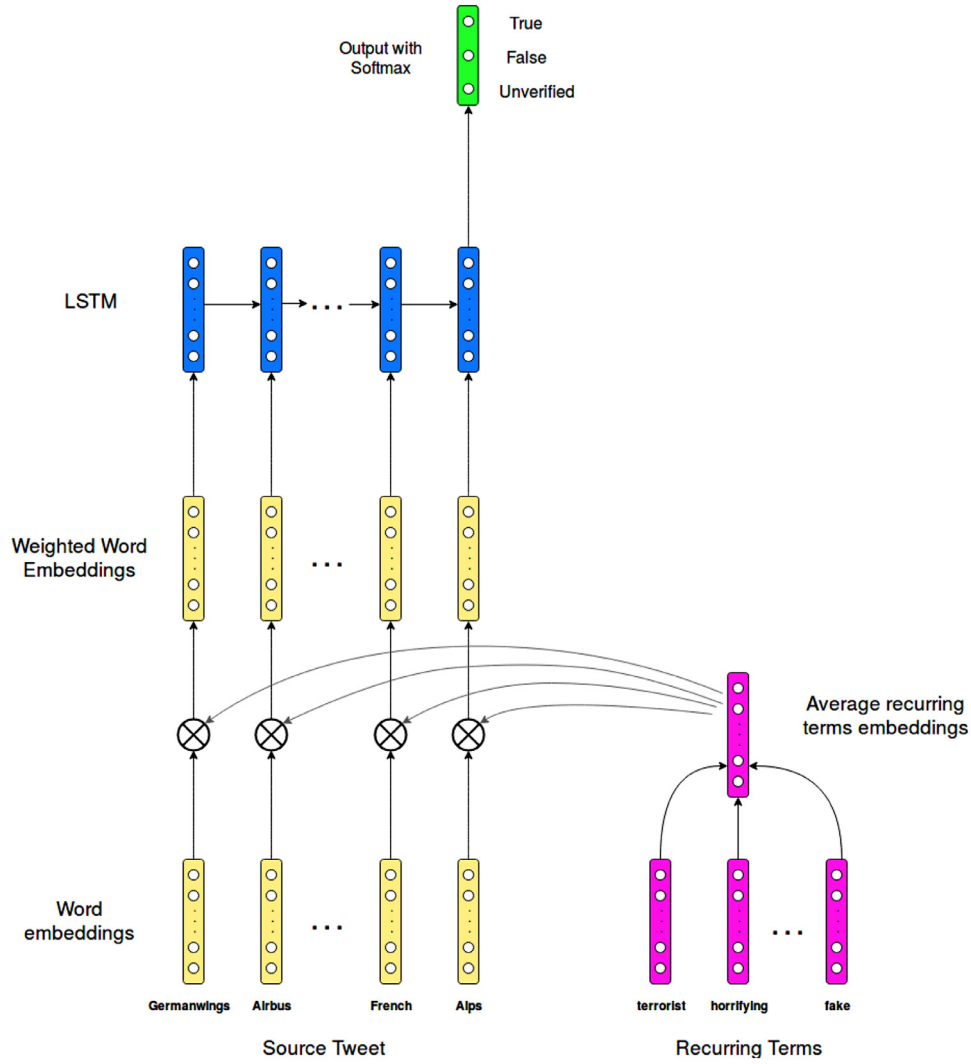


Fig. 2. Inner-attention LSTM network.

from the replying tweets (querying, denying or supporting). A re-implementation of the NileTMRG was performed by Kochkina et al. [13]. The new model NileTMRG\* achieves an accuracy score of 57.1% on the same dataset.

### 5.1.2. Branch LSTM

The branch LSTM (branchLSTM) is a sequential approach which consists of several LSTM layers that are connected to several feed-forward ReLU layers. The prediction layer consists of a softmax layer and outputs the probabilities for the class labels [30]. Kochkina et al. [13] adopted this model for the rumour verification task and applied on the RumourEval2017 dataset.

### 5.1.3. Multi-task learning

Kochkina et al. [13] use a multi-task-learning model in order to predict veracity status of rumours. The model jointly learns to predict for the tasks of rumour detection, stance classification and rumour veracity classification. The joint learning of all the tasks leads to the best performing setting. Using this setting the authors report 57.1% accuracy on the RumourEval2017 test dataset.

### 5.1.4. Vanilla LSTM

In addition to the baselines described above we also use a vanilla LSTM without the inner-attention mechanism. This is to evaluate whether recurring terms applied through the attention

mechanism have an impact on the results or not. Apart from omitting the inner-attention part, the remaining parts of the model shown in Fig. 2 are kept for the vanilla LSTM baseline.

### 5.1.5. LSTM with soft attention

In the last years attention mechanism methods have been applied on the rumour detection task to highlight terms that are helpful to distinguish between rumour/non-rumour representations [27]. To compare our proposed Inner-Attention LSTM we use such an attention mechanism, also referred to as soft attention mechanism [34], as additional baseline. This approach makes use of LSTM working with an attention model. However, unlike our inner attention model, this one does not require recurring term as input to draw attention on the indicative words in the source tweet. LSTM model with soft attention applies attention on top of the LSTM layer in order to determine the important parts of the source tweet after it passes the LSTM layer (see Fig. 3). The outputs of the LSTM layer (at each timestamp) are weighted and aggregated. The weights are learned through the attention mechanism. The result after the aggregation is a context vector  $c$ :

