

Hashtag-Guided Low-Resource Tweet Classification

Shizhe Diao*
The Hong Kong University of Science
and Technology
sdiaooa@connect.ust.hk

Sedrick Scott Keh*
Carnegie Mellon University
skeh@cs.cmu.edu

Liangming Pan
University of California, Santa
Barbara
liangmingpan@ucsb.edu

Zhiliang Tian
The Hong Kong University of Science
and Technology
tianzhilianghit@gmail.com

Yan Song
University of Science and Technology
of China
clksong@gmail.com

Tong Zhang
The Hong Kong University of Science
and Technology
tongzhang@ust.hk

Abstract

Social media classification tasks (*e.g.*, tweet sentiment analysis, tweet stance detection) are challenging because social media posts are typically short, informal, and ambiguous. Thus, training on tweets is challenging and demands large-scale human-annotated labels, which are time-consuming and costly to obtain. In this paper, we find that providing hashtags to social media tweets can help alleviate this issue because hashtags can enrich short and ambiguous tweets in terms of various information, such as topic, sentiment, and stance. This motivates us to propose a novel **Hashtag-guided Tweet Classification** model (**HASHTATION**), which automatically generates meaningful hashtags for the input tweet to provide useful auxiliary signals for tweet classification. To generate high-quality and insightful hashtags, our hashtag generation model retrieves and encodes the post-level and entity-level information across the whole corpus. Experiments show that HASHTATION achieves significant improvements on seven low-resource tweet classification tasks, in which only a limited amount of training data is provided, showing that automatically enriching tweets with model-generated hashtags could significantly reduce the demand for large-scale human-labeled data. Further analysis demonstrates that HASHTATION is able to generate high-quality hashtags that are consistent with the tweets and their labels. The code is available at <https://github.com/shizhediao/HashTation>.

CCS Concepts

• **Information systems** → **Web mining**.

Keywords

social media analysis, tweet classification, hashtag generation, low-resource classification

*Both authors contributed equally to this research.

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.
WWW '23, May 1–5, 2023, Austin, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9416-1/23/04...\$15.00
<https://doi.org/10.1145/3543507.3583194>

ACM Reference Format:

Shizhe Diao, Sedrick Scott Keh, Liangming Pan, Zhiliang Tian, Yan Song, and Tong Zhang. 2023. Hashtag-Guided Low-Resource Tweet Classification. In *Proceedings of the ACM Web Conference 2023 (WWW '23)*, May 1–5, 2023, Austin, TX, USA. ACM, New York, NY, USA, 13 pages. <https://doi.org/10.1145/3543507.3583194>

1 Introduction

Table 1: Examples of how hashtags can provide auxiliary information for better tweet classification.

Original Input	Generated Hashtags	New Hashtag-Guided Input
Abortion IS NOT a political issue. It is a MORAL issue.	AllLivesMatter, ProLife	Abortion IS NOT a political issue. It is a MORAL issue #AllLivesMatter #ProLife
Twitter making everybody mad. It's hilarious	Twitter, hilarious	#Twitter making everybody mad. It's #hilarious
He's the GOAT for sure!	GOAT, NBAFinals	He's the #GOAT for sure! #NBAFinals

Tweet Classification (TC) is an essential task in social media content analysis, which aims to analyze user behaviors and attitudes on Twitter. Typical tasks in this area include stance detection [36], sentiment analysis [37], and hate speech detection [4]. However, these tasks are often difficult because tweets are usually informal, idiosyncratic, and short in length and thus provide limited and ambiguous information. Additional context and background knowledge are often needed to understand the content of a tweet better. Due to this lack of information and the ambiguous nature of tweets, we often need to label a large-scale training corpus in order to train a satisfactory TC model [3, 37]. However, the rapidly changing and evolving nature of social media content makes it challenging to annotate in-domain training data in a timely manner. Furthermore, data annotation is time-consuming and costly. To address the above challenges, we propose a model, HASHTATION, with two novel features: 1) it can automatically enrich the content of social media tweets by *hashtag generation*, and 2) the hashtag-enriched tweet classification model works well under the *low-resource* setting in which only a limited amount of labeled data is available.

Hashtags are commonly contained within tweets or appended to the end of tweets. They not only facilitate rapid lookup for specific themes or web contents, but also contain important information that helps to enrich and disambiguate the contents of tweets. As exemplified in Table 1, we can hardly understand the topic and sentiment of the tweet “Abortion IS NOT a political issue. It is a MORAL

Table 2: Pilot studies on TC. First, we split all the training data and validation data into two sets: tweets containing hashtags (*w.* hashtags) and tweets without hashtags (*w.o.* hashtags). We separately train and evaluate the model’s performance on these two sets over 10 different random seeds. F1 scores on dev. set are reported.

	EMOJI	EMOTION	HATE	IRONY	OFFENSIVE	SENTIMENT	STANCE	AVG.
<i>w.o.</i> HASHTAGS	8.26	55.9	63.7	61.7	64.1	50.2	55.5	51.3
<i>w.</i> HASHTAGS	8.58	64.5	73.5	67.3	65.8	52.0	59.4	55.9

issue.” without its hashtag “#RoeVWade”. Our pilot study (Table 2) finds that hashtags provide important signals for tweet classification; the tweets without hashtags suffer from an average F1-score drop of 4.6% in seven different tasks of TC compared with the tweets containing hashtags. Through this preliminary experiment, we find evidence to support our intuition that hashtags enrich the typically short and ambiguous tweets.

Unfortunately, despite the usefulness of hashtags, a vast majority of tweets do not use them (See Table 3). Motivated by this, our model HASHTATION addresses low-resource tweet classification in two steps: 1) we propose a novel hashtag generation model by collecting hashtag-containing tweets as our training data, and 2) we enrich the hashtag-absent tweets using the generated hashtags from HASHTATION. To generate more accurate and insightful hashtags, we leverage global contexts to consider not only the input tweet but also other relevant tweets, as well as tweets which share the same entities. To do this, we propose two corpus-level attention modules (Tweet Attention Module and Entity Attention Module) to retrieve helpful information across the whole corpus and entity graph. After obtaining generated hashtags, we compose a hashtag-prompted input by combining the hashtags and the original tweet with human-designed templates. We then feed the hashtag-prompted input into a pre-trained TC model.

Experiments on seven diverse TC tasks show that automatically enriching tweets with model-generated hashtags could significantly reduce the demand for large-scale human-labeled data. Further analysis shows that HASHTATION generates high-quality hashtags that are consistent with the tweets and their labels, revealing the effectiveness of leveraging global contexts by retrieval.

Our main contributions are as follows:

- We empirically reveal that hashtags provide crucial signals for classifying social media posts in seven different tasks.
- In light of the empirical evidence, we propose a novel hashtag generation model that leverages multi-grained global contexts to consider not only the input post but also the relevant posts and entities by using a cross-attention module and a graph entity network.
- By enriching social media posts with generated hashtags, we significantly reduce the demand for human annotation for tweet classification. Our model significantly boosts low-resource tweet classification performance in seven different tasks.

2 Approach

Figure 1 shows the architecture and detailed components of HASHTATION. They are as follows: 1) *Hashtag Generator*, which encodes the input tweets and decodes the hashtags via self-attention; 2) *Tweet Attention Module* (TAM), employing a cross-document attention

network to retrieve and incorporate the relevant semantic information at the tweet-level implicitly, which is then fused with the input representation; 3) *Entity Attention Module* (EAM), which employs a graph attention network to retrieve and leverage the relevant semantic information at the entity-level explicitly, which is then fused with the input representation; 4) *Tweet Classifier*, which adopts a Transformer encoder with a classification head to perform tweet classification. Below, we formulate the problem and then describe the details of each module.

