# Rumor Detection in Tweets Using Graph Convolutional Networks

Takumi Takei, Yuichi Sei [a], Yasuyuki Tahara [b] and Akihiko Ohsuga [c]
*Department of Informatics, The University of Electro-Communications, Chofu, Tokyo, Japan*

Keywords:     Rumor Detection, Deep Learning, Natural Language Processing, Twitter.

Abstract:     The recent development of social networking services has made it easier for anyone to get information. On the other hand, rumors which are information whose truth is unverified are not only easy to spread but also can cause damage such as flames, incitement, and slander. Accurate identification of rumors is effective against such problems and may prevent the spread of misinformation. Based on previous research, this study created a dataset of rumors including replies to fact-checked Japanese tweets. Using a GCN-based deep learning classifier, we performed binary classification of whether a tweet is a False rumor or not, and multinomial classification of True rumor, False rumor, and Unclear rumor, varying the amount of propagation information used. The result of binary classification shows that the maximum accuracy is 0.637, and the maximum F value is 0.641, while the result of multinomial classification shows that the maximum accuracy is 0.547, and the maximum F value is 0.460. We discussed the effectiveness of propagation information and deep learning for detecting Japanese rumors.

## 1 INTRODUCTION

In recent years, with the development of the Internet, social networking services (SNS) have been widely developed and many people can easily send and receive information using SNS, but there are also many posts on SNS that are not true or whose truth is unknown. Such information is called rumors, and they can have negative effects such as the spread of misinformation, inflammation, and incitement. As an example, it has been reported that at least 800 people have died due to a rumor that COVID-19 disease can be cured by consuming large amounts of alcohol, and that rumor has resulted in the loss of life (Islam, M. S., et al. 2020).

In such a situation, fact-checking information is being conducted worldwide. English information is mainly fact-checked by Snopes.com[1], and Japanese information is mainly fact-checked by FactCheck Initiative Japan (FIJ)[2]. However, the fact-checking process is mainly conducted manually, which is a time-consuming and labor-intensive task.

To address such a problem, it is assumed that the ability to accurately automatically classify the kind of rumor or identify an obvious False rumor in a post on a social networking service is not only an effective solution to this problem, but also can significantly affect the behavior of the poster before and after the post, and contribute to reducing the number of False rumors on social networking services.

For the purpose of research on the classification of rumors identified as clearly false tweets on Twitter, we target Japanese posts which haven't been studied and conducted a binary classification of whether a rumor is a False rumor or not. Furthermore, considering the practical sides of the problem, we executed a multinomial classification of whether the rumor is True, False, or Unclear. Also, we extract features used in discrimination, not only related to users and tweets but tweet texts were converted into embedded expressions using BERT which has not been used in recent related studies. In addition, propagation information such as replies to discriminated tweets was also incorporated as a feature, and we also checked and examined whether it is also effective for Japanese by differentiating the amount of information used for discrimination.

For our approach, first, we created a Japanese rumor dataset by collecting fact-checked Japanese

[a] https://orcid.org/0000-0002-2552-6717
[b] https://orcid.org/0000-0002-1939-4455
[c] https://orcid.org/0000-0001-6717-7028

[*] https//www.snopes.com/fact-check/
[*] https://fij.info/

postings, graphing the propagation information, and extracting features. Next, we used two types of BERT models to acquire embedded expressions and created three datasets with different amounts of propagated information. Finally, we carried out the classifications using a deep learning classifier based on Graph Convolutional Networks (GCN) for each of the six datasets.

The results showed that the dataset with the most propagation information produced the Accuracy and F value for both classifications. These results indicate that propagation information is effective to some extent in estimating rumors in Japanese. Although the experiments in this paper were conducted on Japanese, we also consider that our classification method, which is based on the acquisition of embedded expressions by BERT and classification by graph convolution using GCN, can be applied not only to Japanese but also to other languages including English. With this in mind, we discussed and examined the classification of rumors for Japanese rumors.

