

# Naive Bayes and BiLSTM Ensemble for Discriminating between Mainland and Taiwan Variation of Mandarin Chinese

Li Yang

Tongji University  
Shanghai  
China

li.yang@tongji.edu.cn

Yang Xiang

Tongji University  
Shanghai  
China

shxiangyang@tongji.edu.cn

## Abstract

Automatic dialect identification is a more challenging task than language identification, requiring ability to discriminate between different varieties within the same language family. In this paper, we propose an ensemble based system, which combine traditional machine learning models trained on bag of n-gram features, with deep learning models trained on word embeddings, to solve the Discriminating between Mainland and Taiwan Variation of Mandarin Chinese (DMT) shared task at VarDial 2019. Our experiments show that a character bigram-trigram combination based Naive Bayes is a very strong model for identifying varieties of Mandarin Chinese. By combining Naive Bayes and BiLSTM using average ensembling, a simple yet effective approach, our system achieved an macro-averaged F1 score of number1% and number2% in two tracks, ranking th and th out of number teams respectively.

## 1 Introduction

Dialect identification can be considered as a special case of language identification, which aims at distinguishing related languages or varieties of a specific language. Being capable of detecting dialect accurately is an important step for many NLP pipelines and applications, such as automatic speech recognition, machine translation and multilingual data acquisition. While language identification can already achieve relatively high performance with a simple model, dialect identification is still a tough problem remained to be tackled. In contrast to the language identification scenario, the linguistic difference between related languages is less obvious. For that reason, dialect identification has drawn many researchers' attention in recent years, becoming an important and attractive research topic.

Term	Mainland China	Taiwan
taxi	出租车	計程車
bicycle	自行车	腳踏車
softmax	软件	軟體
program	程序	程式
kindergarten	幼儿园	幼稚園

Table 1: Different expressions with the same meaning used in Mainland China and Taiwan.

Mandarin Chinese is a group of related varieties of Chinese spoken across many different regions. The group includes *Putonghua*, the official language of Mainland China, and *Guoyu*, another Mandarin variant widely spoken in Taiwan. However related they are, there are still some difference between these two varieties. First, the most notable one is the character set they use. Mainland Chinese uses simplified Chinese characters, as opposed to the traditional Chinese characters for Taiwanese Chinese. Take the phrase “*natural language processing*” as an example, its simplified character form adopted in Mainland China is “自然语言处理”, while the traditional character form in Taiwan is “自然語言處理”. Second, some vocabularies differ. There are some terms can be understood by both varieties to mean the same thing. However, their preferred usage distinguishes. Table ?? lists some examples. Apart from character form and vocabularies, the pronunciation especially in terms of tone is also different. But we don't discuss in this paper since it's less relevant.

The DMT task, first introduced by VarDial evaluation campaign this year, aiming at determining whether a sentence belongs to news articles from Mainland China or from Taiwan. The organizers prepare two versions of the same corpus, traditional and simplified, and ask participants to

predict the labels for text instances in both tracks. For that reason, one can't make use of character form to discriminate these two language varieties. Mainstream approaches to dialect identification is to regard it as a text classification task and use a linear support vector machine (SVM) with bag of n-gram features as input to solve it. However, we want to know what's the best classification algorithm for DMT task. Therefore, we experiment with several classical machine learning models trained on different word or character level n-gram features and feature combinations. Besides, deep learning methods have currently achieved remarkable success in many NLP tasks, including question answering (?), sentiment analysis (?), machine translation (?) and natural language inference (?). To investigate how much can deep neural networks help identify language varieties, we test 7 different deep learning models, varying from CNN based, RNN based to CNN-RNN hybrid models. Fully performance comparison to machine learning models is also conducted. Finally, we explore different ways to combine the classifiers we discuss before.

The rest of this paper is organized as follows. Section ?? introduces related work. In Section ??, we describe the dataset and our proposed solution. Experimental results are presented in Section ??, and we give our conclusion and future work in Section ??.