2.1 Problem Formulation

In our setting, the basic input unit is a tweet $x = \{w_1 \cdots w_n\}$ where w_i is the i -th word of the given tweet. Output $H = \{h_1 \cdots h_m\}$ is a sequence of all the hashtags of x , and y is the target class/label corresponding to the input tweet. The first stage of our model, HASH-GEN, for hashtag generation can be formulated as follows.

$$H = \text{HASH-GEN}(x, \text{TAM}(x, \mathcal{D}), \text{EAM}(x, \mathcal{D})), \quad (1)$$

where HASH-GEN refers to the hashtag generator, TAM refers to the Tweet Attention Module that produces latent topic embeddings (Section 2.2), and EAM is the Entity Attention Module that produces explicit relevant entity information (Section 2.3). Both TAM and EAM exploit a collection of tweets \mathcal{D} . The hashtag generation is then enhanced with the latent topics provided in \mathcal{D} and the explicit entity relations provided in the EAM.

The second stage of our model, TWEET-CLASSIFIER, for tweet classification is formulated as follows.

$$y = \text{TWEET-CLASSIFIER}(\text{Fuse}(H, x)), \quad (2)$$

where TWEET-CLASSIFIER is a Transformer-based text classifier. Details of the HASH-GEN and TWEET-CLASSIFIER, as well as how we integrate them, are described in the following subsections.

2.2 Tweet Attention Module (TAM)

Given an input tweet x , relevant tweets usually share similar topics, which are good references to help determine what could be the optimal hashtags to describe x . For example, for a tweet about school closures, the hashtag ‘#CloseTheSchools’ may be absent in the tweet but may appear in other relevant tweets, providing helpful guidance for hashtag generation in this scenario.

To represent and exploit latent topics from relevant tweets, we first aggregate all tweets from a collection (*i.e.*, the union of both training and validation set), represented as $\mathcal{D} = \{x_1, \dots, x_k, \dots, x_l\}$. Following the Transformer architecture, we define key vectors $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_k, \dots, \mathbf{u}_l\}$ and value vectors $\mathcal{V} = \{\mathbf{v}_1, \dots, \mathbf{v}_k, \dots, \mathbf{v}_l\}$ with \mathbf{u}_k and \mathbf{v}_k corresponding to x_k . Specifically, \mathbf{u}_k is used to compute similarity with the input while \mathbf{v}_k carries x_k ’s encoding information for generating the final output, which acts as the latent topic embedding.

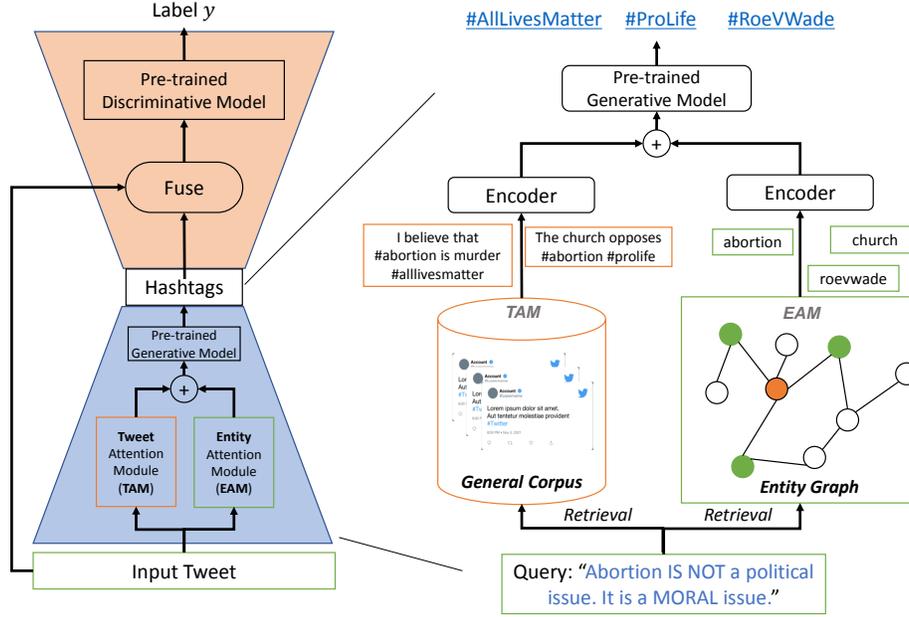


Figure 1: Left: Overview of HASHTATION. There are two main components: the bottom module is the hashtag prompt generator and the top module is the classification model. Right: Illustration of tweet attention module (TAM) and entity attention module (EAM). The hashtag generator encodes the input tweets and decodes the hashtags via self-attention. The TAM retrieves and incorporates the relevant semantic information at the tweet-level implicitly. The EAM retrieves and leverages the relevant semantic information at the entity-level explicitly. The pre-trained discriminative model is adopted with a classification head to perform tweet classification.

Then for each input x , we represent it through its sentential encoding e and use it as the ‘query’ vector to address relevant tweets. In detail, the addressing operation can be formalized as

$$p_k = \frac{\exp(e^\top \cdot u_k)}{\sum_{k=1}^l \exp(e^\top \cdot u_k)}, \quad (3)$$

and for the entire set \mathcal{D} , we have $\mathbf{o} = \sum_{k=1}^l p_k \mathbf{v}_k$, where \mathbf{o} is the output vector of the TAM to represent the latent topics from relevant tweets via a weighted encoding.

2.3 Entity Attention Module (EAM)

Graph Construction. To better aggregate and leverage the entity information across all tweets, the entity attention module offers strong clues in determining what to predict next. We use a named entity recognition model to construct the entity graph by extracting relevant entities across all tweets based on co-occurrence. We first create an empty graph $\mathcal{G} = (\mathcal{N}, \mathcal{E})$, where $\mathcal{N} = \{n_1, \dots, n_i, \dots, n_N\}$ and $\mathcal{E} = \{e_1, \dots, e_j, \dots, e_M\}$ are the node and edge sets, respectively. For each tweet input x , we use our entity recognizer to extract a set of entities. Each entity will be added into \mathcal{G} as a node and linked to those nodes that are extracted from the same tweet. As a result, we obtain an entity graph with multiple nodes and edges, capturing the critical entity relations across relevant tweets.

Graph Encoding. The node embeddings are randomly initialized as $H^{(0)}$ during graph construction. In order to equip them with better semantic relation information, we adopt the graph convolutional networks [22] to provide a better initialization. We then adopt the idea

of graph attention networks (GAT) [46] to obtain the node encoding by aggregating information from their neighbors and update it with dynamic weights. Formally, GAT accepts $H^{(0)}$ and outputs $H^{(L)}$ after conducting L layers of state transitions. Given a sequence of hidden representations $H^{(l)} = [\vec{h}_1^l, \dots, \vec{h}_i^l, \dots, \vec{h}_N^l]$ at layer l with \vec{h}_i^l indicating the feature representation for i -th node, and an adjacency matrix A , the hidden representation at layer $l+1$ is calculated by

$$H^{(l+1)} = \sigma \left(\hat{D}^{-\frac{1}{2}} \hat{A} \hat{D}^{-\frac{1}{2}} H^{(l)} W^{(l)} \right), \quad (4)$$

where $\hat{A} = A + I$, I is the identity matrix and \hat{D} is the diagonal node degree matrix of \hat{A} . $W^{(l)}$ is a weight matrix for the l -th neural network layer, and $\sigma(\cdot)$ is the *ReLU* activation function. Then, we calculate the attention coefficient between nodes n_i and n_j by

$$\alpha_{ij} = \frac{\exp \left([W^Q \vec{h}_i^l W^K \vec{h}_j^l] \right)}{\sum_{k \in \mathcal{M}_i} \exp \left([W^Q \vec{h}_i^l W^K \vec{h}_k^l] \right)}, \quad (5)$$

where W^Q and W^K are weight matrices for feature transformation of query and key, and \mathcal{M}_i is the first-order neighbors of node n_i (including n_i).