$$c = \sum_{t=1}^T \alpha_t h_t \quad (3)$$

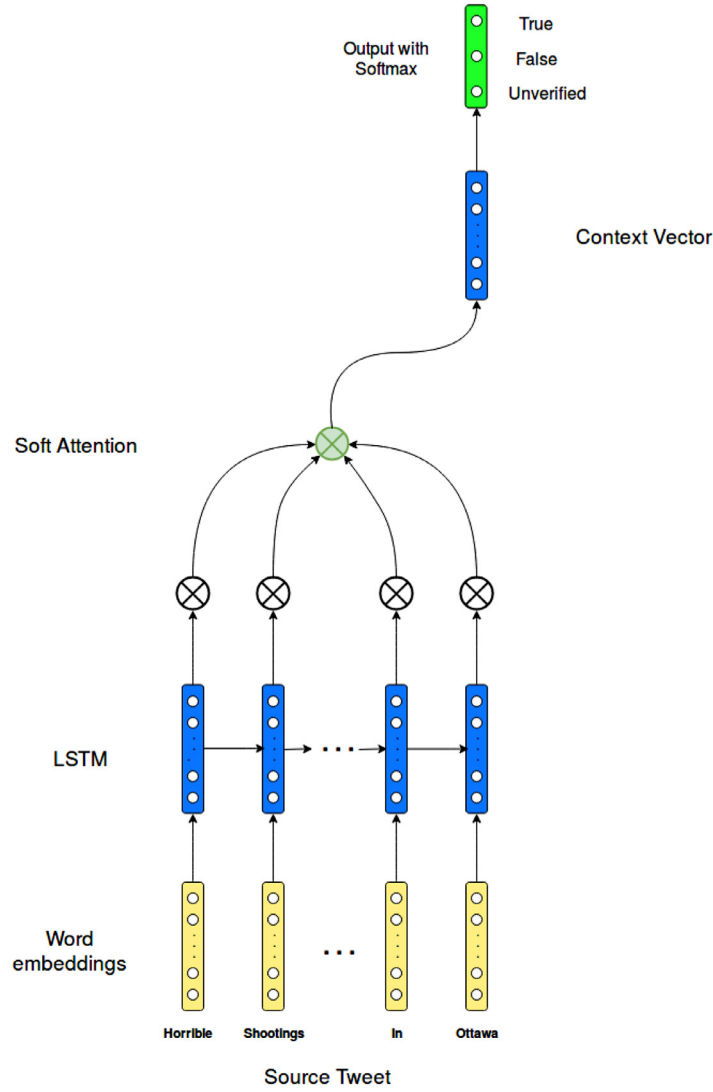


Fig. 3. LSTM network with soft attention.

where  $T$  is the total number of time steps in the hidden state sequence  $h$  and  $\alpha_t$  is an individual weight factor at each time step  $t$ . Each weighting factor  $\alpha_t$  is derived by the following equation.

$$\alpha_t = \sigma(\tanh(W_s h_t)) \quad (4)$$

$\alpha_t$  is determined by a learning function that is composed of its hidden unit  $h_t$  and a trainable weight vector  $W_s$ . Softmax is applied to transform the weighting factors into a probability distribution.

## 5.2. Experimental set-up

### 5.2.1. Experiments on RumourEval2017 and RumourEval2019

For the rumour veracity classification task we conduct two experiments based on different data set-ups. We also conduct an experiment where we apply existing models on the news headlines. In total we have three experiments. As noted above rumours (source tweets) (Section 3.2.1) in our datasets carry some pieces of information put forward about an event and thus require verification. Each such piece of information is independent from the others and thus the rumours about an event do not have timely ordering. This is also the case for the headlines. Therefore in our experiments we treat each rumour/headline independent from

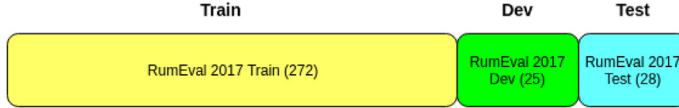
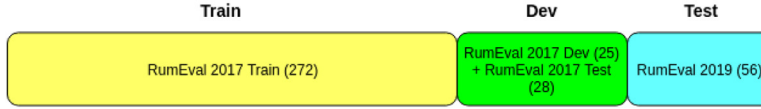
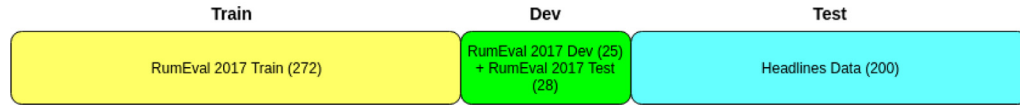
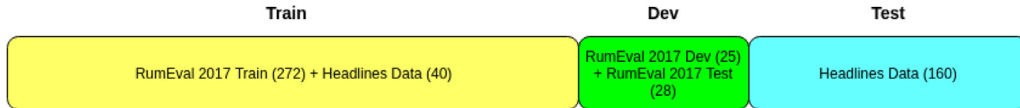
others and assess its veracity only based on the recurring information.

*First experiment:* The first experiment uses the same train/dev/test split as used in the RumourEval2017 Task (cf. Table 1). This means we train our models on the train set, tune the parameters using the dev set and test the tuned model on the test set.

*Second experiment:* The second experiment makes use of the RumourEval2019 dataset. In this second experiment we have two different settings.

*Setting 1:* In the first setting we keep the training from RumourEval2017 dataset as it is, however, merge the dev and test sets together and use this merged dataset for parameter tuning. The entire RumourEval2019 dataset is used as testing data. This tests the performance of a domain independent model (RumourEval2017) on a foreign domain (RumourEval2019).

*Setting 2:* In the second setting, we extend the previously training data with additional RumourEval2019 data (20% – we will release the splits for the community to re-produce our experiments). The merged dev and test sets from RumourEval2017 are still used for tuning the model parameters. This tests the contribution of in-domain data on out-domain trained models. The test happens then on the remaining 80% of the RumourEval2019 data.

**Experiment 1****Experiment 2 (Setting 1)****Experiment 2 (Setting 2)****Experiment 3 (Setting 1)****Experiment 3 (Setting 2)****Fig. 4.** Train/dev/test distribution sets of all experiments and settings.

*Third experiment:* In this experiment, we replicate the settings as done in experiment 2 but replace the testing data with the news headlines. This means we have again two settings:

*Setting 1:* Same as in experiment 2 with setting 1 but the testing data is the entire news headline dataset. The RumourEval2019 dataset is not used.

*Setting 2:* Again same as in experiment 2 with setting 2 but the injected additional data comes from the headline dataset. Also the testing is purely performed with the news headlines.

In all settings the results are macro-averaged. Fig. 4 shows the distribution of train/dev/test sets of all experiments and settings.

### 5.3. Hyperparameters

For the baselines reported by related work (NileTMRG\*, branch LSTM and multi-task learning approaches) we use the models as reported. In the following we describe the hyperparameter tuning performed on the three LSTMs implemented in this work.