## 2 Related work

A number of works have devoted to identify different language varieties or related language, especially since the series of VarDial evaluation campaigns (???). (?) studies on English dialect identification and presents several classification approaches to classify Australian, British and Canadian English. (?) utilizes a character n-gram and a word n-gram language model for automatic classification of two written varieties of Portuguese: European and Brazilian. (?) conducts an initial study on the dialects of Romanian and proposes using the orthographic and phonetic features of the words to build a dialect classifier. (?) uses a majority-vote ensemble of the Naive Bayes, CRF and SVM systems for Swiss German dialects identification. (?) uses two SVM classifiers: one trained on word n-grams features and one trained on Pos n-grams to determine a document to be in Flemish Dutch or Netherlandic Dutch. (?) uses a

unified SVM model based on character and word n-grams features with careful hyperparameter tuning for 4 language/dialect identification tasks.

Methods to discriminate between varieties of Mandarin Chinese haven't been well studied. (?) uses a top-bag-of-word similarity based contrastive approach to reflect distance among three varieties of Mandarin: Mainland China, Singapore and Taiwan. (?) deals with 6 varieties of Mandarin: Mainland, Hong Kong, Taiwan, Macao, Malaysia and Singapore. They discover that character bigram and word segmentation based feature work better than traditional character unigram, and some features such as character form, PMI-based and word alignment-based features can help improve performance. However, a thorough comparison of different algorithms and architectures has yet to be conducted.

## 3 Data and Methodology

### 3.1 Data

The DMT task is provided with labeled sentences belonging to news articles from Mainland China or from Taiwan, which compose of 18770 instances for training set, 2000 for validation set and 2000 for test set. As shown in Table ??, the DMT dataset has a perfectly balanced class distribution. The average sentence length (in word level) of two varieties are almost the same, and short as well. It's worth mentioning that the organizers have prepared two versions of the same dataset: traditional and simplified version, which means we can't utilize character form feature to help us discriminate between these two language varieties. Since the sentences have been tokenized and punctuation removed from the texts, we don't apply any preprocessing on the dataset.

### 3.2 Machine Learning Models with n-gram features

Machine learning models based on feature engineering are the most common methods for dialect identification. In this paper, we experiment with 3 different classifiers: (1) logistic regression (**LR**), (2) linear support vector machine (**SVM**), and (3) multinomial Naive Bayes (**MNB**) based on bag of n-gram features. We also examine other Naive Bayes models such as Gaussian Naive Bayes and Bernoulli Naive Bayes, but they are inferior to multinomial Naive Bayes on the validation set. The bag of n-gram features include word and char-

Variety	Number of instances			Sentence length			
	train	valid	test	min	max	avg	st.dev.
Mainland China	9385	1000	1000	5	66	9.63	3.73
Taiwan	9385	1000	1000	6	48	9.24	3.30

Table 2: Statistics of dataset for each variety. Sentence lengths are calculated based on word-level tokens from training and validation set.

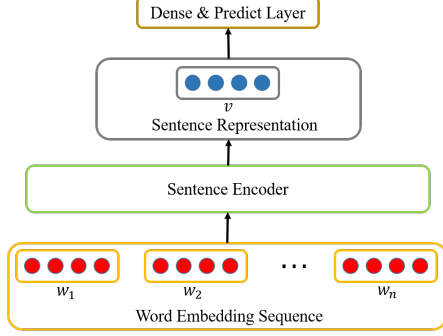


Figure 1: The overall framework of deep models.

acter level n-gram with size ranging from unigram to n-gram of specific order. We conduct a set of experiments to fully explore the most contributing features and feature combinations for DMT task, and the results are shown in next Section.

### 3.3 Deep Learning Models with word embeddings

Deep neural networks (DNNs) are of growing interest for their capacity to learn text representation from data without careful engineering of features. For short-text classification task, Convolution neural network (CNN) and recurrent neural network (RNN) are two mainstream DNN architectures. In this paper, we examine a number of deep models based on the same framework to solve DMT task. Figure ?? shows a high-level view of the framework. Vertically, the figure depicts 3 major components: (1) **Input Embedding Layer**. Suppose the sentence has  $n$  tokens, we use a pre-trained embedding method Word2vec (?) based on training data to represent it in a sequence of word embeddings:

$$S = (w_1, w_2, \dots, w_n) \quad (1)$$

where  $w_i$  is a vector representing a  $d$  dimensional word embedding for the  $i$ -th word in the sentence.  $S$  is thus a sequence that concatenates all word embeddings together. We do try using character embeddings as input and other pre-trained embedding methods such as Glove (?) and Fasttext (?) but observed no further improvement on valida-