Third, after obtaining attention coefficients, the output encoding is a linear combination of the features in its neighbors, computed by

$$\vec{h}_i^{l+1} = \sigma \left(\sum_{n_j \in \mathcal{M}_i} \alpha_{ij} W \vec{h}_j^l \right), \quad (6)$$

where W is a weight matrix and $\sigma(\cdot)$ is the *ReLU* activation function.

Table 3: Dataset statistics of the original dataset and the set whose tweets contain hashtags.

Task	Labels	Original Dataset			With Hashtags		
		Train	Val	Test	Train	Val	Test
Emoji Prediction	20 different emojis	44489	4930	49664	19878	1800	20387
Emotion Recognition	anger, joy, sadness, optimism	3188	363	1384	1394	178	650
Hate Speech Detection	hateful, not hateful	8914	990	2963	2223	228	1394
Irony Detection	irony, not irony	2802	942	781	1017	368	530
Offensive Language Identification	offensive, not offensive	11616	1301	852	1664	194	630
Sentiment Analysis	positive, neutral, negative	45612	2000	12200	8262	347	4571
Stance Detection	in favor, neutral, against	2620	294	1249	2578	288	1249

2.4 Integrating Tweet-level and Entity-level Information

Although RNN-based sequence-to-sequence models are widely used for the hashtag generation task, we use the Transformer [45] as the backbone encoder-decoder framework in this paper. This has been proven to be more effective than RNN-based sequence-to-sequence models in many generation tasks [16, 19, 20, 26, 45]. Once the latent topic embedding \mathbf{o} and graph entity encoding $\tilde{\mathbf{h}}$ are obtained, we combine them with the Transformer encoding-decoding process via the following steps.

First, the tweet input is passed through the Transformer encoder, which results in a hidden state \mathbf{h}_i for each input token. Then we combine \mathbf{h}_i and \mathbf{o} via element-wise addition $\tilde{\mathbf{h}}_i = \mathbf{h}_i + \mathbf{o}$. Third, we enhance $\tilde{\mathbf{h}}_i$ with graph entity information $\tilde{\mathbf{h}}$ by $\tilde{\mathbf{h}}_i = \tilde{\mathbf{h}}_i + \sum_j \tilde{\mathbf{h}}_{i,j}$, where $\tilde{\mathbf{h}}_{i,j}$ is the graph encoding of the j -th entity node associated to the i -th token. $\tilde{\mathbf{h}}_i$ is the final encoding state of i -th token and is sent to the decoding process through each multi-head attention layer to calculate the attention vector $\mathbf{a}^t = \alpha_1^t \dots \alpha_i^t \dots \alpha_n^t$ at each decoding step t . Next, \mathbf{a}^t is used to produce the context vector $\mathbf{c}^t = \sum_{i=1}^n \alpha_i^t \tilde{\mathbf{h}}_i$ by a weighted sum of the encoding hidden states. Later \mathbf{c}^t is concatenated with the decoder output \mathbf{s}^t and then fed into a single linear layer, followed by a softmax function, to produce the vocabulary distribution for the output word at step t

$$\mathbf{d}_o = \frac{\exp(\mathbf{z}^t)}{\sum_V \exp(\mathbf{z}^t)}, \quad (7)$$

where $\mathbf{z}^t = \mathbf{W}_1(\mathbf{W}_2 \cdot (\mathbf{s}^t \oplus \mathbf{c}^t))$, a vector with $|V|$ dimension and V is the predefined vocabulary providing word candidates for hashtag generation. \mathbf{W}_1 and \mathbf{W}_2 are trainable parameters for the two aforementioned linear layers, respectively.

2.5 Tweet Classifier

After obtaining the desired hashtags, we fuse the hashtags $H = h_1, h_2, \dots, h_n$ with the input tweet x to obtain the hashtag-guided tweet $Fuse(H, x)$. There are several ways to implement the $Fuse$ function like simple concatenation, prompting with a manual template, and so on. We explore the effects of different strategies in Section 5.1. Finally, we train a tweet classifier \mathcal{F} with the following objective:

$$\min_{\phi} (\mathcal{L}(\mathcal{F}(Fuse(H, x)), Y)), \quad (8)$$

where \mathcal{L} is the cross-entropy loss function following the standard practice in tweet classification [30, 38].

3 Experimental Settings

In this section, we first introduce the datasets (Section 3.1), followed by the baseline models (Section 3.2) and evaluation metrics (Section 3.3). Lastly, we describe the implementation details (Section 3.4) for various experiments.

3.1 Datasets

We perform our experiments on 7 diverse tweet classification tasks in the TweetEval benchmark datasets [2]. All 7 tasks are taken from previous SemEval tasks (and corresponding datasets) and are as follows: emotion recognition [35], emoji prediction [3], irony detection [44], hate speech detection [4], offensive language identification [57], sentiment analysis [37], and stance detection [36]. Additional details about these tasks can be found in Table 3.

Creating the hashtag generation input and output. For hashtag generation, we train a separate hashtag generator for each task, as the datasets cover significantly different domains. To extract hashtags, we searched for appearances of the octothorpe symbol (#) and considered the contiguous string following it, until a whitespace is reached. These hashtags can either appear mid-tweet or at the end of the tweet. For hashtags that appear mid-tweet, we remove the hashtag symbol (#) but keep the word itself, as removing these words may disrupt the coherence of the sentence. On the other hand, for hashtags that appear at the end of the tweet, we completely remove these from the body of the tweet. If we keep these words at the end of the sentence, the hashtag generator will easily pick up on these false signals and just learn to return the last few words of every sentence without actually learning to perform hashtag generation. By considering this setting of keeping in-sentence hashtags and removing end-of-sentence hashtags, we are able to train the model to perform prediction for both present hashtags and absent hashtags. Examples of this overall process are outlined in Table 9. Similar to Wang et al. [50], we implement some preprocessing steps to clean up our tweets. Links, mentions, and numbers were replaced with “URL”, “MENTION”, and “DIGIT”, respectively. The details of the processed dataset are shown in Table 3.

Low-resource setting. As discussed in Section 1, real-life classification labels are often difficult to acquire for the rapidly changing landscape of tweets. As such, we conduct our experiments in a low-resource setting, where we randomly select a subset from the training set. To keep the sizes of different datasets relatively consistent, we use the following sampling ratios: If the size of the original dataset is < 5000 , we sample 10% of the dataset. If it’s between 5000

Table 4: F1-scores for baseline models (top) and our proposed method and its variants (bottom). Best results are highlighted in bold. We report the average score over ten different random seeds. * denotes HASHTATION has significant differences (p -value <0.05) over baseline models.

	EMOJI	EMOTION	HATE	IRONY	OFFENSIVE	SENTIMENT	STANCE	AVG.
KIM-CNN	3.8	20.1	51.1	39.8	41.9	39.5	32.2	32.6
BiLSTM	5.9	26.8	49.2	37.4	47.6	43.1	28.7	34.1
BERT-BASE	11.2	57.7	51.6	53.7	56.4	56.6	52.3	48.5
ROBERTA-BASE	11.8	58.9	56.7	54.2	59.7	57.3	54.1	50.4
BERTWEET	12.1	59.4	55.5	57.0	61.8	59.0	55.5	51.5
HASHTATION-BT	12.7	61.0	56.4	58.6	63.8*	59.9*	56.3	52.7
w/o TAM	12.5	60.6	56.3	57.7	63.0	59.8	56.2	52.3
w/o EAM	12.6	60.5	55.7	58.6	63.5	59.4	56.3	52.4
TIMELMS	12.4	60.2	56.9	59.6	60.0	57.4	55.9	51.8
HASHTATION-TL	13.0*	61.6*	58.0*	60.9*	62.3	59.6	56.5*	53.1*
w/o TAM	12.7	61.2	57.6	60.6	61.9	58.9	56.1	52.7
w/o EAM	12.7	61.3	57.6	60.7	61.6	59.0	56.2	52.7

and 10000, we sample 5%, and if it’s > 10000, we sample 1%. All experiments are performed over ten different random seeds.