#### 5.3.1. Inner-attention LSTM

Hyperparameter tuning is used for finding the best model configuration. We trained several models based on different LSTM sizes, dropout rates, batch sizes, learning rates and epochs. Our best model is trained for 11 epochs with ADAM optimization algorithm, an LSTM size of 10, a learning rate of 0.01, a dropout rate of 0.1 and a batch size of 16.

We also investigated both GloVe and Word2Vec embeddings. The 25 dimensional GloVe word embeddings performed better than Word2Vec.

#### 5.3.2. Vanilla LSTM

The best Vanilla LSTM model is trained on 10 epochs with ADAM optimization algorithm, an LSTM size of 64, a learning rate of 0.01, a dropout rate of 0.1 and a batch size of 16.

#### 5.3.3. Soft attention LSTM

Our second attention-based model gains best performance results by training for 11 epochs with ADAM optimization algorithm, an LSTM size of 16, a learning rate of 0.01, a dropout rate of 0.1 and a batch size of 16. In terms of hyperparameters this model does not differ much from our inner-attention attention-based model.

### 5.4. Usage of recurrent terms

Our recurrent terms are extracted based on two conditions: (1) We use distinct terms for each class, i.e. we make sure terms occurring frequently in one veracity class never occur in the other veracity classes. (2) A term must appear in at least 2 events. The second condition aims to generalize the term extraction by avoiding terms from one event dominating the entire set of recurring terms which are provided as input to the LSTM model.

We tried several strategies to use the recurring terms in our inner-attention LSTM model. In our first strategy, we extracted

**Table 5**

Results from experiment 1. Comparison of performance of several baselines with inner-attention LSTM model tested on RumourEval2017 dataset.

	Majority vote	Nile TMRG*	Multi task learning	Branch LSTM	Vanilla LSTM	Soft attention LSTM	Inner-attention LSTM
Macro <i>F</i> -1	0.429	0.539	0.558	0.491	0.528	0.496	<b>0.616</b>
Accuracy	0.2	0.571	0.571	0.5	0.537	0.5	<b>0.607</b>

**Table 6**Results from experiment 2, setting 1. Comparison of models implemented in this work and strong baseline from related work (according to *F*-1).

	Multi task learning	Vanilla LSTM	Soft attention LSTM	Inner-attention LSTM
Macro <i>F</i> -1	0.285	0.333	<b>0.413</b>	0.399
Accuracy	0.304	0.482	0.464	<b>0.554</b>

**Table 7**Results from experiment 2, setting 2. Comparison of models implemented in this work and strong baseline from related work (according to *F*-1).

	Multi task learning	Vanilla LSTM	Soft attention LSTM	Inner-attention LSTM
Macro <i>F</i> -1	0.477	0.396	<b>0.633</b>	0.478
Accuracy	0.489	0.636	0.636	<b>0.705</b>

from each rumour class up to 20 frequent terms and used them all in our model. In the second strategy, we included only terms from the *true* and *unverified* classes and omitted terms from the *false* class. Finally we used terms only from the *false* class. Note that the last two settings aim to distinguish the *false* class from the other two classes. We achieve the best results when we guide our LSTM model with the *false* terms only. The results reported in this work are all based on recurring terms extracted from the *false* rumours. These terms are extracted from the RumourEval2017 dataset.<sup>4</sup>

## 6. Results

Table 5 shows the evaluation results of our proposed method compared to the baseline approaches. We first report results on the RumourEval2017 dataset (experiment 1).

First we see that all systems outperform the majority vote baseline which simply outputs the class that has the majority class in the training data. The baselines NileTMRG\* and multi-task learning achieve the same accuracy score of 57.1% and differ in *F*-1 scores (53.9% vs. 55.8% *F*-1). Both systems are currently known state-of-the-art systems on the RumourEval2017 dataset. The branchLSTM performs slightly worse and achieves only 50% accuracy and 49.1% *F*-1. A similar performance is achieved with the soft attention model leading to 50% accuracy and 49.6% *F*-1. The Vanilla LSTM model performs better than the branchLSTM and the soft attention models and achieves 53.7% accuracy and 52.8% *F*-1. On the other hand, the proposed inner-attention based LSTM model outperforms the state-of-the-art models by 3% in the accuracy metric and almost by 6% in *F*-1.<sup>5</sup> This shows that injecting the model with recurrent words indeed helps the model to do better distinction between the classes.

The results of the second experiment (experiment 2), setting 1, are shown in Table 6. Note that in this experiment, we evaluate the performance of an existing model on out-of-domain data. For this experiment, we use the models implemented in this work – the proposed inner-attention model, the re-implementation of the soft attention model and finally the vanilla LSTM model – and compare it with the best performing baseline from related work (according to *F*-1), which is the multi-task learning model.

From the results, we see that there is a performance drop by all the systems due to the out-of-domain application. The biggest drop happens for the multi-task learning system and the least drop is seen for the soft-attention LSTM approach. In case of

performance the proposed inner-attention approach is superior to all other models in terms of accuracy. It is able to classify the instances with 55.4% accuracy. In terms of *F*-1 although it has a huge drop on this metric it is only outperformed by the soft attention approach and achieves better results compared to the other two baselines.

We also tested the scenario of recover by injecting some additional data from the RumourEval2019 testing data to the training data. This is the experiment 2 with setting 2. The results for this experiment are shown in Table 7.

From the results in Table 7, we can see that all systems recover and get a boost in performance after seeing 20% of the new data during training. The trend in terms of performance ranking continues among the systems. In terms of accuracy, the inner-attention approach is superior to all other systems. However, according to *F*-1, the soft attention method shines and outperforms all other systems by a large margin.

