

Zero-Shot Rumor Detection with Propagation Structure via Prompt Learning

Hongzhan Lin^{1*}, Pengyao Yi^{2*}, Jing Ma^{1†}, Haiyun Jiang^{3†},
Ziyang Luo¹, Shuming Shi⁴, Ruifang Liu²

¹Hong Kong Baptist University

²Beijing University of Posts and Telecommunications

³Fudan University

⁴Tsinghua University

{cshzlin, majing, cszylo}@comp.hkbu.edu.hk, jianglu@fudan.edu.cn, {yi.py, lrf}@bupt.edu.cn, ssm01@hotmail.com

Abstract

The spread of rumors along with breaking events seriously hinders the truth in the era of social media. Previous studies reveal that due to the lack of annotated resources, rumors presented in minority languages are hard to be detected. Furthermore, the unforeseen breaking events not involved in yesterday’s news exacerbate the scarcity of data resources. In this work, we propose a novel zero-shot framework based on prompt learning to detect rumors falling in different domains or presented in different languages. More specifically, we firstly represent rumor circulated on social media as diverse propagation threads, then design a hierarchical prompt encoding mechanism to learn language-agnostic contextual representations for both prompts and rumor data. To further enhance domain adaptation, we model the domain-invariant structural features from the propagation threads, to incorporate structural position representations of influential community response. In addition, a new virtual response augmentation method is used to improve model training. Extensive experiments conducted on three real-world datasets demonstrate that our proposed model achieves much better performance than state-of-the-art methods and exhibits a superior capacity for detecting rumors at early stages.

1 Introduction

The spread of rumors emerging along with breaking news is a global phenomenon, which can cause critical consequences for social network users in different lingual contexts. For example, during the unprecedented COVID-19 pandemic, a false rumor claiming that “the vaccine has a chip in it which will control your mind”¹ released by a Muslim cleric, went viral on Facebook and Twitter in different languages. Such misleading claims about vaccines are being shared widely in many countries, which confuse the public and undermine their enthusiasm for vaccination. Due to the barriers of domain and language, even human fact-checkers are poor judges of such rumors. Therefore, it’s imperative to

develop automatic approaches for rumor detection spread in different languages amid unforeseen breaking events.

Social psychology literature defines a rumor as a story or a statement whose truth value is unverified or deliberately false (Allport and Postman 1945). In this study, we focus on detecting rumors on social media, instead of “fake news” strictly defined as a news article published by a news outlet that is verifiably false (Yang et al. 2022). State-of-the-art techniques using deep neural networks (DNNs) (Bian et al. 2020; Lin et al. 2021a; Rao et al. 2021) have promoted the development of rumor detection, but they are all data-driven models that require extensive annotated data for model training. Most corpora are open-domain and presented in English, which makes them not scalable to emerging events in new languages where only few/no labeled data is available. Zero-shot rumor detection task (ZRD) aims to adapt knowledge learned in the source rumor data to the target data without labeled training samples in the target language and domain, as shown in Figure 1. Previous related studies (Du et al. 2021; Tian, Zhang, and Lau 2021) directly utilize pre-trained language models (PLMs) (Devlin et al. 2019) to fine-tune on ZRD task. However, they just formulated the zero-shot rumor detection as a cross-lingual text classification problem and detected the single claim post with a heavy task-specific fine-tuning stage, which makes it deviate from the pre-training target on masked language modeling, even ignoring the domain-invariant interaction of user opinions during the diffusion of rumors. More recently, Lin et al. (2022) propose a contrastive learning framework to detect rumors from different languages and domains, where a small number of target annotation is required. However, it is prone to be poor at emerging events propagated in minority languages without any expertise annotation, especially in some underdeveloped countries and regions. For breaking events with scarce annotated data in different languages, the study on the zero-shot regimes is more urgent and practical for rumor detection on social media.

In this paper, we focus on exploring efficient prompting with language and domain transfer for zero-shot rumor detection. We assume there are no accessible annotations in the target language and domain, so prompt learning mechanisms (Zhao and Schütze 2021) based on existing multi-

*These authors contributed equally.

†Corresponding authors.

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹<https://www.bbc.com/news/55768656>

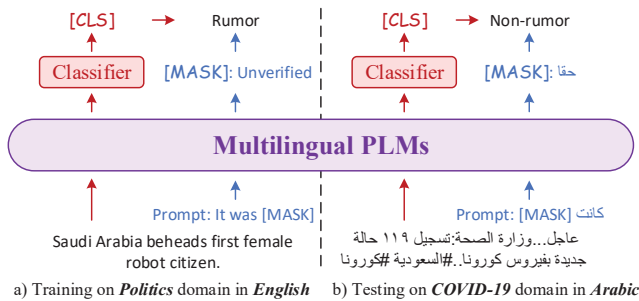


Figure 1: Illustration between the task-specific fine-tuning and the prompt learning paradigms for solving ZRD task.

lingual PLMs can be utilized. However, the standard prompt learning paradigm adopts discrete or soft prompts, where the discrete prompt requires experts of native speakers to design rumor-related templates/rules for different languages, and the soft prompt uses optimized token representations trained on a large dataset. Unlike the standard prompt-tuning paradigm, we propose to decouple shared semantic information from the syntactic bias in specific languages based on multilingual PLMs, which could enhance the semantic interaction between the prompt and rumor data. Besides, as the diffusion of rumors generally follows spatial and temporal relations that provide valuable clues on how a claim is transmitted irrespective of specific domains (Zubiaga et al. 2018), we aim to develop a novel prompt learning mechanism to take such social context into consideration.

To this end, we propose a zero-shot Response-aware Prompt Learning (RPL) framework to detect cross-lingual and cross-domain rumors on social media. More specifically, we firstly rank responsive posts toward the claim to represent diverse propagation threads. Then a hierarchical prompt encoding mechanism is proposed based on multilingual PLMs, which alleviates the effort of prompt designing for different languages. On the other hand, as the propagation structure contains domain-invariant features on how a claim is responded to by users over time, we model the absolute and relative propagation position to capture the latent structure of the propagation thread for better domain adaptation. To further improve the zero-shot model training, we incorporate a new virtual response augmentation mechanism into the prompt learning framework. As there is no public benchmark available for detecting rumors in low-resource languages with propagation threads in tweets, we collected a new rumor dataset corresponding to COVID-19 from Twitter in Cantonese and Arabic languages. Extensive experiments conducted on three real-world rumor datasets corresponding to COVID-19 confirm that (1) our model yields outstanding performance for detecting zero-shot rumors over the state-of-the-art baselines with a large margin; and (2) our method performs particularly well on early rumor detection which is crucial for timely intervention and debunking especially for breaking events.