3.2 Baselines

To verify the effectiveness of our HASHTATION model, we compare its performance against the following baseline classification models: **Kim-CNN** [21], **BiLSTM** [41], **BERT** [10], **RoBERTa** [27], **BERTweet** [39], and **TimeLMs** [30]. BERTweet and TimeLMs are two state-of-the-art models that are pre-trained on 850 million and 124 million English tweets, respectively.

- **Kim-CNN** [21]. A simple convolutional neural networks framework with a classification layer attached to the end. For this CNN-based model, we use kernel sizes of 2,3,4,5 and 64 filters for each, with a dropout of 0.5 before the linear layer.
- **BiLSTM** [41]. A bidirectional long short-term memory network considering the temporal order of words in the tweet. The hidden size and the dropout rate for the LSTM are set to 512 and 0.2, respectively. We use an LSTM hidden size of 512 and an LSTM dropout of 0.2. In addition, we use a dropout of 0.5 and a size of 32 for the final linear classification layer.
- **BERT** [10]. A BERT-base model pre-trained on generic corpus with a classification head fine-tuned on specific task.
- **RoBERTa** [27]. A RoBERTa-base model pre-trained on generic corpus with a classification head fine-tuned on specific task.
- **BERTweet** [39]. A RoBERTa-base model pre-trained on 850M English Tweets with a classification head fine-tuned on specific task.
- **TimeLMs** [30]. A RoBERTa-base model specialized on diachronic Twitter data and pre-trained on 124M English Tweets.
- **HASHTATION** (and its variants) – For the HASH-GEN hashtag generator, we use the Huggingface BART-base model, while for the tweet classifier, we use the BERTweet and TimeLMs. Additionally, for the document embeddings used in the Tweet Attention Module (TAM), we used the `all-MiniLM-L6-v2` model from the `sentence-transformers`¹ library.

3.3 Evaluation Metrics

For evaluation, we use the respective metrics following TweetEval benchmark [2], which are detailed as follows. For EMOJI, EMOTION, HATE, and OFFENSIVE, we use the macro F1-score. For IRONY, we use the F1-score on the positive class. For SENTIMENT, we use the macro recall score. Lastly, for STANCE, we use the average F1 scores of the “against” class and the “favor” class.

3.4 Implementation

For TAM, we utilize sentence-transformer [40] to initialize key vectors u_k and value vectors v_k to guarantee reliable addressing as a warm start for those vectors and they are updated during the training process. Different from u_k and v_k , the sentential encoding e of each input tweet x is represented as the average of its word representations which are randomly initialized to ensure their compatibility with the backbone Transformer’s vector space during training. For EAM, we extract relevant entities with social media NER systems² to construct the graph from the corpus. We use the AdamW [29] optimizer, together with a batch size of 16. We pad or truncate the inputs so that all of them have a length of 64. For the CNN and LSTM-based methods, we use a learning rate of $2e-5$, while for our HASHTATION model, we use a learning rate of $1e-5$ for both the hashtag generator and the tweet classifier. Beam search is applied to generate multiple phrases during hashtag generation with a beam size of 10 and a maximum sequence length of 40. These experiments were performed on Nvidia 2080Ti GPUs with 11 GB memory.

4 Experimental Results

The results on seven benchmark datasets are reported in Table 4. Overall, we observe that our HASHTATION model performs the best, outperforming both BERT-base and RoBERTa-base by a significant margin, as well as displaying gains of 1.2% and 1.3% over the state-of-the-art BERTweet and TimeLMs model, respectively. This shows that the hashtag generator and tweet classifier indeed help the model to make the right predictions. The first two models (Kim-CNN and BiLSTM) are not large language models, and they

¹<https://www.sbert.net/>

²<https://github.com/napsternxg/TwitterNER>

Table 5: Effects of different fusion methods on the TimeLMS model. We explore four methods: no hashtags, standard, pre-pending at the start, appending at the end. F1-scores are reported on seven tasks.

FUSION METHOD	EMOJI	EMOTION	HATE	IRONY	OFFENSIVE	SENTIMENT	STANCE
WITHOUT HASHTAGS	12.4	60.2	56.9	59.6	60.0	57.4	55.9
STANDARD	13.0	61.6	58.0	60.9	62.3	59.6	56.5
START	12.9	61.0	58.0	60.7	61.3	59.3	56.1
END	13.0	61.7	57.6	60.5	61.9	59.6	56.6

perform poorly on this low-resource classification task, with an average performance of 32.6 and 34.1, respectively. The first major performance improvement comes when we introduce the BERT and RoBERTa models, increasing the performance from 48.5 to around 50.4. This is due to the knowledge gained by the language model pre-training, which is helpful in tweet understanding, and consistent with previous language model literature [10, 28]. Next, another significant jump in performance occurs between the RoBERTa-base model and domain-specific models (*i.e.*, BERTweet and TimeLMS). This increase can be attributed to the additional pre-training in a specialized domain (social media texts like tweets). This method is called domain-adaptive continual pre-training and has been verified in other domains, like biomedical domain [23], clinical domain [1], scientific domain [5], social media domain [30, 39] and so on. The final major improvement in performance is from domain-specific pre-trained models to our HASHTATION model. The main differences are the addition of our proposed modules, the TAM and the EAM, as well as the incorporation of the generated hashtags. All of these components complement each other to contribute to the performance improvement. We also ablate certain components of our model to verify the contributions in Table 4. For HASHTATION-TL model, removing the TAM and EAM causes an average of 0.4% performance decrease equally. We observed similar trends on the HASHTATION-BT model as well. These ablation studies confirm the motivation behind generating hashtags by retrieving global information carried by general corpus and entities.

5 Analysis

In this section, we investigate the effects of different fusion methods (Section 5.1), the performance on different datasets (Section 5.2), the quality of the generated hashtags (Section 5.3), and we present a case study on specific examples (Section 5.4).

5.1 Effects of Fusion Methods

There are several different ways to implement the *Fuse* function mentioned in Equation 8. To verify the effects of different fusion methods, we examine the following four templates:

- **Without Hashtags:** $Fuse(H, x) = x$
- **Standard:** For each hashtag, we check whether it is a present or absent hashtag. For present hashtags, we prepend the hashtag symbol (#) to the existing word corresponding to the hashtag; for absent hashtags, we append them to the end of the tweet.
- **Start:** $Fuse(H, x) = [H, x]$
- **End:** $Fuse(H, x) = [x, H]$

Here, $[\cdot, \cdot]$ represents the concatenation operation. The performance on seven datasets is reported in Table 5. We find that there is no

big difference between **Standard** and **End** while a slight performance drop is observed on **Start**. This is consistent with the general structure of tweets, as hashtags usually appear at the end of tweets rather than at the beginning. Nevertheless, all three with-hashtag methods still outperform the without-hashtag fusion, indicating that the mere inclusion of these guiding hashtags is sufficient to increase the performance of the model.

5.2 Performance Across Datasets

Across all 7 datasets, adding the hashtags boosts the performance, as seen in Tables 4 and 5. However, improvement varies across datasets. For instance, with datasets such as IRONY and EMOTION, the improvement of adding hashtags (BERTWEET vs HASHTATION-BT and TIMELEMS vs HASHTATION-TL) is around 2%, which is much higher than other datasets such as HATE or STANCE. This tells us that certain types of data will benefit more from the enriched context than others. For instance, with EMOTION, if the hashtag generation model is able to generate a hashtag such as “happy” or “sad”, then they will serve as very strong signals for the classification model. In contrast, for a dataset such as HATE, the hashtag generation model might predict very useful hashtags in terms of topic (e.g. “elections”, “Ukraine”), but these might not be that useful for our particular downstream task (hate speech vs not hate speech) because both positive and negative classes will share these same topics.