This paper is organized as follows. Chapter 1 describes the introduction of the study such as background and purpose, Chapter 2 discusses related works. We propose the methodology and dataset in Chapter 3, Chapter 4 describes the experiments and results. Finally, Chapters 5 and 6 discuss the results and future prospects.

## 2 RELATED WORKS

In this chapter, we mention the rumor detection studies, which are previous studies of our research. And BERT, a natural language processing model, is used to acquire embedding representations in sentences.

### 2.1 Rumor Detection

In recent years, many Rumor Detection studies that distinguish fact-checked tweets using the information related to user and tweets, and so on have been reported. The research on which the recent research was based has made a dataset including not only user information and tweet information but also information on the propagation of posts on Twitter and Weibo that has been fact-checked on the fact-checking site Snopes.com, and has performed two experiences. One is a binary classification that discriminates between "Non-rumor" or not. Another is the four-value classification that discriminates between "Non-rumor", "True rumor", "False rumor",

and "Unverified rumor" using a machine learning method (Liu, X., et al. 2015, Ma, J., et al. 2016). It has been a benchmark dataset and task for rumor detection research in recent years, and rumor detection research has been conducted in English and Chinese. There are researches using the tree-structured RNN method (Ma, J., et al. 2018), and established detection models such as combining two Graph Convolutional Networks for the above dataset. (Bian, T., et al. 2020). Furthermore, it has reported an approach that points out the issues of the GCN method, it establishes a detection model that learns time series information, propagation information, and textual information separately and combines them (Li, J., et al. 2021), and a detection model that divides linear and nonlinear information, focusing on the fact that each feature of information has a different structure (Lao, A., et al. 2021). Also, there are several approaches based on Generative Adversarial Networks (GAN). It has proposed rumor detection models that can detect rumors without the need for verified data (Cheng, M., et al. 2021), and that generate opinions in favor of or against replies and discriminate based on the results of such generation (Ma, J., et al. 2021). In addition, not only in English and Chinese which could be estimable with the existing dataset, a new rumor detection for Arabic tweets is also conducted (Alzanin, S.M., et al. 2019).

While many studies on Rumor Detection have been reported, most of the studies on Rumor Detection for Japanese have focused on the search for rumors, and there has been no study on Rumor Detection from the viewpoint of how much discrimination can be performed on fact-checked data, which has been the subject of many studies in recent years. Therefore, in this paper, we conducted research on fact-checked posts in Japanese by modifying the extracted features and the discriminant model from the existing studies.

### 2.2 BERT

BERT is a bidirectional transformer-based natural language processing model that requires two types of learning pre-training and fine-tuning, and by adding an output layer, achieved the best results for various language processing tasks (Devlin, J., et al. 2018). In addition, it is also known to be able to acquire an embedding representation of a sentence from the layer before the final layer of output. While the pre-training step takes a great amount of time, pre-trained BERT models using large corpora in each language have been published, and pre-trained models for Japanese have been published as well. (Shibata, T.,

et.al 2019) This study used two Japanese pre-trained BERT models to get sentence embeddings and performed two classifications.

In summary, the focus points of this study are the following two points.

・ A Japanese rumor dataset was created based on the benchmark dataset, and GCN-based the classifier was set up and classified by features using propagation information.

・ By using BERT's sentence embedding method, the content and semantic aspects of tweets are taken into account in the embedding of tweets, rather than the embedding conversion method using feature words like Bag-of-Words.

# 3 APPROACH

## 3.1 Problem Settings

Table 1: Definition of each Category in FIJ.

| | | |
|---|---|---|
| True | Accurate | The claim is factually accurate and not lacking significant elements. |
| | Mostly Accurate | The main elements of the statement are factually accurate, but there are some minor or insignificant errors. |
| Unclear | Misleading | The claim appears not to have a factually inaccurate element, but has a high possibility of causing misunderstanding due to a click-bait style, lack of important facts, or such. |
| | Inaccurate | The claim lacks overall accuracy, but is a mixture of accurate and inaccurate elements. |
| | Unfounded | The claim is not proven to be factually false, but there is very little or no evidence to support the claim. |
| False | False | All or core elements of the claim are factually inaccurate. |
| | Fake | All or core elements of the claim are factually inaccurate. The speaker or writer is strongly suspected of knowing they are inaccurate. |

In this study, tweets whose fact-checking results were shown in the FIJ were treated as rumor data for the study. The definition of each detection category is shown in Table 1, where the seven labels with discrimination results are further grouped into three groups: True, False, and Unclear. While the research

objective of this paper is to determine obvious False rumor or not, considering the practical aspect, classification against only two groups is not appropriate as a problem sets up because there are three groups of tweets.

So, we attempted to evaluate Rumor Detection in Japanese through 1) Binary classification to determine whether a tweet is judged as a False rumor or not, and 2) Multinominal classification, which discriminates between True rumor, False rumor, and Unclear rumor.

## 3.2 Proposed Method

The proposed method outline is in the following Figure 1 and (1)~(3). (1) Posts including replies were converted to graphical data. (2) The feature extraction described in chapter 4 is performed from the modified graph data, and the data is converted into graph data with features. (3) Binary and multinomial classification results were calculated by convolving the entire graph data with the feature values and assigning them to a deep-learning classifier.
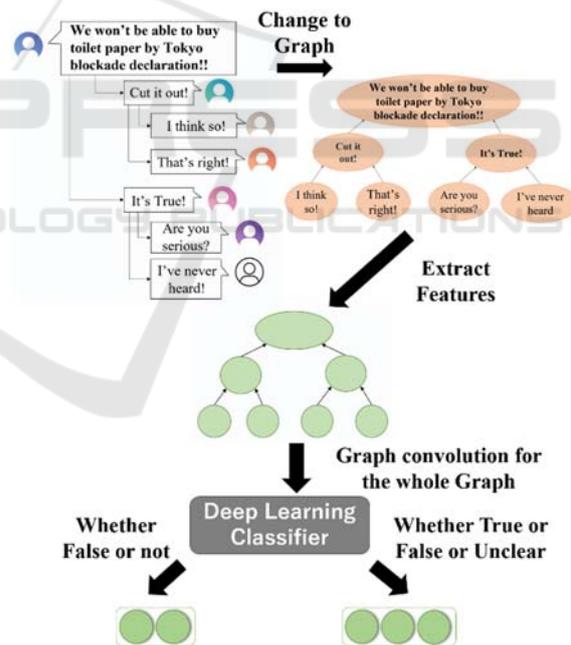


Figure 1: Approach Overview.

### 3.2.1 Change to Graph Data

On Twitter, information is exchanged and spread through replies and retweets to posts. In other words, a tweet starts a topic, and information is propagated through a series of replies to it. This structure is very similar to a tree structure, a type of graph structure,

and can be viewed as a tree structure. Based on these characteristics, we transformed the posted data into a graph as shown in Figure 2, and extracted features to create a graph with features for rumor detection.
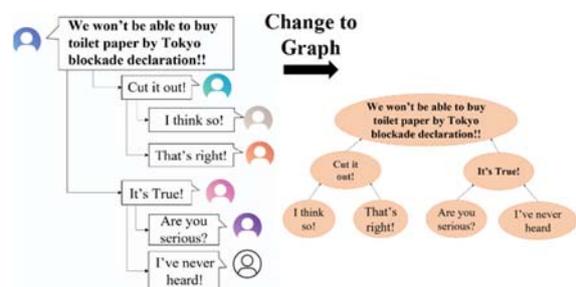


Figure 2: Sample post and replies change to Graph Style *(Post and replies translated into English).*

### 3.2.2 Feature Extraction

We collected the tweet data which were checked in FIJ using Twitter Search API v2. From the collected data, we excluded quoted tweets and extracted the feature. Before extracting features, we performed twelve different cleaning processes on the tweet texts as shown from (1) to (12).