tion set. (2) **Sentence Encoder Layer**. Sentence encoder, specified by different deep models, processes the input word embedding sequence and output a sentence representation:

$$v = \text{encode}(S) \quad (2)$$

(3) **Output Layer**. After obtaining sentence vector, we feed it through one hidden dense layer with 256 unit and a final predict dense layer:

$$\hat{y} = \sigma(W_p \gamma(W_h v + b_h) + b_p) \quad (3)$$

where  $W_h$  and  $b_h$  are the parameters for hidden layer,  $W_p$  and  $b_p$  are the parameters for predict layer,  $\gamma$  and  $\sigma$  are relu and sigmoid activation function respectively,  $\hat{y} \in \mathbb{R}$  represents the predicted score for positive class. During training process, we minimize the binary cross-entropy loss defined as follow:

$$L = -\frac{1}{N} \sum_{i=1}^N (y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)) \quad (4)$$

where  $y_i$  is the ground-truth.

We compare 7 different deep models to encode sentence into fixed-size vector, including cnn-based, rnn-based and cnn-rnn hybrid neural networks.

- **CNN**: First introduced by (?), the convolution network applies a convolution operation with a filter  $W_c \in \mathbb{R}^{h \times d}$  to a window of  $h$  words to produce a new feature.

$$c_i = f(W_c \cdot x_{i:i+h-1} + b) \quad (5)$$

By applying this filter to each possible window of words in the sentence, a feature map can be produced. In this paper, we use 300 filters with size ranging from 2 to 5 to extract four 300 dimensional feature maps. After that, we apply max-over-time pooling operation by taking the highest value for each feature map to capture the most important feature, then concatenate all the feature to represent the input sentence.

- **DCNN:** (?) use a dynamic convolution neural network (DCNN) that alternates wide convolution layers and dynamic  $k$ -Max pooling layers for sentence modeling. By applying  $k$ -Max pooling operation, a feature graph over the sentence can be induced, being capable of explicitly capturing both short and long-range relations.
- **DPCNN:** (?) propose a deep convolutional neural network by stacking convolution blocks (two convolution layers and a shortcut connection) interleaved with pooling layers with stride 2 for downsampling. The 2-stride downsampling reduces the size of the internal representation of each text by half, enabling efficient representation of long-range association in the text. And the shortcut connection ensure training of deep networks. DPCNN has been shown powerful in many text classification task.
- **BiLSTM:** LSTM is an effective neural network for sentence modeling for its ability to capture long-term dependencies. BiLSTM use a forward and a backward LSTM to process sequence, such that each produced hidden state can contain information from context in two opposite direction. Specifically, at each time step  $t$ , hidden state  $h_t$  is the concatenation of results from forward and backward LSTM:

$$\begin{aligned}\vec{h}_t &= \overrightarrow{\text{LSTM}}(\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_t) \\ \overleftarrow{h}_t &= \overleftarrow{\text{LSTM}}(\mathbf{w}_n, \mathbf{w}_{n-1}, \dots, \mathbf{w}_t) \\ h_t &= [\vec{h}_t, \overleftarrow{h}_t]\end{aligned}\quad (6)$$

After obtaining hidden state sequence, we apply max-over-time pooling operation to form a fixed-size vector as sentence representation  $v$ .

- **Self-attentive BiLSTM:** Attention mechanism is most commonly used in sequence-to-sequence models to attend to encoder states (??). In this paper, we make use of attention, more specifically, self-attention (?) to obtain a distribution over features learned from BiLSTM (a.k.a hidden states). Suppose  $H$  is the output hidden vectors of BiLSTM:  $H = (\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n)$ , we can calculate the

attention vector  $\alpha$  and the final sentence representation  $v$  as follows:

$$\begin{aligned}e_t &= \mathbf{U}_a^\top \tanh(\mathbf{W}_a \mathbf{h}_t) \\ \alpha_t &= \frac{\exp(e_t)}{\sum_{i=1}^T \exp(e_i)} \\ v &= \sum_{i=1}^T \alpha_i h_i\end{aligned}\quad (7)$$

where  $\mathbf{W}_a \in \mathbb{R}^{2d \times 2d}$  and  $\mathbf{U}_a \in \mathbb{R}^{2d \times 1}$  are parameters of attention layer (we use  $d$  units for LSTM, thus  $h_t$  is a  $2d$  dimensional vector). Using self-attention allows sentence to attend to itself, thus we can extract the most relevant information.

- **CNN-BiLSTM:** Similar as (?), we first use CNN to extract a higher-level sequence representation from word embedding sequence, and then feed them into BiLSTM to obtain final sentence representation. By combining CNN and BiLSTM, we are able to capture both local features of phrases and global information of sentence.
- **BiLSTM-CNN:** We also try to use the BiLSTM layer as feature extractor first and then feed the hidden states to the CNN layer, which we call BiLSTM-CNN.

### 3.4 Ensemble Models

Classifier ensembles are a way of combining different models with the goal of improving overall performance through enhanced decision making, which has been shown to achieve better results than single classifier. In this paper, we explore 4 ensemble strategies to integrate output (predict label or probability) from models introduced above and reach a final decision.

- **Mean Probability:** Simply take an average of predictions from all the models and use it to make the final prediction.
- **Highest Confidence:** The class label that receives vote with the highest probability is selected as the final prediction.
- **Majority Voting:** Each classifier votes for a single class label. The votes are summed and the label with majority votes (over 50%) wins. In case of tie, the ensemble result falls back to the prediction by the model with highest performance on validation set.

- **Meta-Classifer:** Use the individual classifier outputs along with training labels to train a second-level meta-classifier. The second meta-classifier then serves to predict the final prediction. Meta-Classifer is also referred to as classifier stacking.

While the first three strategies use a simple fusion method to combine models, Meta-Classifier has parameters to train, which attempts to learn the collective knowledge represented by base classifiers. As for choosing estimator for meta-classifier, we test with a wide range of learning algorithms including not only the ones mentioned in Section ??, such as random forest, GBDT and XGBoost. It turns out Gaussian Naive Bayes is the most competitive model, which will be the only meta classifier discussed in next Section.

## 4 Experiments

### 4.1 Experimental setup

We use scikit-learn library<sup>1</sup> for implementation of the n-gram features based models and the ensemble meta-classifier. As for deep models, we implement them using Keras<sup>2</sup> library with Tensorflow backend. We used Adam (?) method as the optimizer, setting the first momentum to be 0.9, the second momentum 0.999 and the initial learning 0.001. The batch size is 32. All hidden states of LSTMs, feature maps of CNNs and word embeddings have 300 dimensions. Word embeddings are fine tuned during training process. All the models are trained separately on traditional and simplified version's dataset, and evaluated using macro-weighted f1 score. Our code for all experiments is publicly available<sup>3</sup>.

### 4.2 Contribution of single n-gram features

To find the most contributing individual n-gram features for discriminating Mandarin Chinese varieties. We run a number of experiments with the three classifiers using one single n-gram at a time, and the results are illustrated in Figure ??. As we can see, for dataset of both simplified and traditional version, performances of 3 models all drop sharply when using a n-gram of larger size, especially for word level n-grams. The most contributing character level ngram is character trigram, to which character bigram is inferior yet

close. Word unigram is the best between word level n-grams, but no better than character bigram or trigram. Among all 3 models, MNB outperforms LR and SVM, despite the trend that SVM has been the most preferred method for dialect identification. Besides, we can also see from the figure that performances on dataset of traditional version are slightly better than simplified version, which is persist with different models.

### 4.3 Combination of n-gram features

Table ?? shows the results of combining individual feature on each dataset. The performances of individual feature are also listed for direct comparison. As indicated from the table, feature combination does bring a performance gain. While 3 kinds of combination achieve a close score, MNB using combination of character bigram and trigram is the best, achieving macro-weighted f1 score of 0.9080 for simplified version dataset and 0.9225 for traditional version.

### 4.4 Performance of deep learning models

To fully compare deep learning methods with machine learning methods for DMT task, 7 deep models are evaluated. Results are listed in Table ??. Among 7 deep models, BiLSTM stands out from the others with macro-weighted score of 0.9000 and 0.9115. All deep models outperform LR and SVM, but are inferior to MNB, which shows again MNB is a very strong classifier for discriminating variations of Mandarin Chinese.

### 4.5 Performance of ensemble models

We also try to achieve a better result by aggregating the outputs of models we have implemented. As presented in Table ??, no single ensemble strategy performs consistently better than the others. The best choice for ensemble model is using MNB and BiLSTM as base classifier, and Mean Probability or Highest Confidence as fusion method. (When there are only 2 base classifiers, results of Mean Probability and Highest Confidence are always the same.)

### 4.6 Results of shared task

To do.

## 5 Conclusion and feature work

To do.

<sup>1</sup><https://scikit-learn.org/stable/>

<sup>2</sup><https://github.com/keras-team/keras>

<sup>3</sup><https://github.com/AlexYangLi/DMT>

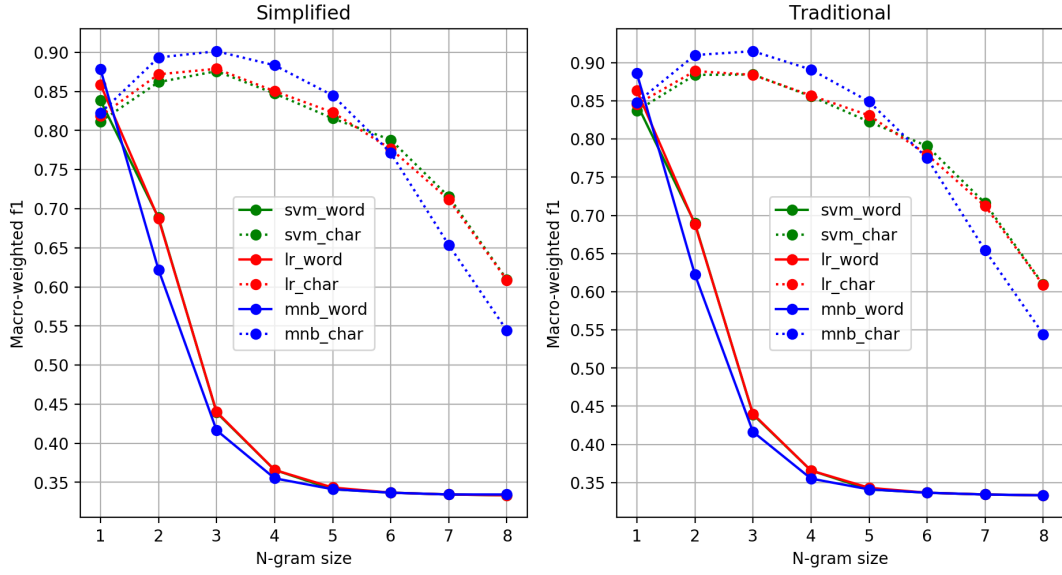


Figure 2: Macro-weighted f1 scores for LR (red lines), SVM (green lines), MNB (blue lines) using character (solid lines) or word level (dotted lines) n-gram of different sizes as input, both on dataset of simplified (left) and traditional (right) version.

Feature	Simplified			Traditional		
	LR	SVM	MNB	LR	SVM	MNB
<b>Individual feature</b>						
word uigram	0.8590	0.8384	0.8784	0.8634	0.8460	0.8860
char bigram	0.8720	0.8620	0.8935	0.8890	0.8840	0.9100
char trigram	0.8790	0.8760	0.9015	0.8840	0.8845	0.9150
char 4gram	0.8504	0.8474	0.8835	0.8570	0.8559	0.8910
<b>Combined feature</b>						
char bigram+trigram	0.8865	0.8830	<b>0.9080</b>	0.8960	0.8925	<b>0.9225</b>
char bigram+trigram+4gram	0.8880	0.8835	0.9030	0.8945	0.8920	0.9170
char bigram+char trigram+word unigram	0.8875	0.8835	0.9055	0.8990	0.8940	0.9200

Table 3: Macro-weighted f1 scores for LR, SVM, MNB using individual or combined features as input, both on dataset of simplified and traditional version.

Model	Simplified	Traditional
<b>CNN-based</b>		
CNN	0.8964	0.9090
DCNN	0.8970	0.9080
DPCNN	0.8925	0.9070
<b>RNN-based</b>		
BiLSTM	<b>0.9000</b>	<b>0.9115</b>
Self-attentive BiLSTM	0.8915	0.9020
<b>CNN-RNN hybrid</b>		
CNN-BiLSTM	0.8935	0.9080
BiLSTM-CNN	0.8950	0.9095

Table 4: Macro-weighted f1 scores for deep models using word embeddings as input, both on dataset of simplified and traditional version.

Ensemble Strategy	Simplified			Traditional		
	all ML*	all DL*	MNB +	all ML*	all DL*	MNB +
			BiLSTM			BiLSTM
Mean Probability	0.9025	0.9050	<b>0.9130</b>	0.9170	0.9215	<b>0.9240</b>
Highest Confidence	0.9080	0.9015	<b>0.9130</b>	0.9225	0.9100	<b>0.9240</b>
Majority Voting	0.8880	0.9060	-	0.8985	0.9195	-
Meta-Classifer	0.8915	0.9025	0.9050	0.906	0.9130	0.9215

Table 5: Macro-weighted f1 scores for 4 ensemble strategies combining different base classifiers, both on dataset of simplified and traditional version. “all ML” and “all DL” refer to combine all machine learning models and deep learning models respectively. All machine learning models use character bigram-trigram combination as input.

## References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. [Neural machine translation by jointly learning to align and translate](#). *CoRR*, abs/1409.0473.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Peng Chen, Zhongqian Sun, Lidong Bing, and Wei Yang. 2017a. [Recurrent attention network on memory for aspect sentiment analysis](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 452–461. Association for Computational Linguistics.
- Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana Inkpen. 2017b. [Enhanced lstm for natural language inference](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1657–1668. Association for Computational Linguistics.
- Alina Maria Ciobanu and Liviu P. Dinu. 2016. [A computational perspective on the romanian dialects](#). In *Proceedings of Language Resources and Evaluation (LREC)*.
- Simon Clematide and Peter Makarov. 2017. [Cluzh at vardial gdi 2017: Testing a variety of machine learning tools for the classification of swiss german dialects](#). In *Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*, pages 170–177. Association for Computational Linguistics.
- Çağrı Çöltekin, Taraka Rama, and Verena Blaschke. 2018. [Tübingen-oslo team at the vardial 2018 evaluation campaign: An analysis of n-gram features in language variety identification](#). In *Proceedings of the Fifth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2018)*, pages 55–65. Association for Computational Linguistics.
- Minwei Feng, Bing Xiang, Michael R. Glass, Lidan Wang, and Bowen Zhou. 2015. [Applying deep learning to answer selection: A study and an open task](#). *CoRR*, abs/1508.01585.
- Chu-Ren Huang and Lung-Hao Lee. 2008. [Contrastive approach towards text source classification based on top-bag-of-word similarity](#). In *Proceedings of the 22nd Pacific Asia Conference on Language, Information and Computation*.
- Rie Johnson and Tong Zhang. 2017. [Deep pyramid convolutional neural networks for text categorization](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 562–570. Association for Computational Linguistics.
- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. [A convolutional neural network for modelling sentences](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 655–665. Association for Computational Linguistics.
- Yoon Kim. 2014. [Convolutional neural networks for sentence classification](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2014. [Adam: A method for stochastic optimization](#). *CoRR*, abs/1412.6980.
- Tim Kreutz and Walter Daelemans. 2018. [Exploring classifier combinations for language variety identification](#). In *Proceedings of the Fifth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2018)*, pages 191–198. Association for Computational Linguistics.
- Zhouhan Lin, Minwei Feng, Cícero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. [A structured self-attentive sentence embedding](#). *CoRR*, abs/1703.03130.



- Marco Lui and Paul Cook. 2013. [Classifying english documents by national dialect](#). In *Proceedings of the Australasian Language Technology Association Workshop 2013 (ALTA 2013)*, pages 5–15.
- Shervin Malmasi, Marcos Zampieri, Nikola Ljubešić, Preslav Nakov, Ahmed Ali, and Jörg Tiedemann. 2016. [Discriminating between similar languages and arabic dialect identification: A report on the third dsl shared task](#). In *Proceedings of the Third Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial3)*, pages 1–14. The COLING 2016 Organizing Committee.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#). *CoRR*, abs/1301.3781.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. [Glove: Global vectors for word representation](#). In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017a. [Attention is all you need](#). In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates, Inc.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017b. [Attention is all you need](#). *CoRR*, abs/1706.03762.
- Fan Xu, Mingwen Wang, and Maoxi Li. 2016. [Sentence-level dialects identification in the greater china region](#). *International Journal on Natural Language Computing (IJNLC)*, 5(6).
- