

# TwHIN-BERT: A Socially-Enriched Pre-trained Language Model for Multilingual Tweet Representations at Twitter

Xinyang Zhang\*  
xz43@illinois.edu  
The University of Illinois at  
Urbana-Champaign  
Urbana, IL, USA

Yury Malkov  
ymalkov@twitter.com  
Twitter Cortex  
San Francisco, CA, USA

Omar Florez  
oflorez@twitter.com  
Twitter Cortex  
San Francisco, CA, USA

Serim Park  
serimp@twitter.com  
Twitter Cortex  
San Francisco, CA, USA

Brian McWilliams  
brimcwilliams@twitter.com  
Twitter Cortex  
San Francisco, CA, USA

Jiawei Han  
hanj@illinois.edu  
The University of Illinois at  
Urbana-Champaign  
Urbana, IL, USA

Ahmed El-Kishky\*  
aekishky@twitter.com  
Twitter Cortex  
San Francisco, CA, USA

## ABSTRACT

Pre-trained language models (PLMs) are fundamental for natural language processing applications. Most existing PLMs are not tailored to the noisy user-generated text on social media, and the pre-training does not factor in the valuable social engagement logs available in a social network. We present TwHIN-BERT, a multilingual language model productionized at Twitter, trained on in-domain data from the popular social network. TwHIN-BERT differs from prior pre-trained language models as it is trained with not only text-based self-supervision but also with a social objective based on the rich social engagements within a Twitter heterogeneous information network (TwHIN). Our model is trained on 7 billion tweets covering over 100 distinct languages, providing a valuable representation to model short, noisy, user-generated text. We evaluate our model on various multilingual social recommendation and semantic understanding tasks and demonstrate significant metric improvement over established pre-trained language models. We open-source TwHIN-BERT and our curated hashtag prediction and social engagement benchmark datasets to the research community<sup>1</sup>.

## CCS CONCEPTS

- **Computing methodologies** → **Natural language processing**;
- **Information systems** → **Social networks**.

## KEYWORDS

language models, social media, social engagement

### ACM Reference Format:

Xinyang Zhang, Yury Malkov, Omar Florez, Serim Park, Brian McWilliams, Jiawei Han, and Ahmed El-Kishky. 2023. TwHIN-BERT: A Socially-Enriched

\* Corresponding authors: xz43@illinois.edu, aekishky@twitter.com

<sup>1</sup><https://github.com/xinyangz/TwHIN-BERT>.

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

*KDD '23, August 6–10, 2023, Long Beach, CA, USA*

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0103-0/23/08...\$15.00

<https://doi.org/10.1145/3580305.3599921>

Pre-trained Language Model for Multilingual Tweet Representations at Twitter. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23), August 6–10, 2023, Long Beach, CA, USA*. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3580305.3599921>

## 1 INTRODUCTION

The proliferation of pre-trained language models (PLMs) [12, 14] based on the Transformer architecture [47] has pushed state of the art across many tasks in natural language processing (NLP). As an application of transfer learning, these models are typically trained on massive text corpora and, when fine-tuned on downstream tasks, have demonstrated state-of-the-art performance.

Despite the success of PLMs in general-domain NLP, fewer attempts have been made in language model pre-training for user-generated text on social media. In this work, we pre-train a language model for Twitter – a prominent social media platform where users post short messages called Tweets. Tweets contain informal diction, abbreviations, emojis, and topical tokens such as hashtags. As a result, PLMs designed for general text corpora may struggle to understand Tweet semantics accurately. Existing works [2, 32] on Twitter LM pre-training do not address these challenges and simply replicate general domain pre-training on Twitter corpora.

A distinctive feature of Twitter social media is the user interactions through Tweet engagements. As seen in Figure 1, when a user visits Twitter, in addition to posting Tweets, they can perform a variety of *social actions* such as “Favoriting”, “Replying” and “Retweeting” Tweets. The wealth of such engagement information is invaluable to Tweet content understanding. For example, the post “bottom of the ninth, two outs, and down by one!!” would be connected to baseball topics by its co-engaged Tweets, such as “three strikes and you’re out!!!”. Without the social contexts, a conventional text-only PLM objective would struggle to build this connection. As an additional benefit, a socially-enriched language model will also vastly benefit common applications on social media, such as social recommendations [53] and information diffusion prediction [10, 40].

We introduce TwHIN-BERT – a multilingual language model for Twitter pre-trained with social engagements. The key idea of our method is to leverage *socially similar Tweets* for pre-training. Building on this idea, TwHIN-BERT has the following features. (1) We



**Figure 1: (a) This mock-up shows a short-text Tweet and social engagements such as Faves, Retweets, Replies, Follows that create a social context to Tweets and signify Tweet appeal to engaging users. (b) Co-engagement is a strong indicator of Tweet similarity.**

construct a Twitter Heterogeneous Information Network (TwHIN) [18] to unify the multi-typed user engagement logs. Then, we run scalable embedding and approximate nearest neighbor search to sift through hundreds of billions of engagement records and mine socially similar Tweet pairs. (2) In conjunction with masked language modeling, we introduce a contrastive social objective that enforces the model to tell if a pair of Tweets are socially similar or not. Our model is trained on 7 billion Tweets from over 100 distinct languages, of which 1 billion have social engagement logs.

We evaluate the TwHIN-BERT model on both social recommendation and semantic understanding downstream evaluation tasks. To comprehensively evaluate on many languages, we curate two large-scale datasets, a social engagement prediction dataset focused on social aspects and a hashtag prediction dataset focused on language aspects. In addition to these two curated datasets, we also evaluate on established benchmark datasets to draw direct comparisons to other available pre-trained language models. TwHIN-BERT achieves state-of-the-art performance in our evaluations with a major advantage in the social tasks.

In summary, our contributions are as follows:

- We build the first ever socially-enriched pre-trained language model for noisy user-generated text on Twitter.
- Our model is the strongest multilingual Twitter PLM so far, covering 100 distinct languages.
- Our model has a major advantage in capturing the social appeal of Tweets.
- We open-source TwHIN-BERT as well as two new Tweet benchmark datasets: (1) hashtag prediction and (2) social engagement prediction.

## 2 TWHIN-BERT

In this section, we outline how we construct training examples for our social pre-training objectives and subsequently train TwHIN-BERT with social and text objectives. As seen in Figure 2, we first construct and embed a user-Tweet engagement network. The resultant Tweet embeddings are then used to mine pairs of socially similar Tweets. These Tweet pairs and others are used to pre-train TwHIN-BERT, which can then be fine-tuned for various downstream tasks.

### 2.1 Mining Socially Similar Tweets

With abundant social engagement logs, we (informally) define socially similar Tweets as *Tweets that are co-engaged by a similar set of users*. The challenge lies in how to implement this social similarity by (1) fusing heterogeneous engagement types, such as “Favorite”, “Reply”, “Retweet”, and (2) efficiently mining billions of similar Tweet pairs.

To address these challenges, TwHIN-BERT first constructs a **Twitter Heterogeneous Information Network (TwHIN)** from the engagement logs, then runs a scalable heterogeneous network embedding method to capture co-engagement and map Tweets and users into a vector space. With this, social similarity translates to embedding space similarity. Subsequently, we mine similar Tweet pairs via ANN search on the Tweet embeddings.

**2.1.1 Constructing TwHIN.** We define and construct TwHIN as:

*Definition 2.1 (TwHIN).* Our Twitter Heterogeneous Information Network is a directed bipartite graph  $G = (U, T, E, \phi)$ , where  $U$  is the set of user nodes,  $T$  is the set of Tweet nodes,  $E = U \times T$  is the set of engagement edges.  $\phi : E \mapsto \mathcal{R}$  is an edge type mapping function. Each edge  $e \in E$  belongs to a type of engagement in  $\mathcal{R}$ .

Our curated TwHIN (Figure 3) consists of approximately 200 million distinct users, 1 billion Tweets, and over 100 billion edges. We posit that our TwHIN encodes not only user preferences but also Tweet social appeal. We perform scalable network embedding to derive a social similarity metric from TwHIN. The network embedding fuses the heterogeneous engagements into a unified vector space that’s easy to operate on.

**2.1.2 Embedding TwHIN Nodes.** We seek to learn shallow embedding vectors (i.e., vector of learnable parameters) for each user ( $u_j$ ) and Tweet ( $t_k$ ) in the TwHIN; we denote these learnable embeddings for users and Tweets as  $\mathbf{u}_j$  and  $\mathbf{t}_k$  respectively. While our approach is agnostic to the exact methodology used to embed TwHIN, we follow the approach outlined in [17, 18]. A user-Tweet pair for a particular relation type  $\phi((u_j, t_k)) = r_m$  is scored with a scoring function of the form  $f(u_j, t_k, r_m)$ . Our training objective seeks to learn  $\mathbf{u}$ ,  $\mathbf{t}$  and  $\mathbf{r}$  parameters that maximize a log-likelihood constructed from the scoring function for  $(u, t) \in G$  and minimize for  $(u, t) \notin G$ .

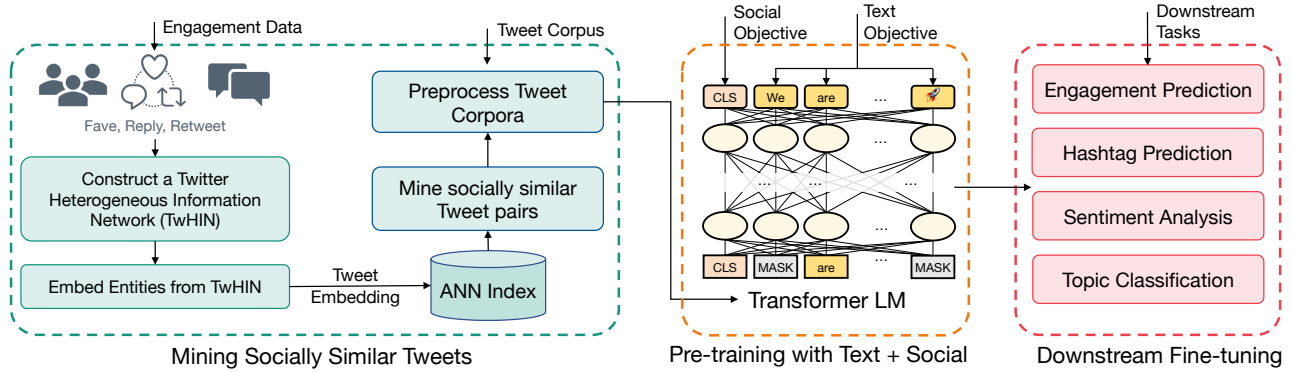
For simplicity, we apply a simple dot product comparison between user and Tweet representations. For a user-tweet edge  $(u_j, t_k)$  of relation  $r_m$ , this operation is defined by:

$$f(e) = f(u_j, t_k, r_m) = (\mathbf{u}_j + \mathbf{r}_m)^\top \mathbf{t}_k \quad (1)$$

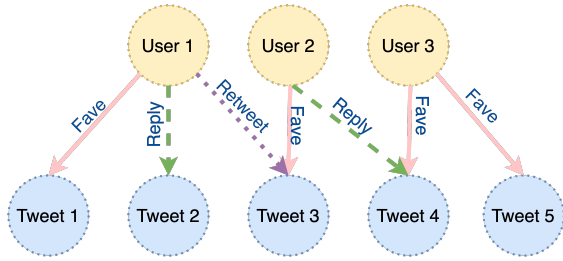
As seen in Equation 1, we co-embed users and Tweets and scoring is performed by applying an engagement-specific translation embedding to user representations and computing the dot product with a Tweet representation. The task is then formulated as an edge (or link) prediction task. Following previous method [18, 20, 30], we maximize the following negative sampling objective:

$$\arg\max_{\mathbf{u}, \mathbf{r}, \mathbf{t}} \sum_{e \in G} \left[ \log \sigma(f(e)) + \sum_{e' \in N(e)} \log \sigma(-f(e')) \right] \quad (2)$$

where  $N(e)$  is a set of negatively sampled edges by corrupting positive edges via replacing either the user or Tweet in an edge with



**Figure 2: We outline the end-to-end TwHIN-BERT process. This three-step process involves (1) mining socially similar Tweet pairs by embedding a Twitter Heterogeneous Information Network (2) training TwHIN-BERT using a joint social and MLM objective and finally (3) fine-tuning TwHIN-BERT on downstream tasks.**



**Figure 3: Twitter Heterogeneous Information Network (TwHIN) capturing social engagements between users and Tweets.**

a negatively sampled user or Tweet. As user-Tweet engagement graphs are very sparse, randomly corrupting an edge in the graph is very likely to be a ‘negative’ edge absent from the graph.

Equation 2 represents the log-likelihood of predicting a binary “real” or “fake” label for the set of edges in the network (real) along with a set of the “fake” negatively sampled edges. To maximize the objective, we learn  $\mathbf{u}$  and  $\mathbf{i}$  parameters to differentiate positive edges from negative, unobserved edges.

We adopt the PyTorch-Biggraph [24] framework for scalability. Following previous approaches, we train for 10 epochs and perform negative sampling both uniformly and proportional to entity prevalence in TwHIN [5, 24]. Optimization is via Adagrad.

Upon learning dense representations of nodes in TwHIN, we utilize the Tweet representations to mine socially similar Tweets.

**2.1.3 Mining Similar Tweet Pairs.** Given the learned TwHIN Tweet embeddings, we seek to identify pairs of Tweets with *similar social appeal* – that is, Tweets that appeal to (i.e., are likely to be engaged with) similar users. We will use these socially-similar Tweet pairs as self-supervision when training TwHIN-BERT. To identify these pairs, we perform an approximate nearest neighbor (ANN) search in the TwHIN embedding space. To efficiently perform the search over 1B+ Tweets, we use the optimized FAISS<sup>1</sup> toolkit [23] to create a compact index of Tweets keyed by their engagement-based TwHIN embeddings. As each Tweet embedding is 256-dimensional, storing billion-scale Tweet embeddings would require more than one TB of

memory. To reduce the size of the index such that it can fit on a 16 A100 GPU node, with each GPU possessing 40GB of memory, we apply product quantization [22] to discretize and reduce embeddings size. The resultant index corresponds to OPQ64, IVF65536, PQ64 in the FAISS index factory terminology.

After creating the FAISS index and populating it with TwHIN Tweet embeddings, we search the index using Tweet embedding queries to find pairs of similar Tweets  $(t_i, t_j)$  such that  $t_i$  and  $t_j$  are close in the embedding space as defined by their cosine distance. To ensure high recall, we query the FAISS index with 2000 probes. Finally, we select the  $k$  closest Tweets with the cosine distance between the query Tweet and retrieved Tweets’ embeddings. These pairs are used in our social objective when pre-training TwHIN-BERT.

## 2.2 Pre-training Objectives

Given the mined *socially similar* Tweets, we describe our language model training process. To train TwHIN-BERT, we first run the Tweets through the language model and then train the model with a joint contrastive social loss and masked language model loss.

**Tweet Encoding with LM.** We use a Transformer language model to encode each Tweet. Similar to BERT [14], given the tokenized text  $\mathbf{w}_t = [w_1, w_2, \dots, w_n]$  of a Tweet  $t$ , we add special tokens to mark the start and end of the Tweet:  $\hat{\mathbf{w}}_t = [\text{CLS}]\mathbf{w}_t[\text{SEP}]$ . As the Tweets are usually shorter than the maximum sequence length of a language model, we group multiple Tweets and feed them together into the language model when possible. We then apply *CLS-pooling*, which takes the [CLS] token embedding of each Tweet. These Tweet embeddings are passed through an MLP projection head for the *social loss* computation.

$$[\mathbf{e}_{t_1}, \mathbf{e}_{t_2}, \dots] = \text{Pool}(\text{LM}([\hat{\mathbf{w}}_{t_1}, \hat{\mathbf{w}}_{t_2}, \dots])) \quad (3)$$

$$\mathbf{z}_t = \text{MLP}(\mathbf{e}_t) \quad (4)$$

**Contrastive Social Loss.** We use a contrastive loss to let our model learn whether two Tweets are socially similar or not. For each batch of  $B$  socially similar Tweet pairs  $\{(t_i, t_j)\}_B$ , we compute the *NT-Xent* loss [8] with in-batch negatives:

$$\mathcal{L}_{\text{social}}(i, j) = -\log \frac{\exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_j)) / \tau}{\sum_{N_B(i)} \exp(\text{sim}(\mathbf{z}_i, \mathbf{z}_k) / \tau)} \quad (5)$$

<sup>1</sup><https://github.com/facebookresearch/faiss>

The negatives  $\mathcal{N}_B(i)$  of Tweet  $t_i$  are the  $(2B - 1)$  other Tweets in the batch that are not paired with  $t_i$ . We use cosine similarity for function  $\text{sim}(\cdot, \cdot)$ .  $\tau$  is the loss temperature.

Our overall pre-training objective is a combination of the contrastive social loss and the masked language model loss [14]:

$$\mathcal{L} = \mathcal{L}_{\text{social}} + \lambda \mathcal{L}_{\text{MLM}} \quad (6)$$

$\lambda$  is a hyperparameter that balances the social and language loss.

## 2.3 Pre-training Setup

**Model Architecture.** We use the same Transformer architecture as BERT [14] for our language model. We adopt the XLM-R [12] tokenizer, which offers good capacity and coverage in all languages. The model has a vocabulary size of 250K. The max sequence length is set to 128 tokens. The detailed model setup can be found in Appendix B. Note that although we have chosen this specific architecture, our social objective can be used in conjunction with a wide range of language model architectures.

**Pre-training Data.** We collect 7 billion Tweets in 100 languages from Jan. 2020 to Jun. 2022. Additionally, we collect 100 billion user-Tweet social engagement data covering 1 billion of our Tweets. We re-sample based on language frequency raised to the power of 0.7 to mitigate the under-representation of low-resource languages.

**Training Procedure.** Our training has two stages. In the first stage, we train the model from scratch using the 6 billion Tweets without user engagement. The model is trained for 500K steps on 16 Nvidia A100 GPUs (a2-megagpu-16g) with a total batch size of 6K. In the second stage, the model is trained for another 500K steps on the 1 billion Tweets with the joint MLM and social loss. We use mixed precision during training. Overall pre-training takes approximately five days for the base model and two weeks for the large model. We refer readers to Appendix B for the detailed hyperparameter setup.

## 3 EXPERIMENTS

In this section, we discuss baseline specifications, evaluation setup, and results from two families of downstream evaluation tasks.

### 3.1 Evaluated Methods

We evaluate TwHIN-BERT against the following baselines.

- **mBERT** [14] is the multilingual language variant of the popular BERT [14] language model. It is a general domain language model trained on Wikipedia dumps.
- **XLM-R** [12] is a state-of-the-art general domain multilingual language model at its sizes. It is trained on over two terabytes of CommonCrawl data.
- **BERTweet** [32] is the previous state-of-the-art English tweet language model. It adopts a monolingual tokenizer trained from scratch on tweets and replicates RoBERTa [27] training from scratch on 845M English tweets.
- **XLM-T** [2] is a multilingual Twitter language model based on XLM-R [12]. It adopts the XLM-R tokenizer and model checkpoint and continues training on over 200M multilingual tweets.
- **TwHIN-BERT-MLM** is an ablation of our model. It is trained on the same corpus and with the same protocol as our main models. It uses only an MLM objective.

We include *base* and *large* sizes of our model train on the same corpus. All baselines are *base* variants (with between 135M to 278M parameters depending on the size of the tokenizer). Our large model has around 550M parameters.

We note that all externally published models we compared against were trained on different quantities of data and the data differed temporally and linguistically. As such, we include these models to demonstrate the gap in performance between widely-used published and open-sourced models and the model we plan on open-sourcing. On the other hand, our *base-MLM* model draws a direct comparison and isolates the effect of social engagement on the resultant model.

### 3.2 Social Engagement Prediction

Our first benchmark task is *social engagement prediction*. This task aims to evaluate how well the pre-trained language models capture the social aspects of user-generated text. In our task, we predict whether users modeled via a user embedding vector will perform a certain social engagement on a given Tweet.

We use different pre-trained language models to generate representations for Tweets, and then feed these representations into a simple prediction model alongside the corresponding user representation. The engagement prediction model is trained to predict whether a user will engage with a specific Tweet. The LM-generated embeddings are fixed when we train the downstream engagement prediction model.

**Dataset.** To curate our Tweet-Engagement dataset, we select the 50 popular languages on Twitter and sample 10,000 (or all if the total number is less than 10,000) Tweets of each language from a fixed time period. All Tweets are available via the Twitter public API. We then collect the user-Tweet engagement records associated with these Tweets. There are, on average, 29K engagement records per language. We ensure that there is no overlap between the evaluation and pre-training datasets.

Each engagement record consists of a pre-trained 256-dimensional user embedding [18] and a Tweet ID that indicates the user has engaged with the given Tweet. To ensure privacy, each user embedding appears only once, however, each tweet may be engaged by multiple users. We split the Tweets into train, development, and test sets with a 0.8/0.1/0.1 ratio, and then collect the respective engagement records for each subset.

**Prediction Model.** Given a pre-trained language model, we use it to generate an embedding for each Tweet  $t$  given its content  $w_t$ :  $e_t = \text{Pool}(\text{LM}(w_t))$ .

We apply the following pooling strategies to calculate the Tweet embedding from the language model. First, we take [CLS] token embedding as the first part. Then, we take the average token embedding of non-special tokens as the second part. The two parts are concatenated to form the *Combined* embedding of a Tweet.

With LM-derived Tweet embeddings, pre-trained user embeddings, and the user-Tweet engagement records, we build an engagement prediction model  $\Theta = (W_t, W_u)$ . Given a user  $u$  and a Tweet  $t$ , the model projects the user embedding  $e_u$  and the Tweet embedding  $e_t$  into the same space, and then calculates the probability of engagement:

**Table 1: Engagement prediction HITS@10 on high, mid, low-resource, and average of all languages.**

Method	High-Resource				Mid-Resource				Low-Resource				All
	en	ja	es	ar	el	ur	tl	nl	no	te	da	ps	Avg.
BERTweet	.1414	-	-	-	-	-	-	-	-	-	-	-	-
mBERT	.0633	.0227	.0575	.0532	.0496	.0437	.0610	.0616	.0731	.0279	.1060	.0522	.0732
XLm-R	.0850	.0947	.0704	.0546	.0628	.0315	.0653	.0650	.1661	.0505	.1150	.0727	.0849
XLm-T	.1181	.1079	.1103	.1403	.0562	.0352	.0877	.0762	.1156	.0728	.1167	.0662	.1043
<b>TwHIN-BERT</b>													
- Base-MLM	.1400	.1413	.1204	.1640	.0801	.0547	.0700	.0965	.1502	.0883	.1334	.0600	.1161
- Base	.1552	.2065	.1618	<b>.2206</b>	.0944	.0627	.1030	<b>.1346</b>	.1920	.1017	.1470	.0799	.1436
- Large	<b>.1585</b>	<b>.2325</b>	<b>.2055</b>	.1989	<b>.1065</b>	<b>.0667</b>	<b>.1053</b>	.1248	<b>.2118</b>	<b>.1654</b>	<b>.1475</b>	<b>.0817</b>	<b>.1497</b>

$$\mathbf{h}_u = \mathbf{W}_u^T \mathbf{e}_u, \quad \mathbf{h}_t = \mathbf{W}_t^T \mathbf{e}_t$$

$$P(t | u) = \sigma(\mathbf{h}_u^T \mathbf{h}_t)$$

We optimize a negative sampling loss on the training engagement records  $R$ . For each engagement pair  $(u, t) \in R$ , the loss is:

$$\log \sigma(\mathbf{h}_u^T \mathbf{h}_t) + \mathbb{E}_{t' \sim P_n(R)} \log \sigma(-\mathbf{h}_u^T \mathbf{h}_{t'})$$

where  $P_n(R)$  is a negative sampling distribution. We use the frequency of each Tweet in  $R$  to the power of  $3/4$  for this distribution.

Our prediction model closely resembles classical link prediction models such as [45]. We keep the model simple, making sure it will not overpower the language model embeddings.

**Evaluation Setup and Metrics.** We conduct a hyperparameter search on the English development dataset and use these hyperparameters for the other languages. The prediction model projects user and Tweet embedding to 128 dimensions. We set the batch size to 512, and the learning rate to  $1e-3$ . The best model on the validation set is selected for test set evaluation.