## 7. Discussion

For results shown in Table 5: Based on the experiment 1 we noted that our inner-attention approach performs substantially better than the current reported state-of-the-art systems such as NileTMRG\* and the multi-task learning approaches. This shows that the recurrent terms indeed help to guide the model in terms of better classification. In addition, it should be noted that these systems use additional information such as stance from the replying tweets. Thus they are forced to wait until some number of replies happen for the source tweet. Our system does not have this limitation and thus have also another advantage that it can be used whenever there is a source tweet. For instance, in the absence of stance information, the multi-task learning system's performance drops to 37% in *F*-1 and 41% in accuracy [13]. Also systems similar to NileTMRG\* like e.g. [23] that focus on feature engineering – but unlike NileTMRG\* omit crowd wisdom such as stance information – achieve an accuracy of only 46.4% on the RumourEval2017 shared task. The proposed approach outperforms also the vanilla LSTM and the soft attention based LSTM approaches with great differences in the scores.

As described in Section 5.4, the inner-attention model uses terms that are selected by conditioning e.g. that they occur in one class (*false*) but do not appear in the other class (*true*). So there is in advance a clear distinction in the e.g. *false*-vs-*true* rumour terms and the results show that exactly this clear distinction helps the model to perform best. The other systems try to learn this distinction on their own, but given the small training data size the distinction stays blurry and they perform poorly.

However, we also injected stance information in our proposed inner attention model to investigate any positive impact the stance would have in the task. The resulting architecture is shown in

<sup>4</sup> Note terms can also be prepared manually and injected to the system. This is especially useful if the manual terms present a better distinction between certain classes of a given problem.

<sup>5</sup> These are results of the model after parameter tuning. We also experimented with models, this is also the case for all subsequent experiments, where no parameter tuning was performed. Results of this experiment are shown in Appendix A.



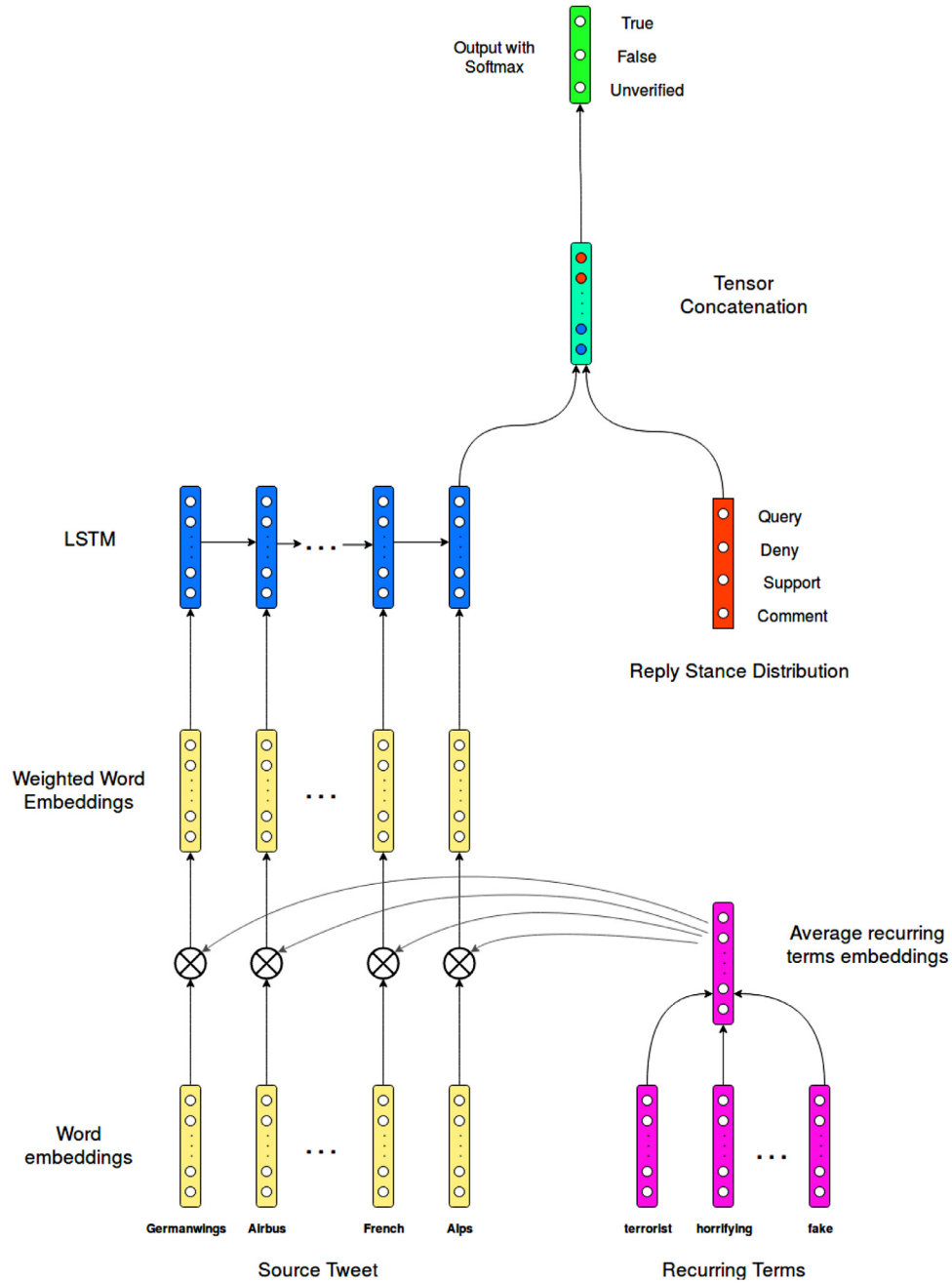


Fig. 5. Inner-attention LSTM network with replies stance distribution.

Fig. 5. In this model we simply aggregate the stance information from the replying tweets (e.g. total number of supporting stances) and merge them with the output of the LSTM before inputting to the softmax layer. This late fusion of stance to the LSTM output aims to emphasize on the importance of stance features and should help to use the stance information effectively.

The late fusion architecture shown in Fig. 5 leads to 56% in macro  $F_1$  and 57.1 accuracy. This is less than what our original model alone achieves. A reason for this, as also shown by Dungs et al. [15], is the way how stance information is passed to the system. Dungs et al. [15] show that stance orientation changes over time and thus is important to capture this development to better capture the wisdom from the crowd and employ it on the verification problem. Our simple aggregation model does not have this capturing characteristics and simply collapses the different stance information given at different time stamps together and employs

the aggregated information on the verification problem. This leads to confusion of the system and affects the results negatively. In our future studies we plan to investigate alternative architectures that capture also stance information well.