2 Related Work

Pioneer studies for automatic rumor detection focused on learning a supervised classifier utilizing features crafted

from post contents, user profiles, and propagation patterns (Castillo, Mendoza, and Poblete 2011; Yang et al. 2012; Liu et al. 2015). Subsequent studies then proposed new features such as those representing rumor diffusion and cascades (Kwon et al. 2013; Friggeri et al. 2014; Hannak et al. 2014). Zhao, Resnick, and Mei (2015) alleviated the engineering effort by using a set of regular expressions to find questing and denying tweets. Deep neural networks such as recurrent neural networks (Ma et al. 2016), convolutional neural networks (Yu et al. 2017), and attention mechanism (Guo et al. 2018) were then employed to learn the features from the stream of social media posts. To extract useful clues jointly from content semantics and propagation structures, some approaches proposed kernel-learning models (Wu, Yang, and Zhu 2015; Ma, Gao, and Wong 2017), tree-structured recursive neural networks (RvNN) (Ma, Gao, and Wong 2018), self-attention models (PLAN (Khoo et al. 2020), STANKER (Rao et al. 2021)) and graph neural networks (BiGCN) (Bian et al. 2020) have been exploited to encode conversation threads for higher-level representations.

Recently, zero-shot transfer learning techniques are applied on PLMs to detect fake news (Du et al. 2021; Schwarz, Theóphilo, and Rocha 2020; De et al. 2021) by downstream task-specific fine-tuning methods. Tian, Zhang, and Lau (2021) utilized PLMs and a self-training loop to adapt the model from the source language to the target language in multi-step iteration. However, these approaches only consider cross-lingual text classification and face problems such as task mismatch between pre-training and fine-tuning, and ignore the domain-invariant propagation patterns from community response. Considering that rumors can be domain-specific and/or presented in different languages, Lin et al. (2022) first introduced supervised contrastive learning for few-shot rumor detection based on propagation structure. However, their few-shot paradigm still relies on a small number of target data for training, which cannot perform well in detecting more minority-language rumor data without any expertise annotation in case of emerging topics.

Prompt learning converts downstream tasks to language modeling tasks via textual prompts, which is found more effective to use PLMs than typical fine-tuning on specific tasks (Brown et al. 2020; Liu et al. 2021a). In recent years, prompt learning has achieved great success in a variety of NLP tasks, such as text classification (Min et al. 2022), semantic parsing (Schucher, Reddy, and de Vries 2021), text generation (Li and Liang 2021), sentiment classification (Seoh et al. 2021) and dialog state tracking (Lee, Cheng, and Ostendorf 2021), etc. Despite the flourish of the research in prompting methods, there is only limited attention being put on the low-resource rumor detection task (Lin et al. 2022). Different from a few previous multilingual work (Zhao and Schütze 2021; Winata et al. 2021; Lin et al. 2021b) on either discrete or soft (Liu et al. 2021b; Lester, Al-Rfou, and Constant 2021) prompts, in this paper, we tune the level-grained models’ parameters for language-agnostic rumor prompts, which further attends to user interactions from community response for zero-shot rumor detection task.

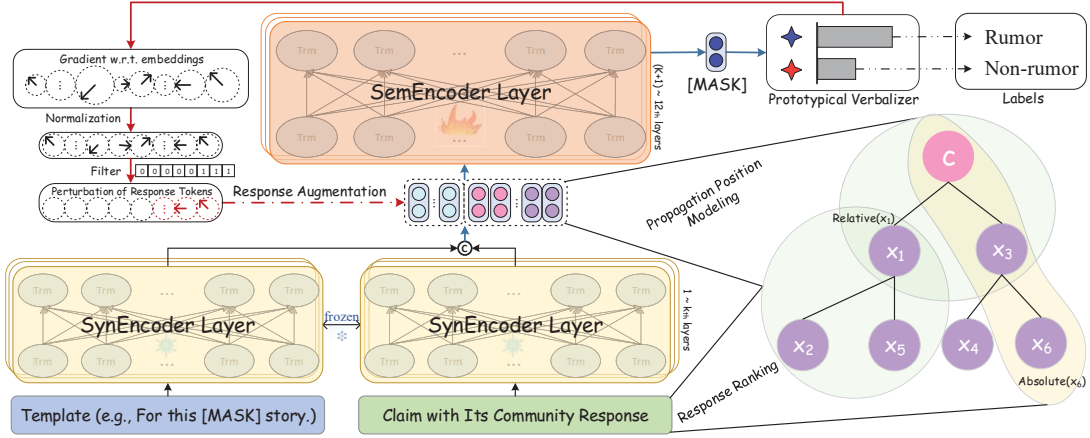


Figure 2: The overall architecture of our proposed method. For the source training data, we first obtain intermediate syntax-independent embeddings after the SynEncoder, then tune the SemEncoder with the prototypical verbalizer paradigm. For target test data, the similarity score between the output states of $[MASK]$ and prototypes would be used to detect rumors.

3 Problem Statement and Background

In this work, we define the zero-shot rumor detection task as: given a dataset as source, classify each event in the target dataset as a rumor or not, where the source and target data are from different languages and domains. Specifically, we define a source dataset for training as a set of events $\mathcal{D}_s = \{C_1^s, C_2^s, \dots, C_M^s\}$, where M is the number of source events. Each event $C^s = (y, c, \mathcal{T}(c))$ is a triplet representing a given claim c which is associated with a veracity label $y \in \{\text{rumor}, \text{non-rumor}\}$, and ideally all its relevant responsive microblog post in chronological order, i.e., $\mathcal{T}(c) = [x_1^s, x_2^s, \dots, x_m^s]$, where m is the number of responsive posts in the conversation thread. We consider the target dataset with a different language and domain from the source dataset for testing $\mathcal{D}_t = \{C_1^t, C_2^t, \dots, C_N^t\}$, where N is the number of target events and each $C^t = (c', \mathcal{T}(c'))$ shares the similar structure as that of the source.