In the test set, we pair each user with 1,000 Tweets: one Tweet they have engaged with, and the rest are randomly sampled negatives. The model ranks the Tweets by the predicted probability of engagement, and we evaluate with HITS@10. We report median results from 6 runs with different initialization.

**Results.** We show summarized results for selected high, mid, and low-resource languages (determined by language frequency on Twitter) in Table 1. Language abbreviations are ISO language codes<sup>2</sup>. We also show the average results from all 50 languages in the evaluation dataset and leave the details in Table 6. Our TwHIN-BERT model demonstrates significant improvement over the baselines on the social engagement task. Comparing our model to the ablation without the social loss, we can see the contrastive social pre-training provides a significant lift over just MLM pre-training for social engagement prediction. An analysis of all 50 evaluation languages shows the *large* model to perform better than the *base* model on average, with more wins than losses. Additionally, we also observe that our method yields the most improvement when using the *Combined* [CLS] token and average non-special token embedding. We believe the [CLS] token embedding from our model

<sup>2</sup><https://www.iso.org/iso-639-language-codes.html>

**Table 2: Text classification dataset statistics. \*Statistics for Hashtag shows the numbers for each language.**

Dataset	Lang.	Label	Train	Dev	Test
SE2017	en	3	45,389	2,000	11,906
SE2018-en	en	20	45,000	5,000	50,000
SE2018-es	es	19	96,142	2,726	9,969
ASAD	ar	3	137,432	15,153	16,842
COVID-JA	ja	6	147,806	16,394	16,394
SE2020-hi	hi+en	3	14,000	3,000	3,000
SE2020-es	es+en	3	10,800	1,200	3,000
Hashtag	multi	500*	16,000*	2,000*	2,000*

captures the social aspects of the Tweet while averaging the other token embeddings captures the semantic aspects of the Tweet. Naturally, utilizing both aspects is essential to better model a Tweet’s appeal and a user’s inclination to engage with a Tweet.

### 3.3 Tweet Classification

Our second collection of downstream tasks is Tweet classification. In these tasks, we take as input the Tweet text and predict discrete labels corresponding to the label space for each task.

**Datasets.** We curate a multilingual Tweet hashtag prediction dataset (available via Twitter public API) to comprehensively cover the popular languages on Twitter. In addition, we evaluate on five external benchmark datasets for tasks such as sentiment classification, emoji prediction, and topic classification in selected languages. We show the dataset statistics in Table 2.

- **Tweet Hashtag Prediction** dataset is a multilingual hashtag prediction dataset we collected from Tweets. It contains Tweets from 50 popular languages. For each language, the 500 most popular hashtags were selected, and 100k tweets with those hashtags were sampled. We ensured each Tweet will only contain one of the 500 candidate hashtags. Similar to the work proposed in Miresghallah et al. [31], the task is to predict the hashtag used in the Tweet.
- **SemEval2017** task 4A [38] is a English Tweet sentiment analysis dataset. The labels are three-point sentiments of “positive”, “negative”, “neutral”.
- **ASAD** [1] is an Arabic Tweet sentiment dataset with the same three-point labels as SemEval2017 T4A.

**Table 3: Multilingual hashtag prediction Macro-F1 on high, mid, low resource, and average of all languages.**

Method	High-Resource				Mid-Resource				Low-Resource				All
	en	ja	es	ar	el	ur	tl	nl	no	te	da	ps	Avg.
BERTweet	59.01	-	-	-	-	-	-	-	-	-	-	-	-
mBERT	54.56	68.43	42.48	38.48	44.00	36.44	52.96	39.75	46.09	49.54	59.54	29.41	50.05
XML-R	53.90	69.07	43.80	37.85	43.94	37.56	52.99	40.85	48.94	51.47	60.35	34.92	50.86
XML-T	55.08	70.55	45.85	42.27	44.15	39.22	54.86	41.01	49.22	52.45	59.97	33.27	51.74
<b>TwHIN-BERT</b>													
- Base-MLM	58.38	72.66	48.41	43.08	46.89	41.53	56.76	42.36	49.60	51.13	61.00	35.37	53.66
- Base	59.31	<b>73.03</b>	48.59	44.24	<b>47.59</b>	42.81	57.33	42.69	51.11	56.66	60.33	36.21	54.62
- Large	<b>60.07</b>	72.91	<b>49.88</b>	<b>45.41</b>	47.43	<b>43.39</b>	<b>59.43</b>	<b>44.80</b>	<b>51.34</b>	<b>57.03</b>	<b>61.56</b>	<b>38.24</b>	<b>55.23</b>

**Table 4: External classification benchmark results.**

Method	SE2017	SE2018		ASAD	COVID-JA	SE2020		Avg.
	en	en	es	ar	ja	hi+en	es+en	
BERTweet	72.97	33.27	-	-	-	-	-	-
mBERT	66.17	27.73	19.19	69.08	80.57	66.55	45.31	53.51
XML-R	71.15	30.94	21.05	79.09	81.67	69.59	48.97	57.49
XML-T	72.01	31.97	21.49	80.70	81.48	70.94	51.06	58.52
<b>TwHIN-BERT</b>								
- Base-MLM	72.10	32.44	21.79	80.48	82.12	72.42	51.67	59.00
- Base	72.30	32.41	22.23	80.73	82.37	71.30	54.32	59.38
- Large	<b>73.10</b>	<b>33.31</b>	<b>22.80</b>	<b>81.19</b>	<b>82.50</b>	<b>73.08</b>	<b>54.47</b>	<b>60.06</b>

- **SemEval2020** task 9 [33] contains code-mixed Tweets of Hindi + English and Spanish + English. We use the three-point sentiment analysis part of the dataset for evaluation.
- **SemEval2018** task 2 [3] is an emoji prediction dataset in both English and Spanish. The objective is to predict the most likely used emoji in a Tweet.
- **COVID-JA** [43] is a Japanese Tweets classification dataset. The objective is to classify each Tweet into one of the six pre-defined topics around COVID-19.

**Setup and Evaluation Metrics.** We use the standard language model fine-tuning method as described in [14] and apply a linear prediction layer on top of the pooled output of the last transformer layer. Each model is fine-tuned for up to 30 epochs, and we evaluate the best model from the training epochs on the test set based on the development set performance. The fine-tuning hyperparameter setup can be found in Appendix B. We report the median results from 3 fine-tuning runs with different random seeds. Results are the evaluation metrics recommended for each benchmark dataset or challenge (Appendix C). For hashtag prediction datasets, we report macro-F1 scores. We conduct data contamination tests using character-level 50-gram overlaps [36] and found 1.56% and 1.10% contamination in the average results reported in Table 3 and Table 4.

**Multilingual Hashtag Prediction.** In Table 3, we show macro F1 scores on selected languages from our multilingual hashtag prediction dataset. We also report the average performance of all 50 languages in the dataset, and leave detailed results in Table 7. We can see that TwHIN-BERT significantly outperforms the baseline

methods at the same *base* size. Our *large* model is slightly better than or on par with the *base* model, with a better overall performance. On the English dataset, our model outperforms the BERTweet monolingual language model trained exclusively on English Tweets and with a dedicated English tokenizer. Comparing our model to the ablation with no social loss, the two models demonstrate similar performance with our model being slightly better. These results show that while our model has a major advantage on social tasks, it retains high performance on semantic understanding applications.

**External Classification Benchmarks.** As shown in Table 4, our TwHIN-BERT matches or outperforms the multilingual baselines on the established classification benchmarks. BERTweet fares better than our *base* model with its dedicated large English tokenizer and monolingual training. Our *large* model outperforms all the baselines. We note that it is not uncommon for a monolingual PLM to perform better than its multilingual counterpart, as observed in [13, 39, 50]. Similar to hashtag prediction, TwHIN-BERT performs on par with or slightly better than the MLM-only ablation.

### 3.4 Varying Downstream Supervision

In this set of experiments, we study how TwHIN-BERT performs when the amount of downstream supervision changes. We fine-tune our model and baseline models on the hashtag prediction dataset (Section 3.3). We select English and Japanese as they are the most popular languages on Twitter. We change the number of training examples given to the models during fine-tuning. It is varied from 2 to 32 labeled training examples per class. We follow the same protocols as Section 3.3 and report macro F1 scores on the test set.