(1) Change alphanumeric characters to half-width characters

(2) Delete the sentence interposed Japanese-style quotation marks, square brackets, lenticular brackets

(3) Count and delete URLs

(4) Count and delete @...

(5) Delete quoted sentence

(6) Delete consecutive characters

(7) Change alphabetic characters to lower

(8) Count and delete hashtags

(9) Delete Emoji, special characters, and date

(10) Delete garbling sentence

(11) Delete space and line break

(12) Change any number to 0

We extract two kinds of features "Basis Features" and "Text Features". In the following, we explain what two features are and how to extract them.

Basis Features are the features obtained from the part other than tweet texts. Specifically, 22 types of information about users and tweets as shown in Table 2 below were calculated from collected data and treated in this study.

Table 2: Basis Features.

| Features | Meaning |
| --- | --- |
| Different Days between Twitter Start and Making Account | The difference in days Twitter Start and Making Account |
| Number of Followers | The number of followers that the user has |
| Number of Follows | The number of follows that the user has |
| Number of Lists | The number of lists that the user sets |
| Number of Tweets | The number of tweets that the user posts |
| Followers / Follows | The ratio of followers to follows |
| Follows / Followers | The ratio of follows to followers |
| Has_Private | Whether the account has a private setting or not |
| Has_Official | Whether the account is an official account or not |
| Has_Abstract | Whether the account sets abstract or not |
| Has_UserLocation | Whether the account has a location setting or not |
| Has_Pinned_User | Whether the account has a pinned user ID or not |
| Language | The language in which the account is set up (Japanese, English, and others) |
| Different Days between Twitter Start and Making Tweet | The difference in days Twitter Start and Making Tweet |
| Number of URLs | The number of URLs in a tweet |
| Number of Hashtags | The number of Hashtags in a tweet |
| Kind_of_Tweet | The type of tweet (None, replied_to, quoted) |
| Has_Link | Whether the tweet has a link or not |
| Has_Media | Whether the tweet has media or not |
| Kind_of_Media | Type of media file with the tweet (None, photo, others) |
| Source of Tweet | Which media is the tweet posted from (Android, iPhone, WebApp, others) |
| Target of Reply | Which target person of the tweet (everyone, others) |

Text Features are the features obtained from the tweet texts which are converted to the vector representation (embedding representation). We tried to use a method

to convert the tweet texts into a vector representation using a pre-trained model of BERT. We used two pre-trained models; one was released from Tohoku University and the other was released from Kyoto University. Both models used ipadic for the dictionary of BERT but used Mecab for morphological analysis in the Tohoku University model, and Juman++ for morphological analysis in the Kyoto University model. Also, we tried to apply the Word2Vec conversion method as a comparison to the methods by BERT.

Hereafter, we will refer to the dataset using Word2Vec as "w2v", the Tohoku University BERT pre-trained model as "Bert-Mecab", and the dataset using the Kyoto University BERT pre-trained model as "Bert-Juman" when the tweets are converted to a vector representation.

### 3.2.3 Classification

In this paper, we use the posted data modified into graph data to perform binary classification of False rumor or not, and multinomial classification of True rumor, False rumor, or Unclear rumor. We set up a graph convolution network-based classification model as shown in Figure 3 below. For both classifications, the layers of the model were not changed, but the number of outputs in the output layer was changed to two for binary classification, and to three for multiclass classification. The label with the highest value among the output values was used as a detection result for comparison and evaluation with the correct data. Comparison and evaluation were made by calculating Accuracy (Acc), Precision (Pre), Recall (Rec), and F-value (F).

In addition, by focusing on the post creation time in the collected data, it is possible to sort the tweets in order of the earlier posting time. In order to consider whether the propagation information is useful for estimating Japanese rumors in the classification process, the dataset was sorted in timestamp order, and we made each dataset 1/6 of the total (Little), 1/2 of the total (Half), and whole of the total (All) extracted, and the classification results were calculated for each of these datasets.
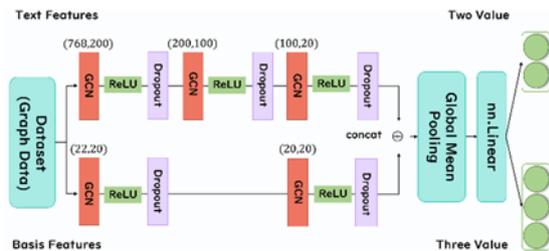


Figure 3: Deep Learning Classifier Model.

## 4 EXPERIMENTS

### 4.1 Dataset Details

The number of events and each tweet in the completed dataset is shown in Table 3 below. The number of source tweets in the table refers to the number of fact-checked tweets, and the number of events refers to the number of topics discussed.

Table 3: Dataset Details.

| | |
|---|---|
| All Tweets | 22,486 |
| Source Tweets | 100 |
| Events | 100 |
| True Source Tweets | 4 |
| Unclear Source Tweets | 46 |
| False Source Tweets | 50 |
| True Tweets (include Replies) | 1,736 |
| Unclear Tweets (include Replies) | 7,761 |
| False Tweets (include Replies) | 12,989 |

### 4.2 Experimental Settings

The complete datasets were sorted in timestamp order, and we made each dataset Little, Half, All. We set up a deep learning classifier using the GCN layer as shown in Figure 3. We performed three times 10-point cross-validation to calculate the average values of Acc, Pre, Rec, and F. Experiments were conducted with CrossEntropyLoss as the loss function, Adam as the optimizer, batch size is 8, Dropout rate is 0.5, and 200 epochs in the deep learning classifier. In summary, we created three datasets with different amounts of reply data for each of the three datasets w2v, Bert-Mecab and Bert-Juman, and compared the results by performing binary and multinominal classification using a setting deep learning classifier.

### 4.3 Results

#### 4.3.1 Binary Classification

The results of the binary classification of False rumor or not are shown in Tables 4, 5, and 6 below. The result of dataset w2v is Table 4, Bert-Mecab is Table 5, and Bert-Juman is Table 6.

Table 4: Binary Classification (w2v).

| | Acc ↑ | Pre ↑ | Rec ↑ | F ↑ |
|---|---|---|---|---|
| Little | 0.527 | 0.507 | 0.512 | 0.505 |
| Half | 0.573 | 0.587 | 0.578 | 0.581 |
| All | 0.597 | 0.611 | 0.616 | 0.608 |

Table 5: Binary Classification (Bert-Mecab).

|  | Acc ↑ | Pre ↑ | Rec ↑ | F ↑ |
|---|---|---|---|---|
| Little | 0.533 | 0.475 | 0.467 | 0.468 |
| Half | 0.613 | 0.635 | 0.619 | 0.626 |
| All | **0.637** | **0.640** | **0.645** | **0.641** |

Table 6: Binary Classification (Bert-Juman).

|  | Acc ↑ | Pre ↑ | Rec ↑ | F ↑ |
|---|---|---|---|---|
| Little | 0.563 | 0.577 | 0.584 | 0.580 |
| Half | 0.573 | 0.578 | 0.575 | 0.576 |
| All | **0.603** | **0.588** | **0.598** | **0.591** |

### 4.3.2 Multinomial Classification

The results of True / False / Unclear multinomial classifications are shown in Tables 7, 8 and 9 below.

The result of dataset w2v is Table 7, Bert-Mecab is Table 8, and Bert-Juman is Table 9.

Table 7: Multinomial Classification (w2v).

|  | Acc ↑ | Pre ↑ | Rec ↑ | F ↑ |
|---|---|---|---|---|
| Little | 0.477 | 0.419 | 0.428 | 0.421 |
| Half | 0.530 | **0.477** | **0.484** | **0.477** |
| All | **0.537** | 0.449 | 0.455 | 0.446 |

Table 8: Multinomial Classification (Bert-Mecab).

|  | Acc ↑ | Pre ↑ | Rec ↑ | F ↑ |
|---|---|---|---|---|
| Little | 0.513 | 0.426 | 0.403 | 0.410 |
| Half | 0.527 | 0.441 | 0.437 | 0.437 |
| All | **0.543** | **0.471** | **0.452** | **0.455** |

Table 9: Multinomial Classification (Bert-Juman).

|  | Acc ↑ | Pre ↑ | Rec ↑ | F ↑ |
|---|---|---|---|---|
| Little | 0.470 | 0.396 | 0.382 | 0.386 |
| Half | 0.490 | 0.446 | 0.433 | 0.434 |
| All | **0.547** | **0.470** | **0.458** | **0.460** |

## 5 DISCUSSION

### 5.1 Experimental Discussion

The results of both binary and multinomial classification experiments show that the larger the

amount of propagation information used, we got the better results. These results indicate that propagation information is effective to some extent in detecting rumors in Japanese. In addition, there were no significant differences in the results between Bert-Mecab and Bert-Juman in terms of text features, but Bert-Mecab calculated better results for Half, which uses half of the reply data arranged in time-series order for both classifications. It can also be seen that the dataset with Bert calculated slightly better results than w2v for both classifications. This suggested that the acquisition of embedding using BERT is somewhat effective for this task. On the other hand, the accuracy of the results was generally poor. We think that there are several possible reasons for this fact.

The first point is the input of the detection model. In this study, the model was designed to divide the input into two parts, basic features, and text features. And, these features were applied convolution using GCN layers and combined with the results. However, since the basic features have category values represented like 0 and 1. The accuracy is expected to be improved by further separating the categorical and non-categorical values from the basic features as input to the detection model.

The second point is related to text features. In this study, the BERT model was used to obtain embedded expressions because we thought this method could get more semantic embedded representations than those obtained by using feature words. However, both BERT models used Japanese Wikipedia data as the pre-training corpus. So, we believe that the results would be better if we utilize the BERT model which pre-trained a large corpus of more colloquial sentences for getting embedded representations. Another possibility for improving the results is that the huge number of embedding dimensions of BERT could have been compressed by principal component analysis (PCA) or other methods before being assigned to the model. Besides, it is thought that the Sentence-BERT model (Reimers, N. et al. 2019) may be more suitable for acquiring sentence vectors than the usual BERT model when viewed in terms of sentence units.

Also, we consider that we can use other embedded methods which have been proposed by other research. Some example methods are text-CNN (Yoon, K., 2014), Rumor2Vec (Tu, K., et al. 2021), and character-level deep learning embedding. (Sato, M., et al. 2018) Besides, although we used the detection model only GCN layers as in this study, we expect that incorporating a recursive layer such as RNN or LSTM that takes into account more time-series

information will improve the accuracy by reflecting more information on tweets posted in the early stage of posting.

## 5.2 Dataset Discussion

We created a Japanese rumor dataset based on previous studies (Liu, X., et al. 2015, Ma, J., et al. 2016) in this paper. However, there are some points to be considered for this dataset as well. Unlike Snopes.com, the FIJ data may make the problem of multinomial classification more difficult because the number of tweets that belong to the True Group is small. In this regard, we thought of an approach to increase the amount of data by creating pseudo extended data using replies as the source tweets by focusing on whether the replies affirm or deny the postings as shown in Figure 4. However, this approach has some issues, such as how to deal with the case of Unclear tweets, and the fact that denial/affirmation does not always correspond to True/False.
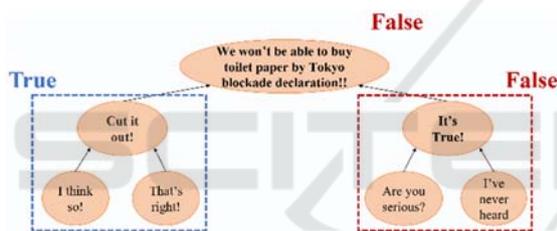


Figure 4: An example of data extension.

In addition, it is necessary to increase the evaluation score of the dataset Little when considering early rumor detection close to the time of source tweet posting. One solution for this problem is to increase the number of feature types. For example, we believe the emotional information of a tweet or the textual information in a summary sentence can be extracted as separate features.

## 6 CONCLUSION

In this study, we created a dataset of Japanese fact-checked tweets on Twitter and performed binary classification which classified a tweet as False or not, and multinomial classification which classified a tweet as True Rumor, Unclear Rumor, or False rumor. Moreover, we examined and discussed Japanese rumor detection. Future prospects are to apply this method to other languages than Japanese and to experiment and confirm the method's accuracy.

Conversely, it is necessary to reconfirm the evaluation score of the current model by comparing the results with existing Rumor Detection models for other languages. Also, one of our future goals is to build a more practical fact-checking support system. To do so, we need to improve the detection evaluation score by building other models, putting in additional features, and studying more effectively.

## REFERENCES

Islam, M. S., Sarkar, T., Khan, S. H., Kamal, A. H. M., Hasan, S. M., Kabir, A., ... and Seale, H. (2020). COVID-19–related infodemic and its impact on public health: A global social media analysis. *The American journal of tropical medicine and hygiene*,

Liu, X., Nourbakhsh, A., Li, Q., Fang, R., Shah, S., (2015). Real-time Rumor Debunking on Twitter, *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pp.1867-1870

Ma, J., Gao, W., Mitra, P., Kwon, S., Jansen, B. J., Wong, K. F., & Cha, M. (2016). Detecting rumors from microblogs with recurrent neural networks., *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI2016)*, pp.3818-3824.

Ma, J., Gao, W., Wong, K. F. (2018). Rumor detection on twitter with tree-structured recursive neural networks., *Association for Computational Linguistics*, pp.1980-1989.

Bian, T., Xiao, X., Xu, T., Zhao, P., Huang, W., Rong, Y., & Huang, J. (2020). Rumor detection on social media with bi-directional graph convolutional networks. In *Proceedings of the AAAI conference on artificial intelligence*, pp. 549-556.

Li, J., Bao, P., Shen, H., & Li, X. (2021). MiSTR: A Multiview Structural-Temporal Learning Framework for Rumor Detection. *IEEE Transactions on Big Data*.

Lao, A., Shi, C., & Yang, Y. (2021). Rumor detection with field of linear and non-linear propagation. In *Proceedings of the Web Conference 2021*, pp. 3178-3187.

Cheng, M., Li, Y., Nazarian, S., Bogdan P., (2021) From rumor to genetic mutation detection with explanations: a GAN approach. *Nature Scientific Reports*,

Ma, J., Li, J., Gao, W., Yang, Y., & Wong, K. F. (2021). Improving Rumor Detection by Promoting Information Campaigns with Transformer-based Generative Adversarial Learning. *IEEE Transactions on Knowledge and Data Engineering*.

Alzanin, S. M., & Azmi, A. M. (2019). Rumor detection in Arabic tweets using semi-supervised and unsupervised expectation–maximization. *Knowledge-Based Systems*,

Devlin, J., Chang, M. W., Lee, K., Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810. 04805*

Shibata, T., Kawahara, D., Kurohashi, S., (2019). Improved accuracy of Japanese parsing with BERT (in Japanese)

Yoon, K., (2014), Convolutional Neural Networks for Sentence Classification, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp.1746-1751

Reimers, N., Gurevych, I., (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP),* pp.3973-3983.

Tu, K., Chen, C., Hou, C., Yuan, J., Li, J., and Yuan, X. (2021). Rumor2vec: a rumor detection framework with joint text and propagation structure representation learning. *Information Sciences*, pp.137-151.

Sato, M., Orihara, R., Sei, Y., Tahara, Y., Ohsuga, A. (2018). Text Classification and Transfer Learning Based on Character-Level Deep Convolutional Neural Networks, *International Conference on Agents and Artificial Intelligence (ICAART)*