This task could be formulated as a supervised classification problem that trains a language/domain-agnostic classifier $f(\cdot)$ transferring the features learned from source datasets to that of the target events, that is, $f(C^t | \mathcal{D}_s) \rightarrow y$.

In this work, we convert the rumor detection as a cloze-style masked language modeling problem. For example, given a cloze-style template p (e.g., “For this $[MASK]$ story.”) as the prompt, spliced with the claim c into \hat{c} , the standard prompt learning leverages PLMs to obtain the hidden state for the $[MASK]$ token, to infer the rumor-indicative words to fill in $[MASK]$. The probability of label y is:

$$\mathcal{P}(y|\hat{c}) = g(\mathcal{P}([MASK] = v|\hat{c})|v \in \mathcal{V}_y) \quad (1)$$

where \mathcal{V} is a set of rumor-related label words, \mathcal{V}_y is the subset of \mathcal{V} corresponding to y and $g(\cdot)$ is a manual verbalizer to transform the probability of label words into that of the label. In this way, we could map predicted words for $[MASK]$ into the veracity label to make a decision on the claim.

4 Our Approach

In this section, we introduce our Response-aware Prompt Learning framework for zero-shot rumor detection. Because

the rumor-related prompt design for different languages can be biased and labor-intensive, we propose to learn language-independent prompts including the template and verbalizer. On another hand, as the responsive posts could provide a domain-invariant propagation structure for representation learning, we explore how to fuse such community response into the prompt learning framework. Figure 2 illustrates an overview of our proposed model: 1) Response Ranking, which presents each event as diverse propagation threads following temporal or spatial relations; 2) Hierarchical Prompt Encoding, which is the backbone to learn language-independent interaction between the prompt and the event with the prior knowledge of multilingual PLMs; 3) Propagation Position Modeling, which equips our proposed prompt-based framework with the latent structure of the propagation thread; and 4) Response Augmentation, which adds noise to responsive posts to improve model training for better robustness.

4.1 Response Ranking

To highlight the social context for enhancing the contextual representation learning for the event, we propose to attend over evidential responses. The core idea is to rank all the responses based on diverse propagation threads.

First, we hypothesize that the attitudes of responsive posts towards the claim will become more inclined as time goes by, thus the responsive posts can be sorted in chronological and inverted order on the time sequence. Specifically, for the chronological order, responsive posts with earlier time stamps are prioritized, i.e., $\mathcal{T}(c) = [x_1, x_2, \dots, x_m]$, and vice versa for the inverted order on the time sequence, i.e., $\mathcal{T}(c) = [x_m, x_{m-1}, \dots, x_1]$.

Besides the perspective of time sequence, inspired by (Ma, Gao, and Wong 2018; Bian et al. 2020), we further represent the propagation thread as a tree structure $\mathcal{T}(c) = \langle \mathcal{G}, \vec{\mathcal{E}} \rangle$, where \mathcal{G} refers to a set of nodes each representing a responsive post of c , and $\vec{\mathcal{E}}$ is a set of directed paths conforming to the responsive relation among the nodes in \mathcal{G} . We scrutinize the optimal search algorithms on the tree

structure to select more evidential posts in depth-first and breadth-first order. Specifically, the depth-first search studies the propagation patterns during information flows from the ancestor to the children nodes while the breadth-first search gives priority to the interaction of user opinions among sibling nodes. Taking the propagation tree in Figure 2 as an example, the depth-first order of the response ranking would be $[x_1, x_2, x_5, x_3, x_4, x_6]$; for the breadth-first order, it would be $[x_1, x_3, x_2, x_4, x_5, x_6]$.

In this way, concerning perspectives of time sequence or propagation tree, we could investigate the importance of different responses $\mathcal{T}(c)$ on the verdict of a claim.

4.2 Hierarchical Prompt Encoding

Generally, it will lead to bias towards syntax in specific languages if we directly utilize the existing tokens from the vocabulary like expertise words or language-specific slang as the template. To bridge the gap between languages in this task, the template shall not depend on any specific language. Although the soft prompt is a potential way to solve this problem, its trainable tokens require enough target rumor data for training, which is challenging in zero-shot regimes. To this end, we hope to implicitly disentangle shared semantic information from different languages with language-specific syntactic knowledge, by leveraging the priors of multilingual PLMs. Previous literature (Jawahar, Sagot, and Seddah 2019; Rao et al. 2021; Huang et al. 2022) has shown that the lower layers of PLMs can capture syntactic-level features while the upper layers of PLMs model the semantic-level features. Therefore, we can present a Hierarchical Prompt Encoding (HPE) mechanism for language-independent representation learning of the template and the event at syntactic and semantic levels. In our approach, we hypothesize that the semantic information could be shared over different languages though the syntax is language-dependent.

SynEncoder Layer. At the syntactic level, to obtain the intermediate syntax-independent embeddings, we copy and froze the parameters of the lower k layers from multilingual PLMs encoders to encode the template and the event, respectively. Specifically, the original template p is syntactically mapped into a shared vector space:

$$X_p = \text{SynEncoder}(p) \quad (2)$$

where $X_p \in \mathbb{R}^{|p| \times d}$ is the template embeddings and d is the dimension of the output state of SynEncoder.

For an event C , as all the responsive posts are presented in the same language and domain as the claim either at the training or testing stages, we could concatenate them in the same frozen SynEncoder to obtain the embeddings of the event:

$$X_{cr} = \text{SynEncoder}([c, \mathcal{T}(c)]) \quad (3)$$

where $[\cdot, \cdot]$ means the splicing operation, $X_{cr} \in \mathbb{R}^{o \times d}$ is the embeddings of the event (i.e., claim with its community response), o is the maximum sequence length of PLMs. Based on the obtained response ranking, the contextually coherent posts could be potentially retained from the perspectives of the temporal and spatial relations, respectively, under the input length restriction of PLMs (Devlin et al. 2019).