Finally, as discussed in Section 4.1 the *false* rumour recurring terms tend to use an exaggerating and the *true* rumours a reporting language (cf. Table 4). In Table 8 we show an example of a false rumour predicted as such. The table shows where the model focuses most (the higher the value the more important is the term for the model for performing the right class prediction). It is clearly recognizable that the words “Muslim”, “big”, “fat” and “mistake” are highly weighted by the model. Most of them also follow an exaggeration pattern and thus the model decides the false label for this particular example. In Table 9 we also show an opposite example with a true rumour that is also correctly classified by our system. The model does not pay almost no attention on stop words,

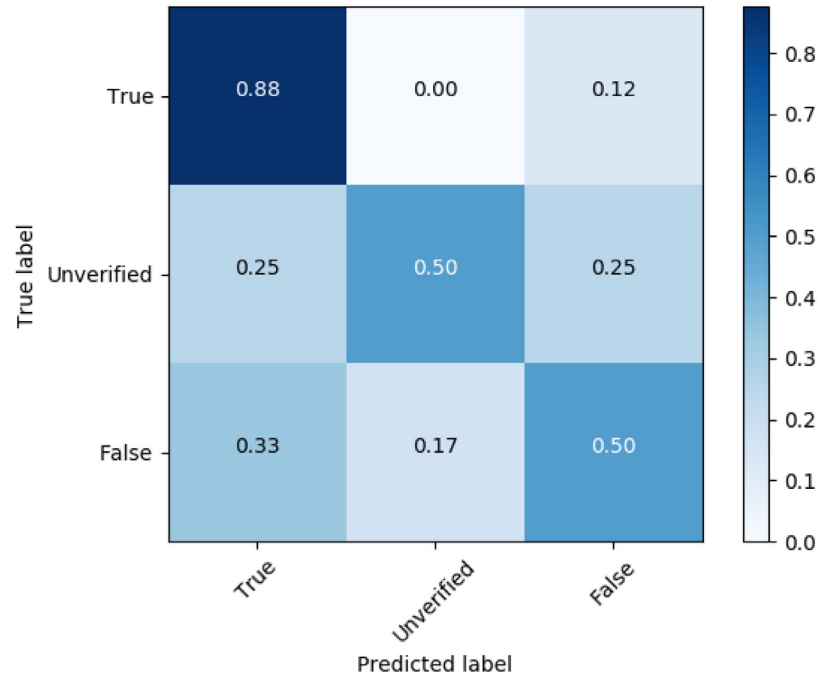


Fig. 6. Confusion matrix of the inner-attention based LSTM model on the RumourEval test dataset.

Table 8

Attention weighting of a correctly classified false rumour post example.

Merkel	ad-	bring-	Mus-	Refugees	was	a	big	fat	mis-
0	mits	ing	lim	0.4	0.2	0.01	0.99	0.99	take
	0.9	0.01	0.99						0.99

Table 10

Attention weighting of a wrongly classified false rumour post example as *unverified*.

Pro-	de-	ejac-	is	mur-	ev-	sperm	cell	is	a	life
Lif-	clare	ulation	0.01	der	ery	0.01	0.01	0.01	0.01	0.99
ers	0.01	0.01		0.99	0.99					
0										

Table 9

Attention weighting of a correctly classified true rumour post example.

New	car	sales	plunge	20	%	in	Septem-
0.81	0.99	0.5	0.2	0.01	0.01	0.01	ber
							0

the number and the month name “September”. It pays attention on the other terms which follow unlike in the false rumour example rather a reporting style. The model makes use of this difference in the style of reporting and classifies this rumour as true.

The *unverified* rumours have vague terms but they tend to look similar to the true rumour terms. This means that the false rumours are lexically different from the other two classes. To understand whether this phenomenon is also reflected in the model behaviour, we computed the confusion matrix (Fig. 6) of the model trained on experimental 1 setting (experiment 1).

From the matrix we can see that the model is able to make the correct prediction for 88% of the time for the true class and confuses them only in 12% of the cases with the false class. This shows that the recurring terms helped to guide the model to perform the class distinction during the classification. It learned that true instances do follow different behaviour than false instances. However, interestingly the false class is confused with the true and the unverified classes. The model miss-classifies 50% of the false messages to the true and unverified categories. We think this is the case when the false rumour starts using the reporting style as

in those two classes and omits indicators of exaggeration. This is e.g. shown in Table 10 with a false rumour that is classified by our model as unverified. The model focuses mostly on the terms like “murder”, “every” and “life” which seem to try exaggerate a phenomenon. However, they are also put into context with reporting style words like “declare”. Because of this the model confuses the false rumour with the unverified class. For real use scenario this means that the proposed system is good in detecting obvious false rumours i.e., rumours that follow no reporting style at all but contain exaggerations, and struggles to capture false rumours following reporting style. Thus in real use scenario we think that the proposed system could work as a pre-filter to detect such obvious false rumours. To also detect false rumours that follow reporting style the system requires to look deeper on other aspects such as user meta information, communication network but also at the replies. We aim to enrich the proposed approach with such aspects in our future work.

Finally, we have also mis-classification in the unverified category (confusion with the true class). As discussed in Section 5.4 the unverified messages are rather vague and we think this is reflected in the results.

For results shown in Tables 6 and 7: In the second experiment (settings 1 and 2) we performed an out-of-domain evaluation. The results showed that overall there is a performance drop in all systems when the model is trained using the RumourEval2017 dataset and tested on the RumourEval2019 testing data. This drops show that all systems are somewhat domain dependent and thus have difficulties to perform decisions on an unseen domain. The most severe case was observed with the multi-task learning approach

**Table 11**

Example recurring terms extracted for each veracity class from RumourEval2019 dataset.

Veracity class	Example terms
False	WTF, wiping, illegal, incredible, never, flooded, infidels, THUGS
True	Weather, islands, hurricanes, guard
Unverified	Nuclear, ocean, drop, helping

and least with the soft-attention LSTM. Overall the best performing systems were the inner-attention approach according to accuracy and the soft attention LSTM according to  $F-1$ . Later in the setting 2 we injected (see Table 7) some out-of-domain data to the training set and could observe recoveries in all systems. The highest recovery was observed in the soft attention LSTM and less in the other models.