5.3 Hashtag Quality

In our work, tweet classification is the downstream task, while hashtag generation is an important intermediate task. In this section, we check the quality of generated hashtags by comparing HASHTATION with three state-of-the-art hashtag generation baselines (Bi-Attn [51], One2Set [56], and SEGTRM [32]). To measure the quality of the generated hashtags, we adopt the precision, recall, and F1 metrics used in document retrieval. Additionally, we make a distinction between “present hashtags” which appear verbatim in the original tweet, and “absent hashtags” which do not match any contiguous word sequence in the original tweet. Results are reported in Table 7. From the table, we conclude that HASHTATION outperforms all baseline hashtag generation models which are specifically designed for this task, illustrating its effectiveness of leveraging tweet-level and entity-level information by TAM and EAM, respectively. Moreover, it is observed that absent hashtags are clearly much harder to predict than present hashtags, as evidenced by the difference in F1-scores (0.381 for present vs 0.122 for absent). This is consistent with our expectation that it is a lot easier to generate hashtags by simply copying important words from the text as opposed to having to find a new (absent) word to use. Given this challenge, it is thus

Table 6: Examples of hashtags generated by HASHTATION, as compared to ground truth hashtags.

Text	Dataset	Classification	Ground Truth Hashtags	Generated Hashtags
Riding with @user and some incredible people. Truly magical.	Emoji	sparkles	#soulcycle #lululemon	#soulcycle
Oh Canada shouldn't be sung like that.	Emotion	Anger	#terrible #MLBTHESHOW17	#terrible #Canada
Both #Frightening and #Demented #Sick. Now it will be our problem.	Hate	No	#Frightening, #Demented, #Sick, #SendThemBack, #KAG	#Sick #BuildThatWall #MAGA #KAG
Oh how I love working in Baltimore	Irony	Yes	#not	#not #Baltimore
@user Ya Obama on trade 2-YEARS AGO: "Trump is just NOT TELLING THE TRUTH How STUPID could our leaders be" God bless Trump!	Offensive	Yes	#MAGA #KAG	#MAGA #Trump
Messi: "To have a good team, everything starts there"	Sentiment	Joy	#fcblive	#fcblive
"Manspreading"? But women hog subway space, too!	Stance (Feminism)	against	#doublestandards	#doublestandards #manspreading

Table 7: Average precision, recall, and F1 scores for the hashtag generation task across all datasets.

Type	Precision	Recall	F1-Score
Bi-Attn [51]	0.202	0.069	0.103
One2Set [56]	0.231	0.095	0.134
SEGTRM [32]	0.352	0.127	0.187
HASHTATION (ours)	0.340	0.135	0.193
Present	0.848	0.234	0.381
Absent	0.171	0.108	0.122

quite impressive that HASHTATION is still able to correctly generate absent hashtags some of the time, considering that most of the input texts are very short. As observed in Table 6, there are certain input tweets from our dataset that are extremely short and seemingly contain no substantial content. For instance, for the tweet "Oh how I love working in Baltimore," HASHTATION is able to connect it with the hashtag #not, even though nowhere in the sentence is the word "not" mentioned. The reason HASHTATION is able to predict this is because there are many other tweets in the corpus which are ironic and also share phrases like "how I love" in them. This is a testament to the effectiveness of our TAM and EAM modules, which allow the HASHTATION model to focus not only on the current text but also on several other relevant texts which may contain more useful information. From Table 7, another key observation is that our model's precision scores are generally much higher than the model's recall scores. This means that when the model generates a hashtag, its prediction is likely to be correct, but there may be some ground-truth hashtags that the model fails to generate, which suggests that our model is generally quite conservative in making its predictions.

5.4 Case Study

In Table 8, we consider a stance detection task with the following example sentence: "3 cases of COVID-19 (coronavirus) in 2 schools in my city both involving teachers coming back from the north of Italy and having contact with the children for almost a week before anything was done." This sentence has a stance in favor of school closures but is not very explicit in its claims. In fact, nowhere in the tweet are the words "school" or "closure" ever mentioned, nor does it use very pointed positive/negative words. Without the relevant surrounding context, it would be difficult for a model to identify this stance that is subtly hidden in a seemingly-neutral sentence.

With the EAM and TAM modules, our HASHTATION model is able to identify some relevant tweets and entities which explicitly discuss school closures and even have hashtags such as #CloseTheSchools and #SchoolsMustShutDown (see Table 8). By identifying these relevant tweets, our model is able to gain a much broader context, providing key information beyond simply looking at the few words of the current tweet. Since our HASHTATION model pays attention to these broader contexts instead of myopically focusing on just the current text, it becomes much easier to generate the correct hashtag, which greatly aids in stance classification.

6 Related Work

6.1 Tweet Classification

Tweet classification aims to classify a piece of the tweet into different categories, which is helpful for understanding public opinion. Popular tweet classification tasks include stance detection [36], sentiment analysis [37] emotion recognition [35], hate speech detection [4], and so on. Several studies explored SVM-based approaches [11, 42], CNN-based models [47, 52] and RNN-based models [43, 58], respectively. Later, with the prevalence of Transformer architecture [45] and pre-trained language models [10, 27], BERT-based approaches [15, 25] attracted a lot of attention and had proven their effectiveness in this task. In addition to architectural improvement, two previous studies [53, 59] have proven the effectiveness of hashtag generation on stance detection, but have not explored mining external knowledge from a large pre-trained language model via hashtag generation task. Given the success of pre-trained language models (PLMs) [10, 12–14, 24, 28, 54, 64], we propose a framework that could leverage the advantages of both generative PLMs and discriminative PLMs for better tweet classification.

6.2 Generic Keyphrase Generation

There is a large body of research focused on the keyphrase generation for generic documents, such as news reports [48], and scientific documents [33], which mainly consist of two methodology streams, extractive and generative approaches. In detail, for the keyphrase generation of scientific document, many previous studies focused on extracting hashtags from documents [18, 31, 34, 48, 61]. Compared to extractive approaches, generative ones have attracted more attention in recent years for their ability to predict absent keyphrases for an input document. For example, [33] proposed CopyRNN, which

Table 8: An example of an input text whose stance is not very obvious and is implied instead of explicitly stated. We highlight the corresponding entities and relevant tweets that are extracted by the EAM and TAM modules.

(a) Example Input Text and Corresponding Hashtags (Entities from EAM highlighted in blue)

Input Text	G. Truth Hashtags	Pred. Hashtags
3 cases of COVID-19 (coronavirus) in 2 schools in my city, both involving teachers coming back from the north of Italy and having contact with the children for almost a week before anything was done.	#coronavirusuk, #CloseTheSchools	#CloseTheSchools

(b) Tweets with Relevant Entities (extracted by EAM)

What a lack of intestinal fortitude MENTION you will have the deaths of Australians on your hands. When we reach 12 pages of obituaries like Italy please look in the mirror for the cause.	#lockusdown, #coronavirus, #CloseTheSchools
At the end of the day, schools with sealed windows, and interior classrooms will have Coronavirus buildup that will increase COVID19 viral load! These building shouldn't be used. Children and teachers are not going to be victims of "risk mitigation"	#Coronavirus, #COVID19, #NotMyChild

(c) Most Relevant Tweets (extracted and ranked by TAM)

We would all feel very different about schools reopening if we had a government that: -Was trustworthy -Had best interest of children as only motivation -Based on Science -provided necessary resources to ALL schools -had a National plan None of these is true	#NotMyChild
Schools should just be closed,the 2020 curriculum will continue once the vaccine is found, even if it means 2020 curriculum get done in 2021-2022, we'll have to adjust and catch-up at some point.	#SchoolsMustShutDown

is an early study with attention and copy mechanism. [7] took correlation among multiple keyphrases into consideration to eliminate duplicate keyphrases. To further enhance keyphrase generation, other studies tried to utilize extra information: [55] proposed to assign synthetic keyphrases to unlabeled documents and then use them to enlarge the training data; [8] retrieved similar documents from the training data for the input document and encoded their keyphrases as external knowledge, while [9] leveraged title information for this task. To increase the diversity of keyphrases, a reinforcement learning approach is introduced by [6] to encourage their model to generate the correct number of keyphrases with an adaptive reward. Although existing models are capable of predicting both present and absent keyphrases, there is still potential to facilitate keyphrase generation with unlabeled data such as relevant documents. In doing so, HASHTATION offers a more effective and efficient solution.

6.3 Hashtag Generation for Social Media Text

For social media text, hashtag generation aims to produce hashtags for a given post automatically. The nature of hashtags makes it difficult for researchers to directly transfer the keyphrase generation methods to this domain mainly because hashtags are short in length and extremely informal. Based on the features of social media text, prior work can be categorized into three types: extractive methods, classification methods, and generative methods. Extractive methods [62, 63] try to identify important words from the input post and extract them as hashtags, which suffers from the lack of ability to deal with absent hashtags. Classification methods [17, 60] formulate the hashtag prediction problem from the classification point of view, which aims to classify the true hashtags from a pre-defined candidate list. However, the hashtags are diverse in social media and a lot of new hashtags are generated by users every day, so we can not include sufficient words in the candidate list. To produce both present and absent hashtags in a more flexible way, generative methods [49, 51] try to generate hashtags from scratch following

a sequence-to-sequence (Seq2Seq) paradigm, which is capable of generating unseen hashtags in the post without needing a candidate list. Because the microblog post is short and lacks sufficient information, previous generative methods use either the topic information from the given post or associated conversation information to enrich the Seq2Seq model. However, only considering topic information or conversation is not enough because they omit the information carried by relevant posts and conversations, and thus the model cannot easily adapt to dynamic semantic context, which is fast changing in social media platforms. Compared to the above studies, HASHTATION provides a solution that explicitly incorporates the associated conversation context and the relevant posts and conversations, in a simple architecture.

7 Conclusion

This paper proposed a novel two-stage framework for hashtag-guided low-resource tweet classification with a **Hashtag Generator** and a **Tweet Classifier**. The hashtag generator introduces relevant tweets and entities, and summarizes them into hashtags, while the tweet classifier leverages the global contexts to classify the tweets. One benefit of our model is the capability of leveraging the pre-trained knowledge of both generative models (*e.g.*, BART) and discriminative models (*e.g.*, BERT). Another benefit is the enriched multi-grained global contexts introduced by the hashtags, especially the absent ones. To conduct our experiments, we extract the hashtags based on seven tweet classification datasets and use these to train our hashtag models. Meanwhile, for the classification models, we consider a low-resource setting by sampling from the training sets. Experiments on a tweet classification benchmark demonstrate the effectiveness of our approach in identifying relevant tweets and entities, as well as in generating present and absent hashtags which provide valuable contexts.

Acknowledgments

We thank the anonymous reviewers for their valuable suggestions. This work was supported by the General Research Fund (GRF) of Hong Kong (No. 16310222 and No. 16201320). Shizhe Diao was supported by the Hong Kong Ph.D. Fellowship Scheme (HKPFS).

References

- [1] Emily Alsentzer, John Murphy, William Boag, Wei-Hung Weng, Di Jindi, Tristan Naumann, and Matthew McDermott. 2019. Publicly Available Clinical BERT Embeddings. In *Proceedings of the 2nd Clinical Natural Language Processing Workshop*. Association for Computational Linguistics, Minneapolis, Minnesota, USA, 72–78. <https://doi.org/10.18653/v1/W19-1909>
- [2] Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa Anke, and Leonardo Neves. 2020. TweetEval: Unified Benchmark and Comparative Evaluation for Tweet Classification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 1644–1650. <https://doi.org/10.18653/v1/2020.findings-emnlp.148>
- [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] Valerio Basile, Cristina Bosco, Elisabetta Fersini, Debora Nozza, Viviana Patti, Francisco Manuel Rangel Pardo, Paolo Rosso, and Manuela Sanguinetti. 2019. SemEval-2019 Task 5: Multilingual Detection of Hate Speech Against Immigrants and Women in Twitter. In *Proceedings of the 13th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, Minneapolis, Minnesota, USA, 54–63. <https://doi.org/10.18653/v1/S19-2007>
- [5] Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. SciBERT: A Pretrained Language Model for Scientific Text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3615–3620. <https://doi.org/10.18653/v1/D19-1371>
- [6] Hou Pong Chan, Wang Chen, Lu Wang, and Irwin King. 2019. Neural Keyphrase Generation via Reinforcement Learning with Adaptive Rewards. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 2163–2174. <https://doi.org/10.18653/v1/P19-1208>
- [7] Jun Chen, Xiaoming Zhang, Yu Wu, Zhao Yan, and Zhoujun Li. 2018. Keyphrase Generation with Correlation Constraints. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 4057–4066. <https://doi.org/10.18653/v1/D18-1439>
- [8] Wang Chen, Hou Pong Chan, Piji Li, Lidong Bing, and Irwin King. 2019. An Integrated Approach for Keyphrase Generation via Exploring the Power of Retrieval and Extraction. 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)*. Association for Computational Linguistics, Minneapolis, Minnesota, 2846–2856. <https://doi.org/10.18653/v1/N19-1292>
- [9] Wang Chen, Yifan Gao, Jiani Zhang, Irwin King, and Michael R. Lyu. 2019. Title-Guided Encoding for Keyphrase Generation. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*. AAAI Press, 6268–6275. <https://doi.org/10.1609/aaai.v33i01.33016268>
- [10] 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, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [11] Kuntal Dey, Ritvik Shrivastava, and Saroj Kaushik. 2017. Twitter stance detection—A subjectivity and sentiment polarity inspired two-phase approach. In *2017 IEEE international conference on data mining workshops (ICDMW)*. IEEE, 365–372.
- [12] Shizhe Diao, Jiabin Bai, Yan Song, Tong Zhang, and Yonggang Wang. 2020. ZEN: Pre-training Chinese Text Encoder Enhanced by N-gram Representations. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 4729–4740. <https://doi.org/10.18653/v1/2020.findings-emnlp.425>
- [13] Shizhe Diao, Xuechun Li, Yong Lin, Zhichao Huang, and Tong Zhang. 2022. Black-box prompt learning for pre-trained language models. *ArXiv preprint abs/2201.08531* (2022). <https://arxiv.org/abs/2201.08531>
- [14] Shizhe Diao, Ruijia Xu, Hongjin Su, Yilei Jiang, Yan Song, and Tong Zhang. 2021. Taming Pre-trained Language Models with N-gram Representations for Low-Resource Domain Adaptation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Online, 3336–3349. <https://doi.org/10.18653/v1/2021.acl-long.259>
- [15] Shalmoli Ghosh, Prajwal Singhania, Siddharth Singh, Koustav Rudra, and Saptarshi Ghosh. 2019. Stance detection in web and social media: a comparative study. In *International Conference of the Cross-Language Evaluation Forum for European Languages*. Springer, 75–87.
- [16] Andrew Hoang, Antoine Bosselut, Asli Celikyilmaz, and Yejin Choi. 2019. Efficient Adaptation of Pretrained Transformers for Abstractive Summarization. *ArXiv preprint abs/1906.00138* (2019). <https://arxiv.org/abs/1906.00138>
- [17] Haoran Huang, Qi Zhang, Yeyun Gong, and Xuanjing Huang. 2016. Hashtag Recommendation Using End-To-End Memory Networks with Hierarchical Attention. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. The COLING 2016 Organizing Committee, Osaka, Japan, 943–952. <https://aclanthology.org/C16-1090>
- [18] Anette Hulth. 2003. Improved Automatic Keyword Extraction Given More Linguistic Knowledge. In *Proceedings of the 2003 Conference on Empirical Methods in Natural Language Processing*. 216–223. <https://aclanthology.org/W03-1028>
- [19] Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. 2019. Ctrl: A Conditional Transformer Language Model for Controllable Generation. *ArXiv preprint abs/1909.05858* (2019). <https://arxiv.org/abs/1909.05858>
- [20] Urvashi Khandelwal, Kevin Clark, Dan Jurafsky, and Lukasz Kaiser. 2019. Sample Efficient Text Summarization Using a Single Pre-Trained Transformer. *ArXiv preprint abs/1905.08836* (2019). <https://arxiv.org/abs/1905.08836>
- [21] Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1746–1751. <https://doi.org/10.3115/v1/D14-1181>
- [22] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net. <https://openreview.net/forum?id=SJU4ayYgl>
- [23] Jinhyuk Lee, Wonjin Yoon, Sungdoon Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: A Pre-Trained Biomedical Language Representation Model for Biomedical Text Mining. *Bioinformatics* 36, 4 (2020), 1234–1240.
- [24] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 7871–7880. <https://doi.org/10.18653/v1/2020.acl-main.703>
- [25] Yingjie Li and Cornelia Caragea. 2021. Target-Aware Data Augmentation for Stance Detection. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online, 1850–1860. <https://doi.org/10.18653/v1/2021.naacl-main.148>
- [26] Yang Liu and Mirella Lapata. 2019. Hierarchical Transformers for Multi-Document Summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 5070–5081. <https://doi.org/10.18653/v1/P19-1500>
- [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. *ArXiv abs/1907.11692* (2019).
- [28] 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. *ArXiv preprint abs/1907.11692* (2019). <https://arxiv.org/abs/1907.11692>
- [29] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net. <https://openreview.net/forum?id=Bkg6RiCqY7>
- [30] Daniel Loureiro, Francesco Barbieri, Leonardo Neves, Luis Espinosa Anke, and Jose Camacho-collados. 2022. TimeLMs: Diachronic Language Models from Twitter. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Dublin, Ireland, 251–260. <https://aclanthology.org/2022.acl-demo.25>
- [31] Yi Luan, Mari Ostendorf, and Hannaneh Hajishirzi. 2017. Scientific Information Extraction with Semi-supervised Neural Tagging. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association

- for Computational Linguistics, Copenhagen, Denmark, 2641–2651. <https://doi.org/10.18653/v1/D17-1279>
- [32] Qianren Mao, Xi Li, Bang Liu, Shu Guo, Peng Hao, Jianxin Li, and Lihong Wang. 2022. Attend and select: A segment selective transformer for microblog hashtag generation. *Knowledge-Based Systems* 254 (2022), 109581.
- [33] Rui Meng, Sanqiang Zhao, Shuguang Han, Daqing He, Peter Brusilovsky, and Yu Chi. 2017. Deep Keyphrase Generation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Vancouver, Canada, 582–592. <https://doi.org/10.18653/v1/P17-1054>
- [34] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order into Text. In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Barcelona, Spain, 404–411. <https://aclanthology.org/W04-3252>
- [35] Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. SemEval-2018 Task 1: Affect in Tweets. In *Proceedings of the 12th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, New Orleans, Louisiana, 1–17. <https://doi.org/10.18653/v1/S18-1001>
- [36] Saif Mohammad, Svetlana Kiritchenko, Parinaz Sobhani, Xiaodan Zhu, and Colin Cherry. 2016. SemEval-2016 Task 6: Detecting Stance in Tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. Association for Computational Linguistics, San Diego, California, 31–41. <https://doi.org/10.18653/v1/S16-1003>
- [37] Preslav Nakov, Sara Rosenthal, Zornitsa Kozareva, Veselin Stoyanov, Alan Ritter, and Theresa Wilson. 2013. SemEval-2013 Task 2: Sentiment Analysis in Twitter. In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*. Association for Computational Linguistics, Atlanta, Georgia, USA, 312–320. <https://aclanthology.org/S13-2052>
- [38] 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*. Association for Computational Linguistics, Online, 9–14. <https://doi.org/10.18653/v1/2020.emnlp-demos.2>
- [39] 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*. Association for Computational Linguistics, Online, 9–14. <https://doi.org/10.18653/v1/2020.emnlp-demos.2>
- [40] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Association for Computational Linguistics, Hong Kong, China, 3982–3992. <https://doi.org/10.18653/v1/D19-1410>
- [41] M. Schuster and K.K. Paliwal. 1997. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing* 45, 11 (1997), 2673–2681. <https://doi.org/10.1109/78.650093>
- [42] Anirban Sen, Manjira Sinha, Sandya Mannarswamy, and Shourya Roy. 2018. Stance classification of multi-perspective consumer health information. In *Proceedings of the ACM India Joint International Conference on Data Science and Management of Data*. 273–281.
- [43] Umme Aymun Siddiqua, Abu Nowshed Chy, and Masaki Aono. 2019. Tweet Stance Detection Using an Attention based Neural Ensemble Model. 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)*. Association for Computational Linguistics, Minneapolis, Minnesota, 1868–1873. <https://doi.org/10.18653/v1/N19-1185>
- [44] Cynthia Van Hee, Els Lefever, and Véronique Hoste. 2018. SemEval-2018 Task 3: Irony Detection in English Tweets. In *Proceedings of The 12th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, New Orleans, Louisiana, 39–50. <https://doi.org/10.18653/v1/S18-1005>
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.), 5998–6008. <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>
- [46] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net. <https://openreview.net/forum?id=IjXmpkCZ>
- [47] Prashanth Vijayaraghavan, Ivan Sysoev, Soroush Vosoughi, and Deb Roy. 2016. DeepStance at SemEval-2016 Task 6: Detecting Stance in Tweets Using Character and Word-Level CNNs. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. Association for Computational Linguistics, San Diego, California, 413–419. <https://doi.org/10.18653/v1/S16-1067>
- [48] Xiaojun Wan and Jianguo Xiao. 2008. Single Document Keyphrase Extraction Using Neighborhood Knowledge. In *Proceedings of the 23rd National Conference on Artificial Intelligence - Volume 2*. 855–860.
- [49] Yue Wang, Jing Li, Hou Pong Chan, Irwin King, Michael R. Lyu, and Shuming Shi. 2019. Topic-Aware Neural Keyphrase Generation for Social Media Language. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 2516–2526. <https://doi.org/10.18653/v1/P19-1240>
- [50] Yue Wang, Jing Li, Irwin King, Michael R. Lyu, and Shuming Shi. 2019. Microblog Hashtag Generation via Encoding Conversation Contexts. 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)*. Association for Computational Linguistics, Minneapolis, Minnesota, 1624–1633. <https://doi.org/10.18653/v1/N19-1164>
- [51] Yue Wang, Jing Li, Irwin King, Michael R. Lyu, and Shuming Shi. 2019. Microblog Hashtag Generation via Encoding Conversation Contexts. 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)*. Association for Computational Linguistics, Minneapolis, Minnesota, 1624–1633. <https://doi.org/10.18653/v1/N19-1164>
- [52] Wan Wei, Xiao Zhang, Xuqin Liu, Wei Chen, and Tengjiao Wang. 2016. pkudblab at SemEval-2016 Task 6 : A Specific Convolutional Neural Network System for Effective Stance Detection. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. Association for Computational Linguistics, San Diego, California, 384–388. <https://doi.org/10.18653/v1/S16-1062>
- [53] Jason Weston, Sumit Chopra, and Keith Adams. 2014. #TagSpace: Semantic Embeddings from Hashtags. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, 1822–1827. <https://doi.org/10.3115/v1/D14-1194>
- [54] Ze Yang, Wei Wu, Can Xu, Xinnian Liang, Jiaqi Bai, Liran Wang, Wei Wang, and Zhoujun Li. 2020. StyleDGPT: Stylized Response Generation with Pre-trained Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 1548–1559. <https://doi.org/10.18653/v1/2020.findings-emnlp.140>
- [55] Hai Ye and Lu Wang. 2018. Semi-Supervised Learning for Neural Keyphrase Generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 4142–4153. <https://doi.org/10.18653/v1/D18-1447>
- [56] Jiacheng Ye, Tao Gui, Yichao Luo, Yige Xu, and Qi Zhang. 2021. One2Set: Generating Diverse Keyphrases as a Set. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 4598–4608.
- [57] Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. SemEval-2019 Task 6: Identifying and Categorizing Offensive Language in Social Media (OffensEval). In *Proceedings of the 13th International Workshop on Semantic Evaluation*. Association for Computational Linguistics, Minneapolis, Minnesota, USA, 75–86. <https://doi.org/10.18653/v1/S19-2010>
- [58] Guido Zarrella and Amy Marsh. 2016. MITRE at SemEval-2016 Task 6: Transfer Learning for Stance Detection. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. Association for Computational Linguistics, San Diego, California, 458–463. <https://doi.org/10.18653/v1/S16-1074>
- [59] Guido Zarrella and Amy Marsh. 2016. MITRE at SemEval-2016 Task 6: Transfer Learning for Stance Detection. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*. Association for Computational Linguistics, San Diego, California, 458–463. <https://doi.org/10.18653/v1/S16-1074>
- [60] Qi Zhang, Jiawen Wang, Haoran Huang, Xuanjing Huang, and Yeyun Gong. 2017. Hashtag Recommendation for Multimodal Microblog Using Co-Attention Network. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, Carles Sierra (Ed.), ijcai.org, 3420–3426. <https://doi.org/10.24963/ijcai.2017/478>
- [61] Qi Zhang, Yang Wang, Yeyun Gong, and Xuanjing Huang. 2016. Keyphrase Extraction Using Deep Recurrent Neural Networks on Twitter. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 836–845. <https://doi.org/10.18653/v1/D16-1080>
- [62] Qi Zhang, Yang Wang, Yeyun Gong, and Xuanjing Huang. 2016. Keyphrase Extraction Using Deep Recurrent Neural Networks on Twitter. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Austin, Texas, 836–845. <https://doi.org/10.18653/v1/D16-1080>
- [63] Yingyi Zhang, Jing Li, Yan Song, and Chengzhi Zhang. 2018. Encoding Conversation Context for Neural Keyphrase Extraction from Microblog Posts. In *Proceedings of the 2018 Conference of the North American Chapter of the Association*

- for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 1676–1686. <https://doi.org/10.18653/v1/N18-1151>
- [64] Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. DIALOGPT : Large-Scale Generative Pre-training for Conversational Response Generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Online, 270–278. <https://doi.org/10.18653/v1/2020.acl-demos.30>

A Details of Datasets

We perform our experiments on 7 diverse tweet classification tasks in the TweetEval benchmark datasets [2]. All 7 tasks are taken from previous SemEval tasks (and corresponding datasets) and are as follows: emotion recognition [35], emoji prediction [3], irony detection [44], hate speech detection [4], offensive language identification [57], sentiment analysis [37], and stance detection [36]. Additional details about these tasks can be found in Table 3.

Creating the hashtag generation input and output. For hashtag generation, we train a separate hashtag generator for each task, as the datasets cover significantly different domains. To extract hashtags, we searched for appearances of the octothorpe symbol (#) and considered the contiguous string following it, until a whitespace is reached. These hashtags can either appear mid-tweet or at the end of the tweet. For hashtags that appear mid-tweet, we remove the hashtag symbol (#) but keep the word itself, as removing these words may disrupt the coherence of the sentence. On the other hand, for hashtags that appear at the end of the tweet, we completely remove these from the body of the tweet. If we keep these words at the end of the sentence, the hashtag generator will easily pick up on these false signals and just learn to return the last few words of every sentence without actually learning to perform hashtag generation. By considering this setting of keeping in-sentence hashtags and removing end-of-sentence hashtags, we are able to train the model to perform prediction for both present hashtags and absent hashtags. Examples of this overall process are outlined in Table 9. Similar to Wang et al. [50], we implement some preprocessing steps to clean up our tweets. Links, mentions, and numbers were replaced with “URL”, “MENTION”, and “DIGIT”, respectively. The details of the processed dataset are shown in Table 3.

Low-resource setting. As discussed in Section 1, real-life classification labels are often difficult to acquire for the rapidly changing landscape of tweets. As such, we conduct our experiments in a low-resource setting, where we randomly select a subset from the training set. To keep the sizes of different datasets relatively consistent, we use the following sampling ratios: If the size of the original dataset is < 5000, we sample 10% of the dataset. If it’s between 5000 and 10000, we sample 5%, and if it’s > 10000, we sample 1%. All experiments are performed over ten different random seeds.

B Details of Baselines

To verify the effectiveness of our HASHTATION model, we compare its performance against the following baseline classification models: **Kim-CNN** [21], **BiLSTM** [41], **BERT** [10], **RoBERTa** [27], **BERTweet** [39], and **TimeLMs** [30]. BERTweet and TimeLMs are two state-of-the-art models that are pre-trained on 850 million and 124 million English tweets, respectively.

- **Kim-CNN** [21]. A simple convolutional neural networks framework with a classification layer attached to the end. For this CNN-based model, we use kernel sizes of 2,3,4,5 and 64 filters for each, with a dropout of 0.5 before the linear layer.
- **BiLSTM** [41]. A bidirectional long short-term memory network considering the temporal order of words in the tweet. The hidden size and the dropout rate for the LSTM are set to 512 and 0.2, respectively. We use an LSTM hidden size of 512 and an LSTM

dropout of 0.2. In addition, we use a dropout of 0.5 and a size of 32 for the final linear classification layer.

- **BERT** [10]. A BERT-base model pre-trained on generic corpus with a classification head fine-tuned on specific task.
- **RoBERTa** [27]. A RoBERTa-base model pre-trained on generic corpus with a classification head fine-tuned on specific task.
- **BERTweet** [39]. A RoBERTa-base model pre-trained on 850M English Tweets with a classification head fine-tuned on specific task.
- **TimeLMs** [30]. A RoBERTa-base model specialized on diachronic Twitter data and pre-trained on 124M English Tweets.
- **HASHTATION** (and its variants) – For the HASH-GEN hashtag generator, we use the Huggingface BART-base model, while for the tweet classifier, we use the BERTweet and TimeLMs. Additionally, for the document embeddings used in the Tweet Attention Module (TAM), we used the `all-MiniLM-L6-v2` model from the `sentence-transformers`³ library.

³<https://www.sbert.net/>

Table 9: Examples of how we generate the training data and ground truth for our hashtag seq2seq model. The blue text represents the present hashtags, while the red text represents the absent hashtags.

Original Tweet	Processed Tweet (New Input)	Hashtags (Target Output)
Going to #BigApple tomorrow! Loving #NJTransit #commute! See you all!	Going to BigApple tomorrow! Loving NJ-Transit commute! See you all!	BigApple, NJTransit, commute
Abortion IS NOT a political issue it is a MORAL issue. #Catholic #Christian #Conservative	Abortion IS NOT a political issue it is a MORAL issue.	Catholic, Christian, Conservative
Latest #crypto developments: Top 10 coins to watch #NotFinancialAdvice #DoYourOwnResearch	Latest crypto developments: Top 10 coins to watch	crypto, NotFinancialAdvice, DoYourOwnResearch