SemEncoder Layer. At the semantic level, we initialize a trainable semantic encoder with the $(k+1)^{th}$ layer to the top layer of the PLMs. Then we concatenate and refine the output states of the template and the event, on top of the frozen SynEncoder, to further model the semantic interaction between the template and the event:

$$H = \text{SemEncoder}([X_p, X_{cr}]) \quad (4)$$

In summary, we map the simple English discrete prompt into a shared embedding space at the syntactic level by the prior knowledge of SynEncoder, which is then fed into the SemEncoder for semantic interaction with the event. On top of the SemEncoder, we present a prototypical verbalizer to map the output states H^m of $[MASK]$ token into the label y without manual rumor-related label words for specific languages, which would be depicted in Sec. 4.5.

4.3 Propagation Position Modeling

To bridge the prompt learning and propagation structures for zero-shot rumor detection on social media, we further propose *Absolute* and *Relative* Propagation Position Modeling, to inject the propagation information into the tunable SemEncoder for domain-invariant structural features extraction at the semantic level.

For Absolute Propagation Position, we exploit the propagation path of a responsive post in the propagation tree, which is complementary to its sequential counterpart (Devlin et al. 2019). Specifically, given a token q from a post x_i , we treat the claim c of the event as the root and use the distance of the responsive path from the current post to the root as the absolute propagation position: $abs_{pro}(q) = \text{distance}_{tree}(x_i, c)$, where $tree$ is the propagation structure $\mathcal{T}(c) = \langle \mathcal{G}, \vec{\mathcal{E}} \rangle$. Note that in this work, we make the tokens in the same post share the propagation position of the post in the propagation tree. Thus we update the input representation of the token q for the tunable SemEncoder by summing the corresponding token embeddings in X_{cr} and its absolute position embeddings, where the absolute position embeddings are trained with learnable parameters (Gehring et al. 2017).

For Relative Propagation Position, we mainly focus on the local context of a responsive post in the propagation tree as its relative propagation position. As each post in the propagation tree may trigger a set of responsive posts, we aim to capture the relative user opinions among responsive posts in such a subtree structure. Specifically, towards a post x_i , we consider the relative posts with five relationships in the subtree as the relative propagation position: 1) *Parent*⁽⁺⁾; 2) *Children*⁽⁻⁾; 3) *Siblings*⁽⁺⁾; 4) *Siblings*⁽⁻⁾; 5) *Itself*, where $+/-$ denotes the relative post comes earlier/later than the current post in the subtree. We then extend the self-attention computation to consider the pairwise relationships among posts in the same subtree and project the relative propagation position into the SemEncoder by drawing the practice of Shaw, Uszkoreit, and Vaswani (2018). In this way, the relative propagation patterns in a local subtree can be captured explicitly as users share opinions towards the same subtree root, to cross-check the inaccurate information.

4.4 Response Augmentation

Since the model could suffer from noisy responses, we propose to enhance the prompt learning by creating additional adversarial examples. We present a new virtual response augmentation algorithm, ViRA, a variant of the virtual adversarial algorithm (Miyato et al. 2018). To create an adversarial example, we apply Fast Gradient Value (Rozsa, Rudd, and Boulton 2016) to approximate a worst-case perturbation, where the gradient is normalized to represent the direction that significantly decreases the model’s performance, and a norm is used to ensure the approximation is reasonable. However, the value ranges (norms) of the embedding vectors vary among different data and models. The variance gets larger for bigger models with billions of parameters, leading to some instability of adversarial training. To this end, we first apply layer normalization (Ba, Kiros, and Hinton 2016) on top of the frozen SynEncoder to normalize the embeddings into stochastic vectors, and then perform a mask operation to filter out the template and claim embeddings, lastly add the perturbation to the normalized embedding vectors of responsive posts. Adversarial noise enables the model to handle extensive noisy responsive posts and can be regarded as a response augmentation mechanism.

4.5 Model Training

On top of the SemEncoder where the template and a event sample (i.e., a claim and its responsive posts) could be transformed into a shared semantic latent space, inspired by Prototypical Networks (Snell, Swersky, and Zemel 2017; Lin, Yan, and Chen 2021), we further introduce a prototypical verbalizer paradigm to prevent the rumor-related label words from heavily relying on the language-specific expert words. The core idea is to utilize the representative features of instances from the same classes for encapsulating event-level semantic features instead of the language-dependent label words.

Given the $[MASK]$ token representation H_i^m of a training example C_i , we minimize a prototypical loss as follows:

$$\mathcal{L}_{proto} = -\log \frac{e^{\mathcal{S}(H_i^m, l_y)}}{\sum_{y'} e^{\mathcal{S}(H_i^m, l_{y'})}} \quad (5)$$

where y is the ground truth of H_i^m , \mathcal{S} denotes the normalized cosine similarity score. l_y denotes the learnable prototype vectors of the class y , which is the cluster representative of the embedded support points belonging to the class. By optimizing the above objective function \mathcal{L}_{proto} , rumor features can be close to corresponding rumor prototype in semantic space and be away from the non-rumor prototype.

In addition, we adopt the contrastive loss to pull up the intra-class variance and down the inter-class variance of instances in a batch:

$$\mathcal{L}_{con} = -\frac{1}{B_{y_i} - 1} \sum_j \mathbb{1}_{[i \neq j]} \mathbb{1}_{[y_i = y_j]} \frac{e^{\mathcal{S}(H_i^m, H_j^m)}}{\sum_{j'} \mathbb{1}_{[i \neq j']} e^{\mathcal{S}(H_i^m, H_{j'}^m)}} \quad (6)$$

where B_{y_i} is the number of source examples with the same

label y_i in the event C_i in a batch, and $\mathbb{1}$ is the indicator.

We jointly train the model with the prototypical contrastive objectives: $\mathcal{L} = \alpha \mathcal{L}_{proto} + (1 - \alpha) \mathcal{L}_{con}$, where α is a trade-off parameter set as 0.5 in our experiments. So we generate a pseudo augmented example for C_i based on response augmentation, which is again fed into the tunable SemEncoder to compute the new loss $\tilde{\mathcal{L}}$. Finally, we use the average loss $\mathcal{L}_{avg} = \text{mean}(\mathcal{L} + \tilde{\mathcal{L}})$ for the back-propagation (Collobert et al. 2011) with the AdamW optimizer (Loshchilov and Hutter 2018). We set the layer number k of the SynEncoder as 6. The learning rate is initialized as 1e-5. Early stopping (Yao, Rosasco, and Caponnetto 2007) is applied to avoid overfitting.