These behaviours are connected with how the models are focused on the task. The multi-task learning as well as the vanilla LSTM approaches have no attention mechanisms integrated. They have only the LSTM units to adopt themselves to the task – learn what terms are useful to distinguish between the classes. However, this can only happen with a lot of data which is unfortunately not available. This is why these models do not have much chance to adopt themselves to the task.

The soft attention model has also an LSTM in its core but it has the advantage that it uses an attention part on top of the LSTM which helps the entire model to focus on terms that are indicative for the different classes. It is an extra help for the model to do something more quickly. The recovery results from Tables 6 to 7 show this phenomenon.

In case of the inner-attention, the selection of the recurrent terms is very important. This is the extra bit that helps the model to quickly focus itself to the right classes. However, as described in Section 5.4 our terms are taken from the false rumours and also only from the training data. In experiment 2 with setting 1 all the false rumours come from the RumourEval2017 data and more than 90% in case of setting 2 (4 false rumours from RumourEval2019 have been injected to the training data). This means the recurrent terms still represent the RumourEval2017 domain and how false rumours are expressed in man made disasters. However, unlike the man made disasters the natural catastrophes follow rather a different language usage and have little in common with the false terms. For instance, recurrent terms extracted from the RumourEval2019 entire dataset contain in the false category are terms like *wiping*, *illegal*, *incredible*, *never* (see Table 11). These are certainly different from those extracted from the RumourEval2017 dataset (see Table 4) and thus the model has little information from the attention part to be guided. This is why we think that the inner attention approach struggles with out-of-domain data.

*The weakness is maybe a strength on something else:* In Section 4.1 we discussed that the true/unverified recurrent terms follow a reporting and that the terms extracted from the false message follow an exaggerating style. From the results shown in Table 5, we have seen that our inner attention model has learned to use them well to distinguish between the different classes. A logical question which arises with this is whether such a model is also able to distinguish between news headlines of legitimate news (which follow rather reporting style) and fake news headlines where the language is rather exaggerated. To analyse this we performed the experiment 3 with settings 1 and 2. For these experiments the headlines dataset described in Section 3 is used.

According to experiment 3 setting 1 we first apply the RumourEval2017 trained model on the entire headlines dataset – this means the entire headlines dataset is treated as testset. In the second case (setting 2) we inject 20% from the headlines data to the

**Table 12**

Some models trained on RumourEval2017 dataset and tested on headlines data (Experiment 3, setting 1).

	Vanilla LSTM	Soft attention LSTM	Inner-attention LSTM
Macro $F-1$	0.295	0.323	<b>0.539</b>
Accuracy	0.485	0.545	<b>0.805</b>

**Table 13**

Some models trained on RumourEval2017 dataset with 20% reserved headlines data and tested on remaining headlines (Experiment 3, setting 2).

	Vanilla LSTM	Soft attention LSTM	Inner-attention LSTM
Macro $F-1$	0.49	0.344	<b>0.811</b>
Accuracy	0.734	0.581	<b>0.813</b>

RumourEval2017 training set and test the resulting model on the remaining 80% headlines. We use the vanilla, soft attention and the inner-attention LSTM approaches. Note, these are the only models which do not use additional information such as stance and thus are directly applicable to the headlines use case. The other systems require also stance information from replies which do not exist in the headlines case.

From the results shown in Table 12, we see that the inner-attention LSTM approach outperforms the other two baselines significantly ( $p < 0.001$ ).<sup>6</sup> It clearly learned how to distinguish between reporting style (true messages) and exaggerated language (false messages). It achieves an accuracy of 80.5% and 53.4%  $F-1$ . Furthermore, it learns to adopt itself more on the task after seeing 20% of the headlines (results shown in Table 13). In this case its  $F-1$  performance gets boosted to 81.1% and gains slightly further on accuracy. In the two other models the difference between reporting and exaggeration is not captured and thus they fail to perform well on this second out-of-domain data.

## 8. Conclusion and outlook

We have proposed a rumour veracity system that solely works on the content of source rumour post without requiring reply messages and additional crowd wisdom such as stance information. By leveraging the fact that false rumours are reported with different recurring terms in contrast to true as well as unverified rumours, we use the recurring false terms and inject them in our inner-attention LSTM in order to guide it in the distinction between the rumour classes. We compare the results of the proposed approach to vanilla LSTM and soft attention model where no attention on these (false) indicators is performed along with three baselines which reported state-of-the-art results on known benchmark data. In the RumourEval2017 experiment our model outperforms the state-of-the-art classifiers leading to 60.7% accuracy and 61.6%  $F-1$ . We also showed that simple aggregation of the stance information from the replies and merging it with the inner-attention network does not lead to better performance but rather confuses the system slightly. Additionally we conducted an experiment using the RumourEval2019 as test set and showed that our model is competitive compared to the baselines, achieves even best results in terms of accuracy but fails to take the lead in terms of  $F-1$ . Finally we showed that our proposed model shines in the headline experiments and clearly outperforms baselines with large margin.

In the future we aim to investigate alternative architectures that capture both recurring patterns and stance information. However,

<sup>6</sup> We performed an independent 2-samples  $t$ -test with Bonferroni correction.

we also plan to investigate user information as well as communication network as additional features in our system. Next, we aim to also apply our model to different problems such as hate speech and help seeking/offering detection in disaster situations. In both cases, users reporting in social media will tend to use recurring terms. We think that capturing and using these terms will help to tackle these tasks.

### Declaration of Competing Interest

We confirm that this paper is no submitted elsewhere. There is also no conflict of interest.

### Acknowledgements

This work was partially supported by the European Union under grant agreements No. 654024 SoBigData, No. 825297 WeVerify and the Deutsche Forschungsgemeinschaft (DFG, [German Research Foundation](#)) – [GRK 2167](#), Research Training Group “User-Centred Social Media”.

### Appendix A

Table 14 shows the performance of our inner attention model with and without parameter tuning.