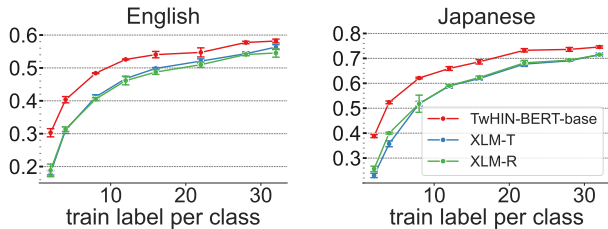


Figure 4: Macro-F1 on English and Japanese hashtag prediction datasets w.r.t. the number of labeled training examples per class.

Table 5: Feature-based classification on hashtag prediction datasets (Macro-F1).

Method	English	Japanese	Arabic
BERTweet	48.56	-	-
XLM-R	30.88	41.14	21.55
XLM-T	41.66	51.56	32.46
TwHIN-BERT-base	51.16	<b>64.12</b>	37.20
TwHIN-BERT-large	<b>54.12</b>	64.03	<b>38.78</b>

Figure 4 shows the results. TwHIN-BERT holds significant performance gain across different amount of downstream supervision. Note that when supervision is scarce, e.g., two labeled training examples per class given, our model has an even larger relative performance improvement over the baselines. The results indicate that our model may empower weakly supervised applications on Tweet natural language understanding.

### 3.5 Feature-based Classification

In addition to language model fine-tuning experiments, we evaluate TwHIN-BERT’s performance as a feature extractor. We use the hashtag prediction datasets (Section 3.3) and select three popular languages with different scripts. We use our model and the baseline models to embed each Tweet into a feature vector and train a Logistic Regression classifier with the fixed feature vectors as input.

Table 5 shows TwHIN-BERT outperforming the baselines with a wide margin on all languages. This not only shows TwHIN-BERT has learned superior Tweet representations but also showcases its potential in other feature-based downstream applications.

## 4 RELATED WORKS

**Pre-trained Language Models.** Since their introduction [14, 35], pre-trained language models have enjoyed tremendous success in all aspects of natural language processing. Follow-up research has advanced PLMs by further scaling them with respect to the number of parameters and training data. PLM models have grown considerably in their sizes, from millions [14, 51] of parameters to billions [6, 37, 41] and even trillion-level [19]. Another avenue of improvement has been improving the training objectives used to train PLMs. A broad spectrum of pre-training objectives have been explored with different levels of success. Notable examples include masked language modeling [14], auto-regressive causal language modeling [51], model-based denoising [11], and corrective language modeling [29]. Despite these innovations in scaling and

pre-training objectives, the vast majority of the work has focused on text-only training objectives applied to general domain corpora, e.g., Wikipedia and CommonCrawl. In this paper, we deviate from most previous works by exploring PLM training using solely Twitter in-domain data and training our model based on both text-based and social-based objectives.

**Tweet Language Models.** While a majority of PLMs are trained on general domain corpora, a few language models have been proposed specifically for Twitter and other social media platforms. BERTweet [32] mirrors BERT training on 850 million English Tweets. TimeLMs [28] trains a set of RoBERTa [27] models for English Tweets on different time ranges. XLM-T [2] continues the pre-training process from an XLM-R [12] checkpoint on 198 million multilingual Tweets. These methods mostly replicate existing general domain PLM methods and simply substitute the training data with Tweets. However, our approach utilizes additional social engagement signals to enhance the pre-trained Tweet representations.

**Enriching PLMs with Additional Information.** Several existing works use additional information for language model pre-training. ERNIE [54] and K-BERT [25] inject entities and their relations from knowledge graphs to augment the pre-training corpus. OAG-BERT [26] appends metadata of a document to its raw text, and designs objectives to jointly predict text and metadata. These works focus on bringing metadata and knowledge by injecting training instances, while our work leverages the rich social engagements embedded in the social media platform for text relevance. Recent work [52] has utilized document hyperlinks for LM pre-training, but does so with a simple three-way classification objective.

**Network Embedding.** Network embedding has emerged as a valuable tool for transferring information from relational data to other tasks [16]. Early network embedding methods such as DeepWalk [34] and node2vec [21] embed homogeneous graphs by performing random walks and applying SkipGram modeling. With the introduction of heterogeneous information networks [42] as a formalism to model rich multi-typed, multi-relational networks, many heterogeneous network embedding approaches were developed [7, 9, 15, 44, 49]. However, many of these techniques are difficult to scale to very large networks. In this work, we apply knowledge graph embeddings [5, 46, 48], which have been shown to be both highly scalable and flexible enough to model multiple node and edge types.

## 5 CONCLUSIONS

In this work we introduce TwHIN-BERT, a multilingual language model trained on a large Tweet corpus. Unlike previous BERT-style language models, TwHIN-BERT is trained using two objectives: (1) a standard MLM pre-training objective and (2) a contrasting social objective. We perform a variety of downstream tasks using TwHIN-BERT on Tweet data. Our experiments demonstrate that TwHIN-BERT outperforms previously released language models on both semantic and social engagement prediction tasks. We release TwHIN-BERT<sup>34</sup> to the academic community to further research in social media NLP.

<sup>3</sup><http://huggingface.co/Twitter/twhin-bert-base>

<sup>4</sup><http://huggingface.co/Twitter/twhin-bert-large>

**Table 6: Full social engagement prediction results (HITS@10) on all evaluation Languages.**

Language	TwHIN-BERT					
	mBERT	XLM-R	XLM-T	Base-MLM	Base	Large
English (en)	.0633	.0850	.1181	.1400	.1552	<b>.1585</b>
Japanese (ja)	.0227	.0947	.1079	.1413	.2065	<b>.2325</b>
Turkish (tr)	.0348	.0476	.1180	<b>.1268</b>	.1204	.0547
Spanish (es)	.0575	.0704	.1103	.1204	.1618	<b>.2055</b>
Arabic (ar)	.0532	.0546	.1403	.1640	<b>.2206</b>	.1989
Portuguese (pt)	.0731	.1285	.1709	.1201	<b>.1924</b>	.1915
Persian (fa)	.0556	.1621	.1754	.1903	.2065	<b>.2097</b>
Korean (ko)	.0275	.1105	.1446	.1675	.3611	<b>.3714</b>
French (fr)	.0488	.0635	.0805	.0700	.1030	<b>.1053</b>
Russian (ru)	.0889	.1482	.1530	.0990	<b>.1726</b>	.1704
German (de)	.0852	.1071	.3019	.2189	<b>.3020</b>	.2621
Thai (th)	.0659	.1027	.1056	.1196	<b>.2083</b>	.2004
Italian (it)	.0586	.0769	.1237	.1478	.1699	<b>.1706</b>
Hindi (hi)	.0870	.0838	.1140	.1054	.1737	<b>.1751</b>
Indonesian (id)	.0809	.0735	.0921	.1014	.1021	<b>.1115</b>
Polish (pl)	.0867	.0835	.1031	.1402	<b>.1696</b>	.1633
Urdu (ur)	.0437	.0315	.0352	.0547	.0627	<b>.0667</b>
Filipino (tl)	.0610	.0653	.0877	.1045	.1332	<b>.1400</b>
Egpt. Arabic (arz)	.0669	.0749	.1049	.0943	<b>.1159</b>	.1122
Greek (el)	.0496	.0628	.0562	.0801	.0944	<b>.1065</b>
Serbian (sr)	.1013	.1144	.1359	.1394	<b>.1647</b>	.1556
Dutch (nl)	.0616	.0650	.0762	.0965	<b>.1346</b>	.1248
Hebrew (he)	.0392	.0433	.0441	.0499	.0550	<b>.0577</b>
Ukrainian (uk)	.0497	<b>.0842</b>	.0669	.0711	.0811	<b>.0842</b>
Catalan (ca)	.1339	.1364	.1650	.1930	<b>.1955</b>	.1713
Swedish (sv)	.0942	.0716	.1161	.1342	<b>.1467</b>	.1462
Tamil (ta)	.0556	.0691	.0929	.1005	.1037	<b>.1060</b>
Finnish (fi)	.0876	.1067	.1317	.1529	.1710	<b>.1809</b>
Czech (cs)	.1155	.0904	.0766	.0997	.1062	<b>.1308</b>
Nepali (ne)	.0421	.0555	.0486	.0589	.0787	<b>.0851</b>
Azerbaijani (az)	.1561	.1148	.1702	.1576	.1712	<b>.1839</b>
Marathi (mr)	.0506	.0600	.0519	.0597	.0780	<b>.0906</b>
Bangla (bn)	.1361	.1350	.1320	.1601	.1649	<b>.1675</b>
Norwegian (no)	.0731	.1661	.1156	.1502	.1920	<b>.2118</b>
Telugu (te)	.0279	.0505	.0728	.0883	.1017	<b>.1654</b>
Pashto (ps)	.0522	.0727	.0662	.0600	.0799	<b>.0817</b>
Danish (da)	.1060	.1150	.1167	.1334	.1470	<b>.1475</b>
Vietnamese (vi)	.0929	.1060	.1085	.1216	.1417	<b>.1809</b>
Cen. Kurdish (ckb)	.0725	.0699	.0946	.1023	.1023	<b>.1185</b>
Gujarati (gu)	.0666	.0676	.0676	.0793	.1054	<b>.1057</b>
Macedonian (mk)	.0685	.0945	.0534	.0973	.1089	.1041
Cebuano (ceb)	.1222	.1267	.1767	.1900	.2003	<b>.2334</b>
Romanian (ro)	.1718	.1493	.1991	.2071	<b>.2264</b>	<b>.2264</b>
Kannada (kn)	.0552	.1355	.0814	.1098	.1282	<b>.2113</b>
Latvian (lv)	.0480	.0297	.0493	.0642	.0655	<b>.0750</b>
Bulgarian (bg)	.1953	.0448	.0702	.1790	.2248	<b>.2269</b>
Sinhala (si)	.0504	.0142	.0378	.0630	<b>.0709</b>	.0661
Icelandic (is)	.0319	.0341	.0466	.0364	.0387	<b>.0603</b>
Sindhi (sd)	.0619	.0288	.0553	.0885	.0951	<b>.0973</b>
Amharic (am)	.0293	.0663	.0491	.0543	.0698	<b>.0818</b>
Average	.0732	.0849	.1043	.1161	.1436	<b>.1497</b>

**Table 7: Full hashtag prediction results (Macro-F1) on all evaluation languages.**

Language	TwHIN-BERT					
	mBERT	XLM-R	XLM-T	Base-MLM	Base	Large
English (en)	54.56	53.90	55.08	58.38	59.31	<b>60.07</b>
Japanese (ja)	68.43	69.07	70.55	72.66	<b>73.03</b>	72.91
Turkish (tr)	42.87	46.37	47.14	48.72	49.31	<b>51.12</b>
Spanish (es)	42.48	43.80	45.85	48.41	48.59	<b>49.88</b>
Arabic (ar)	38.48	37.85	42.27	43.08	44.24	<b>45.41</b>
Portuguese (pt)	47.81	50.33	51.98	52.15	52.98	<b>56.08</b>
Persian (fa)	43.39	45.04	45.25	46.02	47.46	<b>47.94</b>
Korean (ko)	75.46	77.73	78.45	79.49	79.11	<b>80.02</b>
French (fr)	40.37	40.81	41.89	44.43	45.40	<b>47.01</b>
German (de)	40.80	41.42	41.11	41.32	41.38	<b>42.59</b>
Thai (th)	44.10	56.27	57.40	58.25	58.80	<b>59.46</b>
Italian (it)	42.36	41.82	42.76	45.11	44.18	<b>45.72</b>
Hindi (hi)	49.84	51.92	52.58	55.17	55.28	<b>57.29</b>
Chinese (zh)	72.88	72.54	72.40	73.85	<b>73.94</b>	72.30
Polish (pl)	48.97	50.20	50.50	51.20	51.81	<b>54.49</b>
Urdu (ur)	36.44	37.56	39.22	41.53	42.81	<b>43.39</b>
Filipino (tl)	52.96	52.99	54.86	56.76	57.33	<b>59.43</b>
Greek (el)	44.00	43.94	44.15	46.89	<b>47.59</b>	47.43
Serbian (sr)	42.50	42.32	40.71	44.22	45.95	<b>47.45</b>
Dutch (nl)	39.75	40.85	41.01	42.36	42.69	<b>44.80</b>
Catalan (ca)	48.61	47.85	48.79	51.72	52.60	<b>52.90</b>
Swedish (sv)	47.79	47.80	47.31	49.39	51.28	<b>51.44</b>
Tamil (ta)	48.04	49.67	50.65	52.85	54.14	<b>54.92</b>
Finnish (fi)	45.28	45.28	44.03	43.98	45.59	<b>46.42</b>
Czech (cs)	53.03	52.60	52.89	55.01	55.93	<b>56.02</b>
Nepali (ne)	44.58	47.00	46.94	49.83	<b>51.57</b>	51.12
Marathi (mr)	50.85	48.40	51.44	54.18	<b>55.76</b>	55.31
Malayalam (ml)	38.43	42.20	42.77	44.72	<b>45.86</b>	44.36
Bangla (bn)	57.79	57.08	56.74	59.11	60.32	<b>60.92</b>
Hungarian (hu)	60.29	59.94	60.08	61.86	<b>63.81</b>	62.68
Slovenian (sl)	58.79	59.68	59.13	61.18	62.34	<b>62.74</b>
Norwegian (no)	46.09	48.94	49.22	49.60	51.11	<b>51.34</b>
Telugu (te)	49.54	51.47	52.45	55.13	56.66	<b>57.03</b>
Pashto (ps)	29.41	34.92	33.27	35.37	36.21	<b>38.24</b>
Danish (da)	59.54	60.35	59.97	61.00	60.33	<b>61.56</b>
Cen. Kurdish (ckb)	40.28	37.59	39.06	42.89	<b>45.65</b>	45.26
Gujarati (gu)	52.55	54.09	54.24	57.36	<b>58.59</b>	58.54
Romanian (ro)	71.24	71.53	72.34	73.17	73.25	<b>73.58</b>
Kannada (kn)	54.19	55.76	56.68	59.09	<b>61.34</b>	60.19
Estonian (et)	57.81	58.10	59.00	61.95	61.22	<b>62.61</b>
Latvian (lv)	58.03	55.18	57.53	58.43	59.52	<b>61.47</b>
Bulgarian (bg)	65.20	65.52	66.45	67.42	66.94	<b>68.35</b>
Sinhala (si)	37.71	40.17	42.77	45.72	<b>47.54</b>	47.06
Icelandic (is)	51.32	48.53	50.16	53.39	<b>55.53</b>	54.61
Sindhi (sd)	27.28	26.46	31.08	32.42	35.05	<b>35.28</b>
Basque (eu)	58.55	56.55	56.78	59.56	60.62	<b>61.04</b>
Amharic (am)	24.10	28.57	35.01	34.20	<b>37.47</b>	36.87
Lithuanian (lt)	71.31	69.50	69.65	72.26	72.43	<b>73.09</b>
Welsh (cy)	58.36	58.50	57.56	59.66	59.95	<b>60.24</b>
Haitian Creole (ht)	68.13	67.05	67.97	70.70	71.33	<b>71.41</b>
Average	50.05	50.86	51.74	53.66	54.62	<b>55.23</b>



## REFERENCES

- [1] Basma Alharbi, Hind Alamro, Manal Abdulaziz Alshehri, Zuhair Khayyat, Manal Kalkatawi, Inji Ibrahim Jaber, and Xiangliang Zhang. 2020. ASAD: A Twitter-based Benchmark Arabic Sentiment Analysis Dataset. *ArXiv abs/2011.00578* (2020).
- [2] Francesco Barbieri, Luis Espinosa Anke, and José Camacho-Collados. 2021. XLM-T: A Multilingual Language Model Toolkit for Twitter. *ArXiv abs/2104.12250* (2021).
- [3] Francesco Barbieri, Jose Camacho-Collados, Francesco Ronzano, Luis Espinosa-Anke, Miguel Ballesteros, Valerio Basile, Viviana Patti, and Horacio Saggion. 2018. SemEval 2018 Task 2: Multilingual Emoji Prediction. In *Proceedings of The 12th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, New Orleans, Louisiana, 24–33. <https://doi.org/10.18653/v1/S18-1003>
- [4] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics* 5 (2017), 135–146. [https://doi.org/10.1162/tacl\\_a\\_00051](https://doi.org/10.1162/tacl_a_00051)
- [5] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Duran, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. *Advances in neural information processing systems* 26 (2013).
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.
- [7] S. Chang, W. Han, J. Tang, G. Qi, C. Aggarwal, and T. Huang. 2015. Heterogeneous network embedding via deep architectures. In *SIGKDD*. 119–128.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A Simple Framework for Contrastive Learning of Visual Representations. In *Proceedings of the 37th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 119)*, Hal Daumé III and Aarti Singh (Eds.). PMLR, 1597–1607. <https://proceedings.mlr.press/v119/chen20j.html>
- [9] T. Chen and Y. Sun. 2017. Task-guided and path-augmented heterogeneous network embedding for author identification. In *WSDM*. 295–304.
- [10] Justin Cheng, Lada A. Adamic, P. Alex Dow, Jon M. Kleinberg, and Jure Leskovec. 2014. Can cascades be predicted?. In *23rd International World Wide Web Conference, WWW '14, Seoul, Republic of Korea, April 7-11, 2014*, Chin-Wan Chung, Andrei Z. Broder, Kyuseok Shim, and Torsten Suel (Eds.). ACM, 925–936. <https://doi.org/10.1145/2566486.2567997>
- [11] Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: Pre-training Text Encoders as Discriminators Rather than Generators. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. <https://openreview.net/forum?id=r1xMH1BtvB>
- [12] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised Cross-lingual Representation Learning at Scale. In *ACL*.
- [13] Alexis Conneau and Guillaume Lample. 2019. Cross-lingual Language Model Pretraining. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 7057–7067. <https://proceedings.neurips.cc/paper/2019/hash/c04c19c2c247dbf5f7ac4372c5b9af1-Abstract.html>
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, 4171–4186. <https://doi.org/10.18653/v1/n19-1423>
- [15] Y. Dong, N. Chawla, and A. Swami. 2017. metapath2vec: Scalable representation learning for heterogeneous networks. In *SIGKDD*. 135–144.
- [16] Ahmed El-Kishky, Michael Bronstein, Ying Xiao, and Aria Haghighi. 2022. Graph-based Representation Learning for Web-scale Recommender Systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4784–4785.
- [17] Ahmed El-Kishky, Thomas Markovich, Kenny Leung, Frank Portman, and Aria Haghighi. 2022. kNN-Embed: Locally Smoothed Embedding Mixtures For Multi-interest Candidate Retrieval. *arXiv preprint arXiv:2205.06205* (2022).
- [18] Ahmed El-Kishky, Thomas Markovich, Serim Park, Chetan Verma, Baekjin Kim, Ramy Eskander, Yury Malkov, Frank Portman, Sofia Samaniego, Ying Xiao, and Aria Haghighi. 2022. TwHIN: Embedding the Twitter Heterogeneous Information Network for Personalized Recommendation. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, Aidong Zhang and Huzefa Rangwala (Eds.). ACM, 2842–2850. <https://doi.org/10.1145/3534678.3539080>
- [19] William Fedus, Barret Zoph, and Noam Shazeer. 2021. Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. *CoRR abs/2101.03961* (2021). [arXiv:2101.03961](https://arxiv.org/abs/2101.03961) <https://arxiv.org/abs/2101.03961>
- [20] Y. Goldberg and O. Levy. 2014. word2vec Explained: deriving Mikolov et al.'s negative-sampling word-embedding method. *arXiv preprint arXiv:1402.3722* (2014).
- [21] A. Grover and J. Leskovec. 2016. node2vec: Scalable feature learning for networks. In *SIGKDD*. 855–864.
- [22] Herve Jegou, Matthijs Douze, and Cordelia Schmid. 2010. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence* 33, 1 (2010), 117–128.
- [23] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [24] Adam Lerer, Ledell Wu, Jiajun Shen, Timothee Lacroix, Luca Wehrstedt, Abhijit Bose, and Alex Peysakhovich. 2019. Pytorch-bigggraph: A large-scale graph embedding system. *arXiv preprint arXiv:1903.12287* (2019).
- [25] Weijie Liu, Peng Zhou, Zhe Zhao, Zhiruo Wang, Qi Ju, Haotang Deng, and Ping Wang. 2020. K-BERT: Enabling Language Representation with Knowledge Graph. In *AAAI*.
- [26] Xiao Liu, Da Yin, Jingnan Zheng, Xingjian Zhang, P. Zhang, Hongxia Yang, Yuxiao Dong, and Jie Tang. 2022. OAG-BERT: Towards a Unified Backbone Language Model for Academic Knowledge Services. *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (2022).
- [27] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR abs/1907.11692* (2019). <http://arxiv.org/abs/1907.11692>
- [28] Daniel Loureiro, Francesco Barbieri, Leonardo Neves, Luis Espinosa Anke, and José Camacho-Collados. 2022. TimeLMs: Diachronic Language Models from Twitter. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics, ACL 2022 - System Demonstrations, Dublin, Ireland, May 22-27, 2022*, Valerio Basile, Zornitsa Kozareva, and Sanja Stajner (Eds.). Association for Computational Linguistics, 251–260. <https://doi.org/10.18653/v1/2022.acl-demo.25>
- [29] Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia Song. 2021. COCO-LM: Correcting and Contrasting Text Sequences for Language Model Pretraining. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 23102–23114. <https://proceedings.neurips.cc/paper/2021/hash/c2c2a04512b35d13102459f8784f1a2d-Abstract.html>
- [30] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. 2013. Distributed representations of words and phrases and their compositionality. *NeurIPS* 26 (2013).
- [31] Fatemehsadat Mirehghallah, Nikolai Vogler, Junxian He, Omar Florez, Ahmed El-Kishky, and Taylor Berg-Kirkpatrick. 2022. Non-Parametric Temporal Adaptation for Social Media Topic Classification. *arXiv preprint arXiv:2209.05706* (2022).
- [32] Dat Quoc Nguyen, Thanh Vu, and Anh Tuan Nguyen. 2020. BERTweet: A pre-trained language model for English Tweets. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, EMNLP 2020 - Demos, Online, November 16-20, 2020*, Qun Liu and David Schlangen (Eds.). Association for Computational Linguistics, 9–14. <https://doi.org/10.18653/v1/2020.emnlp-demos.2>
- [33] Parth Patwa, Gustavo Aguilar, Sudipta Kar, Suraj Pandey, Srinivas PYKL, Björn Gambäck, Tanmoy Chakraborty, Thamar Solorio, and Amitava Das. 2020. SemEval-2020 Task 9: Overview of Sentiment Analysis of Code-Mixed Tweets. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation, SemEval@COLING 2020, Barcelona (online), December 12-13, 2020*, Aurélie Herbelot, Xiaodan Zhu, Alexis Palmer, Nathan Schneider, Jonathan May, and Ekaterina Shutova (Eds.). International Committee for Computational Linguistics, 774–790. <https://doi.org/10.18653/v1/2020.semeval-1.100>
- [34] B. Perozzi, R. Al-Rfou, and S. Skiena. 2014. Deepwalk: Online learning of social representations. In *SIGKDD*. 701–710.
- [35] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2018, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 1 (Long Papers)*, Marilyn A. Walker, Heng Ji, and Amanda Stent (Eds.). Association for Computational Linguistics, 2227–2237. <https://doi.org/10.18653/v1/n18-1202>
- [36] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [37] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *J. Mach. Learn. Res.* 21 (2020), 140:1–140:67. <http://jmlr.org/papers/v21/20-074.html>

- [38] Sara Rosenthal, Noura Farra, and Preslav Nakov. 2019. SemEval-2017 Task 4: Sentiment Analysis in Twitter. *CoRR* abs/1912.00741 (2019). arXiv:1912.00741 <http://arxiv.org/abs/1912.00741>
- [39] Phillip Rust, Jonas Pfeiffer, Ivan Vulic, Sebastian Ruder, and Iryna Gurevych. 2021. How Good is Your Tokenizer? On the Monolingual Performance of Multilingual Language Models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, 3118–3135. <https://doi.org/10.18653/v1/2021.acl-long.243>
- [40] Aravind Sankar, Xinyang Zhang, Adit Krishnan, and Jiawei Han. 2020. Inf-VAE: A Variational Autoencoder Framework to Integrate Homophily and Influence in Diffusion Prediction. In *WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020*, James Caverlee, Xia (Ben) Hu, Mounia Lalmas, and Wei Wang (Eds.). ACM, 510–518. <https://doi.org/10.1145/3336191.3371811>
- [41] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism. *CoRR* abs/1909.08053 (2019). arXiv:1909.08053 <http://arxiv.org/abs/1909.08053>
- [42] Y. Sun and J. Han. 2013. Mining heterogeneous information networks: a structural analysis approach. *Acm Sigkdd Explorations Newsletter* (2013).
- [43] Yu Suzuki. 2019. Filtering Method for Twitter Streaming Data Using Human-in-the-Loop Machine Learning. *J. Inf. Process.* 27 (2019), 404–410. <https://doi.org/10.2197/ipsjip.27.404>
- [44] J. Tang, M. Qu, and Q. Mei. 2015. Pte: Predictive text embedding through large-scale heterogeneous text networks. In *SIGKDD*. 1165–1174.
- [45] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. 2015. LINE: Large-scale Information Network Embedding. In *Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015*, Aldo Gangemi, Stefano Leonardi, and Alessandro Panconesi (Eds.). ACM, 1067–1077. <https://doi.org/10.1145/2736277.2741093>
- [46] T. Trouillon, J. Welbl, S. Riedel, É. Gaussier, and G. Bouchard. 2016. Complex embeddings for simple link prediction. In *ICML*. PMLR, 2071–2080.
- [47] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [48] Q. Wang, Z. Mao, B. Wang, and L. Guo. 2017. Knowledge graph embedding: A survey of approaches and applications. *TKDE* 29, 12 (2017), 2724–2743.
- [49] L. Xu, X. Wei, J. Cao, and P. Yu. 2017. Embedding of embedding (EOE) joint embedding for coupled heterogeneous networks. In *WSDM*. 741–749.
- [50] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A Massively Multilingual Pre-trained Text-to-Text Transformer. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tür, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (Eds.). Association for Computational Linguistics, 483–498. <https://doi.org/10.18653/v1/2021.naacl-main.41>
- [51] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime G. Carbonell, Ruslan Salakhutdinov, and Quoc V. Le. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 5754–5764. <https://proceedings.neurips.cc/paper/2019/hash/dc6a7e655d7e5840e66733e9ee67cc69-Abstract.html>
- [52] Michihiro Yasunaga, Jure Leskovec, and Percy Liang. 2022. LinkBERT: Pretraining Language Models with Document Links. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 8003–8016. <https://doi.org/10.18653/v1/2022.acl-long.551>
- [53] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*, Yike Guo and Faisal Farooq (Eds.). ACM, 974–983. <https://doi.org/10.1145/3219819.3219890>
- [54] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, Anna Korhonen, David R. Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, 1441–1451. <https://doi.org/10.18653/v1/p19-1139>

## A DISTRIBUTION OF LANGUAGES IN TRAINING DATASET

Figure 5 shows the distribution of languages in our pre-training dataset. Some languages with different variations (e.g., Hindi and Hindi Romanized) are represented with the same ISO language code. We run fastText [4] language identification model `lid.176.bin`<sup>5</sup> to detect languages.

We deem a language “high-resource” if we have more than  $10^8$  Tweets during pre-training *after* frequency-based re-sampling (Section 2.3); “mid-resource” if we have more than  $10^7$  and less than  $10^8$  Tweets; “low-resource” if we have less than  $10^7$  Tweets.

## B HYPERPARAMETERS FOR PRE-TRAINING AND FINE-TUNING

Table 8 shows the pre-training hyperparameters. The model architecture and hyperparameters not shown in the table are the same as RoBERTa [27].

Table 9 shows the hyperparameters for classification fine-tuning. We do hyperparameter selection on the development datasets and share the same set of hyperparameters for the base models, as we find them to perform well with this setting. The weight decay for base models is set to zero. A different set of hyperparameters were necessary for the large model because it behaves differently from the base models in terms of convergence.

## C EVALUATION METRICS FOR EXTERNAL CLASSIFICATION BENCHMARKS

The recommended evaluation metrics that we report in Table 4 are as follows. Average recall for ASAD, SemEval 2017 datasets; Macro-F1 for SemEval 2018 English and Spanish datasets; Accuracy for COVID-JA, SemEval 2020 datasets.

## D ENGAGEMENT PREDICTION RESULTS ON ADDITIONAL LANGUAGES

Table 6 shows the engagement prediction results on all available evaluation languages. Some languages have more examples than other languages due to data availability.

## E HASHTAG PREDICTION RESULTS ON ADDITIONAL LANGUAGES

Table 7 shows the hashtag prediction results on all available evaluation languages. A small number of languages have less examples than shown in Table 2 due to data availability. The Russian language is not evaluated as the XLM-T baseline fails on some Russian characters in our dataset.

<sup>5</sup><https://fasttext.cc/docs/en/language-identification.html>

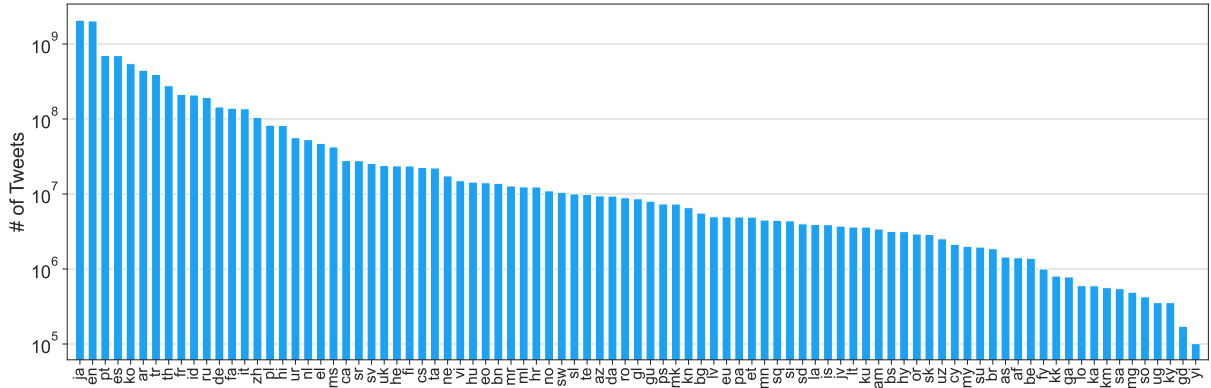


Figure 5: The number of Tweets in the pre-training dataset for each language. Languages are marked by ISO language codes.

Table 8: Hyperparameters for pre-training TwHIN-BERT.

Hyperparameter	TwHIN-BERT-base	TwHIN-BERT-large
Max sequence length	128	128
Precision	BF16	BF16
<b>Stage 1: MLM</b>		
Total batch size	6K	8K
Gradient accumulation steps	1	4
Peak learning rate	$2e-4$	$2e-4$
Warmup steps	30K	30K
Total steps	500K	500K
<b>Stage 2: MLM + Social</b>		
Total batch size	6K	6K
Gradient checkpointing	No	Yes
Peak learning rate	$1e-4$	$1e-4$
Warmup steps	30K	30K
Total steps	500K	500K
Contrastive projection head	[768, 768]	[1024, 512]
Contrastive loss temperature	0.1	0.1
Loss balancing $\lambda$	0.05	0.05

Table 9: Hyperparameters for fine-tuning TwHIN-BERT and the baselines for classification.

Hyperparameter	Hashtag	SE2017	SE2018	ASAD	COVID-JA	SE2020
<b>Base models</b>						
Learning rate	$4e-5$	$4e-5$	$1e-5$	$1e-5$	$2e-5$	$2e-5$
Batch size	128	128	128	128	128	128
<b>TwHIN-BERT-large</b>						
Learning rate	$2e-5$	$2e-5$	$1e-5$	$1e-5$	$1e-5$	$1e-5$
Weight decay	0	0	$5e-4$	$5e-4$	0	$5e-4$
Batch size	128	128	128	128	128	128