5 Experiments

5.1 Datasets

We utilize FOUR public datasets TWITTER, WEIBO (Ma et al. 2016), Twitter-COVID19 and Weibo-COVID19 (Lin et al. 2022) for experiments. TWITTER and Twitter-COVID19 are English rumor datasets with conversation thread in tweets while WEIBO and Weibo-COVID19 are Chinese rumor datasets with the similar composition structure. Furthermore, as there are no public benchmarks available for detecting rumors in low-resource languages with propagation structure in tweets, we organized and constructed a new low-resource rumor dataset CatAr-COVID19. Specifically, we resort to two COVID-19 rumor datasets (Alam et al. 2021; Ke et al. 2020), which only contain multilingual textual claims in Cantonese and Arabic without propagation thread. We extend each claim by collecting its propagation thread via Twitter academic API in python. Finally, we annotated the claim tweets by referring to the labels of the events from the original datasets².

5.2 Experimental Setup

We compare our model with several state-of-the-art zero-shot rumor detection systems: 1) **Vanilla-Ft**: Fine-tune the model for classification by adding a task-specific linear layer with the $[CLS]$ token on top of PLMs (Devlin et al. 2019); 2) **Translate-Ft**: Utilize rumor data in source language for training and translate the claim into target languages for testing (Du et al. 2021); 3) **Contrast-Ft**: We employ and extend an existing few-shot learning technique, supervised contrastive learning (Lin et al. 2022), to fine-tuning on the source data in the zero-shot scenario; 4) **Adapter**: Fix the parameters of PLMs and add only a few trainable parameters per task within a residual adapter (Houlsby et al. 2019); 5) **Parallel-Adpt**: An adapter-based variant (He et al. 2021) by transferring the parallel insertion of prefix tuning into adapters; 6) **Source-Ppt**: A prompt-based tuning method (Lin et al. 2021b) both trains and tests the model by prompt in source languages; 7) **Translate-Ppt**: Train on prompts in the source language and test on the target-lingual prompts after translation (Zhao and Schütze

²Our code and resources will be available at https://github.com/PengyaoYi/zeroRumor_AAAI

Source	TWITTER								WEIBO							
Target	Weibo-COVID19				CatAr-COVID19				Twitter-COVID19				CatAr-COVID19			
Model	Acc.	Mac- F_1	R F_1	NR F_1	Acc.	Mac- F_1	R F_1	NR F_1	Acc.	Mac- F_1	R F_1	NR F_1	Acc.	Mac- F_1	R F_1	NR F_1
Vanilla-Ft	0.623	0.585	0.711	0.459	0.518	0.402	0.583	0.220	0.603	0.602	0.619	0.585	0.481	0.481	0.479	0.474
Translate-Ft	0.639	0.567	0.745	0.388	0.523	0.457	0.637	0.277	0.634	0.574	0.653	0.495	0.505	0.512	0.528	0.496
Contrast-Ft	0.656	0.582	0.759	0.405	0.584	0.458	0.720	0.196	0.653	0.644	0.699	0.590	0.562	0.561	0.571	0.551
Adapter	0.644	0.600	0.737	0.463	0.558	0.438	0.665	0.211	0.652	0.612	0.736	0.487	0.548	0.556	0.605	0.508
Parallel-Adpt	0.651	0.598	0.730	0.467	0.567	0.450	0.701	0.198	0.667	0.653	0.731	0.574	0.579	0.585	0.636	0.534
Source-Ppt	0.664	0.648	0.722	0.574	0.589	0.564	0.460	0.669	0.670	0.616	0.760	0.472	0.599	0.565	0.688	0.441
Translate-Ppt	0.650	0.489	0.776	0.201	0.573	0.568	0.519	0.617	0.674	0.651	0.740	0.562	0.604	0.542	0.374	0.711
Soft-Ppt	0.652	0.574	0.756	0.392	0.590	0.565	0.446	0.683	0.685	0.652	0.758	0.546	0.609	0.575	0.518	0.633
RPL-Cho	0.713	0.675	0.786	0.563	0.613	0.581	0.455	0.707	0.715	0.689	0.778	0.601	0.634	0.633	0.616	0.650
RPL-Inv	0.728	0.666	0.810	0.521	0.601	0.592	0.473	0.711	0.733	0.710	0.788	0.632	0.647	0.640	0.586	0.693
RPL-Dep	0.732	0.689	0.805	0.574	0.640	0.619	0.530	0.708	0.723	0.711	0.771	0.650	0.657	0.636	0.547	0.724
RPL-Bre	0.745	0.719	0.804	0.634	0.631	0.617	0.544	0.689	0.727	0.697	0.793	0.601	0.672	0.664	0.614	0.714

Table 1: Rumor detection results on the target test datasets.

2021); 8) **Soft-Ppt**: Instead of discrete tokens, tunable tokens (Lester, Al-Rfou, and Constant 2021) are utilized as the prompt; 9) **RPL-***: Our proposed response-aware prompt learning framework with the diverse propagation threads, i.e., chronological (Cho) and inverted (Inv) order in time sequence, depth-first (Dep) and breadth-first (Bre) order in tree structure.

In this work, we consider the most challenging case, i.e., detecting events (i.e., target) from a new domain and language. Specifically, we use the well-resourced TWITTER (Ma, Gao, and Wong 2017) and WEIBO (Ma et al. 2016) datasets as the source data, and Weibo-COVID19, Twitter-COVID19 and CatAr-COVID19 datasets as the target. We use accuracy and macro-averaged F1, as well as class-specific F1 scores as the evaluation metrics.

5.3 Rumor Detection Performance

Table 1 shows the performance of our proposed method versus all the compared methods on the Weibo-COVID19, Twitter-COVID19 and CatAr-COVID19 datasets with pre-determined training datasets. From Table 1, it is observed that the performance of the baselines in the first group are obviously poor due to heavy reliance on downstream classification objectives with a task-related linear layer added on top of PLMs, which is randomly initialized and too easily overfit the source data to generalize to the target.

The prompt-based baselines in the third group are relatively better than the adapter-based baselines in the second group though Soft-Ppt is somewhat related to the adapter style in the form of parameter tuning (He et al. 2021). However, their performance are still limited to the following reasons: 1) Source-Ppt lacks cross-lingual transferability. Generally, the multilingual PLMs cannot deal well with the cross-lingual combination between the template in the source language and the claim post in the target language, where such data format is rarely seen in the pre-training stage. 2) Translate-Ppt easily suffers from error propagation of the machine translation quality, and the language-agnostic knowledge is not decoupled and transferred from the source template to the target. 3) Soft-Ppt requires abundant target rumor data for sufficient optimization, which cannot be sat-

isfied with the zero-shot setting.

In contrast, our proposed RPL-based approaches achieve superior performance among all the baselines, which suggests their strong generalization for zero-shot transfer between different languages and different domains. It’s observed that the performance of RPL-Inv is relatively better than that of RPL-Cho. We speculate that the reason is that questioning posts at the later stage of propagation could indicate a higher tendency that the claim is rumor or not. Although achieving promising performance, RPL-Dep does not achieve the expected best performance because with the propagation of the claim there is more semantic and structural information but the noisy information is increased simultaneously, especially in the vein of relatively deep conversation or argument. Overall, RPL-Bre obtains stable and excellent performance generally among the four RPL-based variants by making full use of the subtree-structure property via breadth-first ranking and propagation position modeling for response fusion, which verifies that inaccurate information on social media can be “self-checked” by making a comparison with responsive posts towards the same topic.

5.4 Ablation Study

We perform ablation studies by discarding some important components of our best-performed approach RPL-Bre on CatAr-COVID19, which include 1) *w/o RR*: We simply encode the claim without the Response Ranking (RR) strategies that consider the social contexts in community response. 2) *w/o APP*: We discard the Absolute Propagation Position as presented in Eq.?? 3) *w/o RPP*: We discard the Relative Propagation Position (RPP). 4) *w/o ViRA*: We neglect the Virtual Response Augmentation (ViRA) mechanism. 5) *w/o HPE*: Instead of our proposed Hierarchical Prompt Encoding (HPE) mechanism, we devise our backbone as two tiers of transformers: one for encoding all the responsive posts independently, and another for processing the sequence of posts using representations from the first transformer (i.e., PLMs), where the second-tier transformer has a similar architecture to PLMs, but has only 2 layers and its parameters are initialized randomly. 6) *w/o PV*: We design a manual verbalizer for label mapping, to replace the

Source	TWITTER		WEIBO	
Model	Acc.	Mac- F_1	Acc.	Mac- F_1
RPL-Bre	0.631	0.617	0.672	0.664
RPL-Bre w/o RR	0.605	0.598	0.613	0.611
RPL-Bre w/o APP	0.622	0.607	0.626	0.624
RPL-Bre w/o RPP	0.610	0.601	0.633	0.632
RPL-Bre w/o ViRA	0.626	0.612	0.644	0.634
RPL-Bre w/o HPE	0.571	0.451	0.581	0.433
RPL-Bre w/o PV	0.592	0.589	0.621	0.617

Table 2: Ablation studies on our proposed model.

Prototypical Verbalizer (PV) for model training.

As demonstrated in Table 2, the ablative models suffer different degrees of such performance degradation, indicating the effectiveness of our proposed components for adapting features learned from source rumor data to that of the target. Specifically, *RPL-Bre*'s performance significantly decreases without response ranking due to the lack of collective wisdom on social media. Both *w/o APP* and *w/o RPP* also achieve worse performance than *RPL-Bre*, suggesting that both perspectives of propagation position modeling are comparably helpful to the domain-variant propagation patterns extraction in zero-shot regimes; *RPL-Bre* makes improvements over *w/o ViRA*, which implies the promoting role of ViRA that enables our approach hardly compromised when the input length is limited and there may be noise in response. Moreover, *w/o HPE* leads to much performance degradation, which implies the prompt encoding framework ingeniously reserves the prior syntactic and semantic knowledge from the PLMs and contributes more accurate zero-shot rumor predictions with language disentanglement. Compared with *RPL-Bre*, the performance of *w/o PV* also significantly decreases, highlighting the importance and complementary of the prototypical paradigm in our framework for language and domain adaptation.

5.5 Early Detection

Early alerts of rumors can prevent the wide-spreading of rumorous contents. By setting detection checkpoints of "delays" that can be either the count of reply posts or the time elapsed since the first posting, only contents posted no later than the checkpoints is available for model evaluation. The performance is evaluated by Macro F1 obtained at each checkpoint. To satisfy each checkpoint, we incrementally scan test data in order of time until the target time delay or post volume is reached.

Figure 3 shows the early detection performance of our approach versus Soft-Ppt, PLAN, STANKER, BiGCN and RvNN at various deadlines. To make fair comparisons, the inputs of all baselines are encoded with the same multilingual PLM. We observe that our proposed RPL-based approach outperforms other baselines throughout the whole lifecycle, and reaches a relatively high Macro F1 score at a very early period after the initial broadcast. One interesting phenomenon is that our method only needs about 20 posts on CatAr-COVID19 and 4 hours on Twitter-COVID19, to achieve the saturated performance, indicating the advanced response fusion strategy and remarkably superior early detection performance of our method.

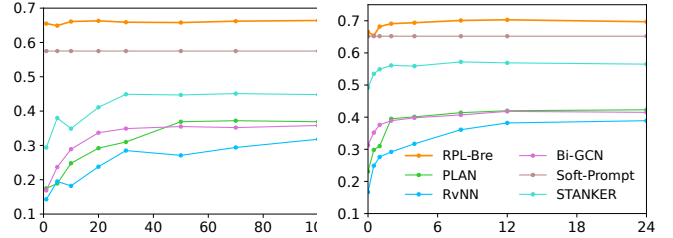


Figure 3: Early detection performance at different checkpoints of posts count on CatAr-COVID19 (left) and elapsed time on Twitter-COVID19 (right) datasets.

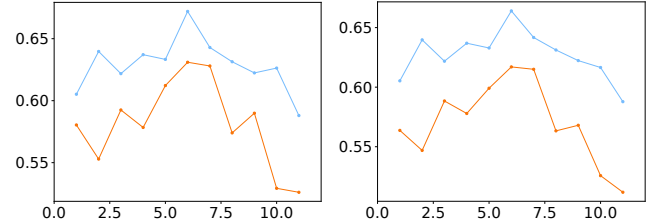


Figure 4: Effect of the layer number k of SynEncoder with Accuracy (left) and Macro F1 (right).