**Table 14**

Performance comparison of inner-attention LSTM model w/ and w/o parameter tuning. The first number refers to macro F1-score while the second number represents accuracy.

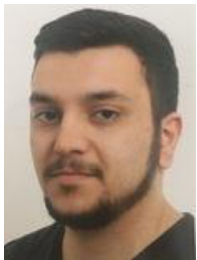
	Inner-attention LSTM	Inner-attention LSTM w/o parameter tuning
Experiment 1	0.607/0.616	0.534/0.536
Experiment 2 (Setting 1)	0.399/0.554	0.373/0.536
Experiment 2 (Setting 2)	0.478/0.705	0.467/0.659
Experiment 3 (Setting 1)	0.539/0.805	0.505/0.76
Experiment 3 (Setting 2)	0.811/0.813	0.778/0.781

### References

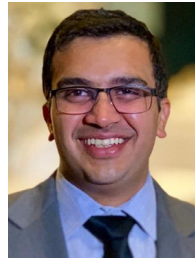
- [1] A. Whiting, D. Williams, Why people use social media: a uses and gratifications approach, *Qual. Market Res. Int. J.* 16 (4) (2013) 362–369.
- [2] A. Zubiaga, A. Aker, K. Bontcheva, M. Liakata, R. Procter, Detection and resolution of rumours in social media: a survey, *ACM Comput. Surv. (CSUR)* 51 (2) (2018) 32.
- [3] M. Hu, S. Liu, F. Wei, Y. Wu, J. Stasko, K.-L. Ma, Breaking news on twitter, in: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, ACM, 2012, pp. 2751–2754.
- [4] V. Qazvinian, E. Rosengren, D.R. Radev, Q. Mei, Rumor has it: identifying misinformation in microblogs, in: *Proceedings of the 2011 EMNLP*, 2011, pp. 1589–1599.
- [5] C. Castillo, M. Mendoza, B. Poblete, Information credibility on twitter, in: *Proceedings of the 20th International Conference on World Wide Web*, ACM, 2011, pp. 675–684.
- [6] S. Kwon, M. Cha, K. Jung, W. Chen, Y. Wang, Prominent features of rumor propagation in online social media, in: *Proceedings of the 13th IEEE International Conference on Data Mining*, IEEE, 2013, pp. 1103–1108.
- [7] S. Vosoughi, Automatic detection and verification of rumors on Twitter. PhD diss., Massachusetts Institute of Technology, 2015.
- [8] K. Wu, S. Yang, K.Q. Zhu, False rumors detection on Sina Weibo by propagation structures, in: *Proceedings of the 31st IEEE International Conference on Data Engineering*, IEEE, 2015, pp. 651–662.
- [9] J. Ma, W. Gao, Z. Wei, Y. Lu, K.-F. Wong, Detect rumors using time series of social context information on microblogging websites, in: *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*, ACM, 2015, pp. 1751–1754.
- [10] M. Lukasik, P. Sriji, D. Vu, K. Bontcheva, A. Zubiaga, T. Cohn, Hawkes processes for continuous time sequence classification: an application to rumour stance classification in twitter, in: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, Association for Computational Linguistics, 2016, pp. 393–398.
- [11] X. Liu, A. Nourbakhsh, Q. Li, R. Fang, S. Shah, Real-time rumor debunking on twitter, in: *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*, ACM, 2015, pp. 1867–1870.
- [12] O. Enayet, S.R. El-Beltagy, NileTMRG at SemEval-2017 Task 8: determining rumour and veracity support for rumours on Twitter, in: *Proceedings of the 2017 SemEval*, ACL, 2017.
- [13] E. Kochkina, M. Liakata, A. Zubiaga, All-in-one: multi-task learning for rumour verification, in: *Proceedings of the 27th International Conference on Computational Linguistics*, 2018, pp. 3360–3370.
- [14] A. Sen, K. Rudra, S. Ghosh, Extracting situational awareness from microblogs during disaster events, in: *Proceedings of the 7th International Conference on Communication Systems and Networks (COMSNETS)*, IEEE, 2015, pp. 1–6.
- [15] S. Dungs, A. Aker, N. Fuhr, K. Bontcheva, Can rumour stance alone predict veracity? in: *Proceedings of the 27th International Conference on Computational Linguistics*, 2018, pp. 3360–3370.
- [16] S. Kwon, M. Cha, K. Jung, Rumor detection over varying time windows, *PLoS One* 12 (1) (2017) e0168344.
- [17] F. Yang, Y. Liu, X. Yu, M. Yang, Automatic detection of rumor on Sina Weibo, in: *Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics*, ACM, 2012, p. 13.
- [18] S. Wang, T. Terano, Detecting rumor patterns in streaming social media, in: *Proceedings of the 2015 IEEE International Conference on Big Data (Big Data)*, IEEE, 2015, pp. 2709–2715.
- [19] W. Chen, C.K. Yeo, C.T. Lau, B.S. Lee, Behavior deviation: an anomaly detection view of rumor preemption, in: *Proceedings of the 7th IEEE Annual on Information Technology, Electronics and Mobile Communication Conference (IEMCON)*, IEEE, 2016, pp. 1–7.
- [20] C. Chang, Y. Zhang, C. Szabo, Q.Z. Sheng, Extreme user and political rumor detection on twitter, in: *Proceedings of the 12th International Conference on Advanced Data Mining and Applications, ADMA 2016*, Springer, Gold Coast, QLD, Australia, 2016, pp. 751–763. December 12–15, 2016.
- [21] G. Tong, W. Wei, G. Ling, L. Deying, L. Cong, L. Bin, D. Ding-Zhu, An efficient randomized algorithm for rumor blocking in online social networks, in: *IEEE Transactions on Network Science and Engineering*, 2017.
- [22] L. Derczynski, K. Bontcheva, M. Liakata, R. Procter, G. Wong Sak Hoi, A. Zubiaga, SemEval-2017 Task 8: RumourEval: determining rumour veracity and support for rumours, in: *Proceedings of the 2017 SemEval*, ACL, 2017.
- [23] F. Wang, M. Lan, Y. Wu, ECNU at SemEval-2017 task 8: rumour evaluation using effective features and supervised ensemble models, in: *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, 2017, pp. 491–496.
- [24] V. Singh, S. Narayan, M.S. Akhtar, A. Ekbal, P. Bhattacharya, IITP at SemEval-2017 Task 8: a supervised approach for rumour evaluation, in: *Proceedings of the 2017 SemEval*, 2017.
- [25] Y.-C. Chen, Z.-Y. Liu, H.-Y. Kao, IKM at SemEval-2017 Task 8: convolutional neural networks for stance detection and rumor verification, in: *Proceedings of the 2017 SemEval*, ACL, 2017.
- [26] A. Srivastava, R. Rehm, J. Moreno Schneider, DFKI-DKT at SemEval-2017 Task 8: rumour detection and classification using cascading heuristics, in: *Proceedings of the 2017 SemEval*, ACL, 2017, pp. 486–490.
- [27] T. Chen, X. Li, H. Yin, J. Zhang, Call attention to rumors: deep attention based recurrent neural networks for early rumor detection, in: *Proceedings of the 2018 Pacific-Asia Conference on Knowledge Discovery and Data Mining*, Springer, 2018, pp. 40–52.
- [28] Z. Zhao, P. Resnick, Q. Mei, Enquiring minds: early detection of rumors in social media from enquiry posts, in: *Proceedings of the 24th International Conference on World Wide Web*, International World Wide Web Conferences Steering Committee, 2015, pp. 1395–1405.
- [29] G. Genevieve, K. Bontcheva, L. Derczynski, E. Kochkina, M. Liakata, and A. Zubiaga, RumourEval 2019: Determining Rumour Veracity and Support for Rumours. arXiv preprint arXiv:1809.06683 (2018).
- [30] A. Zubiaga, E. Kochkina, M. Liakata, R. Procter, M. Lukasik, K. Bontcheva, T. Cohn, I. Augenstein, Discourse-aware rumour stance classification in social media using sequential classifiers, *Inf. Process. Manag.* 54 (2) (2018) 273–290.
- [31] K. Shu, D. Mahudeswaran, S. Wang, D. Lee, H. Liu, Fakenewsnet: A data repository with news content, social context and dynamic information for studying fake news on social media. arXiv preprint arXiv:1809.01286 (2018).
- [32] J. Gerhards, M.S. Schäfer, Is the internet a better public sphere? Comparing old and new media in the usa and germany, *New Med. Soc.* 12 (1) (2010) 143–160.
- [33] K. Starbird, Examining the alternative media ecosystem through the production of alternative narratives of mass shooting events on twitter, in: *Proceedings of the Eleventh International AAAI Conference on Web and Social Media*, 2017.
- [34] D. Bahdanau, K. Cho, Y. Bengio, Neural Machine Translation by Jointly Learning to Align and Translate, arXiv:1409.0473 (2014).
- [35] B. Wang, K. Liu, J. Zhao, Inner attention based recurrent neural networks for answer selection, in: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1, 2016, pp. 1288–1297.
- [36] S. Christian, I. Gurevych, Cross-topic argument mining from heterogeneous sources using attention-based neural networks. arXiv preprint arXiv:1802.05758 (2018).
- [37] T. Mikolov, I. Sutskever, K. Chen, G.S. Corrado, J. Dean, Distributed representations of words and phrases and their compositionality, in: *Proceedings of the Advances in Neural Information Processing Systems*, 2013, pp. 3111–3119.
- [38] J. Pennington, R. Socher, C.D. Manning, Glove: global vectors for word representation, in: *Proceedings of the 2014 EMNLP*, 14, 2014, pp. 1532–1543.
- [39] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural Comput.* 9 (8) (1997) 1735–1780.



**Dr. Ahmer Aker** is a post-doctoral research fellow at the University of Duisburg-Essen, Germany. He has organized one workshop at LREC 2014 on Building and Using Comparable Corpora. His main research area is multi-lingual argument mining and argument retrieval from social media. The retrieval of reliable arguments is an important factor and thus Dr Aker is also concerned with fake news as well as mis/disinformation spread in social media. He has published several papers on these topics as well as developed software solutions for tracking mis/disinformation in social media. His current stance classification solution achieved the state-of-the-art results on the SemEval stance task. Dr Aker has published papers at ACL, EMNLP, CIKM, ECIR as well as journals such as JASIST, LRE. He has more than 50 peer-reviewed publications. He is also acting as PC member for several topic relevant workshops and international conferences.



**Alfred Sliwa** is PhD student at the University of Duisburg-Essen, Germany. His research interests are interactive information retrieval, machine learning and information diffusion in social networks. He obtained a master degree in Computer Science at the University of Duisburg-Essen. The topic of his master thesis was about fake news detection and rumour verification. He published conference and workshop papers concerning argument mining at LREC and EMNLP.



**Fahim Dalvi** is an experienced Software Engineer with a demonstrated history of working in the research industry and is currently employed at the Qatar Computing Research Institute. Fahim's research is centered around the intersection of Natural Language Processing and Deep Learning, and he has worked on wide variety of problems in these fields including Machine Translation, Language Modelling and Explainability in Deep Neural Networks. He also spends his time converting research into practical applications, with a focus on scalable web applications. Fahim also spends some time every year mentoring and teaching Deep Learning at Fall and Summer schools.



**Prof. Kalina Bontcheva** conceived and led the PHEME project (<http://pheme.eu>), which (in January 2014) was amongst the first to study computational methods for detecting and tracking of online misinformation. Following on from PHEME, Kalina is now the scientific coordinator of the WeVerify project, which is developing intelligent content verification tools for journalists and citizens, alongside new AI methods for detecting online misinformation and analysing its spread and impact.

Prof. Bontcheva and her team have also carried out real-time analysis of political engagement on social media in the run up to the 2015 and 2017 UK general elections, as well as three joint projects with BuzzFeed UK: analysing political discourse on Twitter during the UK EU membership referendum; online abuse towards politicians; and the impact of Russian cyborg accounts. In December 2017 she was called as an expert witness by the UK Digital, Culture, Media and Sport Committee enquiry in fake news. Between September 2018 and March 2019 Kalina co-authored a STOA report for the European Parliament on the challenges of technological responses to online disinformation.