5.6 Discussion

Figure 4 shows the effect of layer number (from 1 to 12) of the SynEncoder on zero-shot rumor detection performance, with the CatAr-COVID19 as the target, TWITTER (in orange) and WEIBO (in blue) as the source data, respectively. We can observe that when the SynEncoder is initialized with the lower 4 layers of PLMs, it is still biased to specific languages due to the surface features mainly learned. Since PLMs could unearth rich linguistic features at the lower 6 layers, the best performance is obtained when k is set to 6 (i.e., the setting in our model), which is in line with the finding of Jawahar, Sagot, and Seddah (2019). After that, as k continues to increase, although the capacity to decouple shared semanteme from specific linguistic features is enhanced, the number of SemEncoder layers with prior semantic knowledge activated for the interaction of prompts and events decreases, thus the generalization ability of the model to rumor data in different domains is limited, resulting in a fluctuated decline of performance.

6 Conclusion and Future Work

In this paper, we propose a zero-shot Response-aware Prompt Learning framework to bridge language and domain gaps in rumor detection. We present a prompt-based approach to avoid the reliance on language-specific rumor prompt engineering, with effective response fusion strategies to incorporate influential and structural propagation threads for domain adaptation. Results on three real-world benchmarks confirm the advantages of our zero-shot detection model. For future work, we plan to study specialized PLMs for rumor detection to better utilize the wisdom of crowds and circumvent the sequence length limit, then collect and apply our model to more languages and domains.

Acknowledgments

Work partially done when Hongzhan Lin was an intern at Tencent AI Lab. This work was partially supported by Hong Kong RGC ECS (Ref. 22200722) and HKBU One-off Tier 2 Start-up Grant (Ref. RCOFSGT2/20-21/SCI/004).

References

- Alam, F.; Shaar, S.; Dalvi, F.; Sajjad, H.; Nikolov, A.; Mubarak, H.; Da San Martino, G.; Abdelali, A.; Durrani, N.; Darwish, K.; et al. 2021. Fighting the COVID-19 Infodemic: Modeling the Perspective of Journalists, Fact-Checkers, Social Media Platforms, Policy Makers, and the Society. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 611–649.
- Allport, G. W.; and Postman, L. J. 1945. Section of psychology: The basic psychology of rumor. *Transactions of the New York Academy of Sciences*, 8(2 Series II): 61–81.
- Ba, J. L.; Kiros, J. R.; and Hinton, G. E. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Bian, T.; Xiao, X.; Xu, T.; Zhao, P.; Huang, W.; Rong, Y.; and Huang, J. 2020. Rumor Detection on Social Media with Bi-Directional Graph Convolutional Networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*.
- Castillo, C.; Mendoza, M.; and Poblete, B. 2011. Information credibility on twitter. In *Proceedings of the 20th international conference on World wide web*, 675–684.
- Collobert, R.; Weston, J.; Bottou, L.; Karlen, M.; Kavukcuoglu, K.; and Kuksa, P. 2011. Natural language processing (almost) from scratch. *Journal of machine learning research*, 12(ARTICLE): 2493–2537.
- De, A.; Bandyopadhyay, D.; Gain, B.; and Ekbal, A. 2021. A Transformer-Based Approach to Multilingual Fake News Detection in Low-Resource Languages. *Transactions on Asian and Low-Resource Language Information Processing*.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186.
- Du, J.; Dou, Y.; Xia, C.; Cui, L.; Ma, J.; and Philip, S. Y. 2021. Cross-lingual COVID-19 Fake News Detection. In *2021 International Conference on Data Mining Workshops (ICDMW)*, 859–862. IEEE.
- Friggeri, A.; Adamic, L.; Eckles, D.; and Cheng, J. 2014. Rumor cascades. In *Eighth international AAAI conference on weblogs and social media*.
- Gehring, J.; Auli, M.; Grangier, D.; Yarats, D.; and Dauphin, Y. N. 2017. Convolutional sequence to sequence learning. In *International conference on machine learning*. PMLR.
- Guo, H.; Cao, J.; Zhang, Y.; Guo, J.; and Li, J. 2018. Rumor detection with hierarchical social attention network. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, 943–951.
- Hannak, A.; Margolin, D.; Keegan, B.; and Weber, I. 2014. Get Back! You Don’t Know Me Like That: The Social Mediation of Fact Checking Interventions in Twitter Conversations. In *ICWSM*.
- He, J.; Zhou, C.; Ma, X.; Berg-Kirkpatrick, T.; and Neubig, G. 2021. Towards a Unified View of Parameter-Efficient Transfer Learning. In *International Conference on Learning Representations*.
- Houlsby, N.; Giurgiu, A.; Jastrzebski, S.; Morrone, B.; De Laroussilhe, Q.; Gesmundo, A.; Attariyan, M.; and Gelly, S. 2019. Parameter-efficient transfer learning for NLP. In *International Conference on Machine Learning*, 2790–2799. PMLR.
- Huang, L.; Ma, S.; Zhang, D.; Wei, F.; and Wang, H. 2022. Zero-shot Cross-lingual Transfer of Prompt-based Tuning with a Unified Multilingual Prompt. *arXiv preprint arXiv:2202.11451*.
- Jawahar, G.; Sagot, B.; and Seddah, D. 2019. What does BERT learn about the structure of language? In *ACL 2019-57th Annual Meeting of the Association for Computational Linguistics*.
- Ke, L.; Chen, X.; Lu, Z.; Su, H.; and Wang, H. 2020. A novel approach for cantonese rumor detection based on deep neural network. In *2020 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, 1610–1615. IEEE.
- Khoo, L. M. S.; Chieu, H. L.; Qian, Z.; and Jiang, J. 2020. Interpretable rumor detection in microblogs by attending to user interactions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 8783–8790.
- Kwon, S.; Cha, M.; Jung, K.; Chen, W.; and Wang, Y. 2013. Prominent features of rumor propagation in online social media. In *2013 IEEE 13th International Conference on Data Mining*, 1103–1108. IEEE.
- Lee, C.-H.; Cheng, H.; and Ostendorf, M. 2021. Dialogue State Tracking with a Language Model using Schema-Driven Prompting. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 4937–4949.
- Lester, B.; Al-Rfou, R.; and Constant, N. 2021. The Power of Scale for Parameter-Efficient Prompt Tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 3045–3059.
- Li, X. L.; and Liang, P. 2021. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*.
- Lin, H.; Ma, J.; Chen, L.; Yang, Z.; Cheng, M.; and Guang, C. 2022. Detect Rumors in Microblog Posts for Low-Resource Domains via Adversarial Contrastive Learning. In *Findings of the Association for Computational Linguistics: NAACL 2022*.

- Lin, H.; Ma, J.; Cheng, M.; Yang, Z.; Chen, L.; and Chen, G. 2021a. Rumor detection on twitter with claim-guided hierarchical graph attention networks. In *Proceedings of the 2021 conference on empirical methods in natural language processing (EMNLP)*.
- Lin, H.; Yan, Y.; and Chen, G. 2021. Boosting Low-Resource Intent Detection with in-Scope Prototypical Networks. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*.
- Lin, X. V.; Mihaylov, T.; Artetxe, M.; Wang, T.; Chen, S.; Simig, D.; Ott, M.; Goyal, N.; Bhosale, S.; Du, J.; et al. 2021b. Few-shot Learning with Multilingual Language Models. *arXiv preprint arXiv:2112.10668*.
- Liu, P.; Yuan, W.; Fu, J.; Jiang, Z.; Hayashi, H.; and Neubig, G. 2021a. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.
- Liu, X.; Ji, K.; Fu, Y.; Du, Z.; Yang, Z.; and Tang, J. 2021b. P-Tuning v2: Prompt Tuning Can Be Comparable to Fine-tuning Universally Across Scales and Tasks. *arXiv preprint arXiv:2110.07602*.
- Liu, X.; Nourbakhsh, A.; Li, Q.; Fang, R.; and Shah, S. 2015. Real-time rumor debunking on twitter. In *Proceedings of the 24th ACM International Conference on Information and Knowledge Management*, 1867–1870.
- Loshchilov, I.; and Hutter, F. 2018. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*.
- Ma, J.; Gao, W.; Mitra, P.; Kwon, S.; Jansen, B. J.; Wong, K.-F.; and Cha, M. 2016. Detecting rumors from microblogs with recurrent neural networks. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, 3818–3824.
- Ma, J.; Gao, W.; and Wong, K.-F. 2017. Detect Rumors in Microblog Posts Using Propagation Structure via Kernel Learning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*.
- Ma, J.; Gao, W.; and Wong, K.-F. 2018. Rumor Detection on Twitter with Tree-structured Recursive Neural Networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*.
- Min, S.; Lewis, M.; Hajishirzi, H.; and Zettlemoyer, L. 2022. Noisy Channel Language Model Prompting for Few-Shot Text Classification. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*.
- Miyato, T.; Maeda, S.-i.; Koyama, M.; and Ishii, S. 2018. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*.
- Rao, D.; Miao, X.; Jiang, Z.; and Li, R. 2021. STANKER: Stacking Network based on Level-grained Attention-masked BERT for Rumor Detection on Social Media. In *EMNLP*.
- Rozsa, A.; Rudd, E. M.; and Boulton, T. E. 2016. Adversarial diversity and hard positive generation. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 25–32.
- Schucher, N.; Reddy, S.; and de Vries, H. 2021. The Power of Prompt Tuning for Low-Resource Semantic Parsing. *arXiv preprint arXiv:2110.08525*.
- Schwarz, S.; Theóphilo, A.; and Rocha, A. 2020. Emet: Embeddings from multilingual-encoder transformer for fake news detection. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing*.
- Seoh, R.; Birle, I.; Tak, M.; Chang, H.-S.; Pinette, B.; and Hough, A. 2021. Open Aspect Target Sentiment Classification with Natural Language Prompts. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 6311–6322.
- Shaw, P.; Uszkoreit, J.; and Vaswani, A. 2018. Self-Attention with Relative Position Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*.
- Snell, J.; Swersky, K.; and Zemel, R. 2017. Prototypical networks for few-shot learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 4080–4090.
- Tian, L.; Zhang, X.; and Lau, J. H. 2021. Rumour Detection via Zero-Shot Cross-Lingual Transfer Learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 603–618. Springer.
- Winata, G. I.; Madotto, A.; Lin, Z.; Liu, R.; Yosinski, J.; and Fung, P. 2021. Language Models are Few-shot Multilingual Learners. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, 1–15.
- Wu, K.; Yang, S.; and Zhu, K. Q. 2015. False rumors detection on sina weibo by propagation structures. In *2015 IEEE 31st international conference on data engineering*. IEEE.
- Yang, F.; Liu, Y.; Yu, X.; and Yang, M. 2012. Automatic detection of rumor on sina weibo. In *Proceedings of the ACM SIGKDD workshop on mining data semantics*, 1–7.
- Yang, Z.; Ma, J.; Chen, H.; Lin, H.; Luo, Z.; and Chang, Y. 2022. A Coarse-to-fine Cascaded Evidence-Distillation Neural Network for Explainable Fake News Detection. In *Proceedings of the 29th International Conference on Computational Linguistics*, 2608–2621.
- Yao, Y.; Rosasco, L.; and Caponnetto, A. 2007. On early stopping in gradient descent learning. *Constructive Approximation*, 26(2): 289–315.
- Yu, F.; Liu, Q.; Wu, S.; Wang, L.; and Tan, T. 2017. A Convolutional Approach for Misinformation Identification. In *IJCAI*.
- Zhao, M.; and Schütze, H. 2021. Discrete and Soft Prompting for Multilingual Models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Zhao, Z.; Resnick, P.; and Mei, Q. 2015. Enquiring minds: Early detection of rumors in social media from enquiry posts. In *Proceedings of the 24th international conference on world wide web*, 1395–1405.
- Zubiaga, A.; Aker, A.; Bontcheva, K.; Liakata, M.; and Procter, R. 2018. Detection and resolution of rumours in social media: A survey. *ACM Computing Surveys (CSUR)*.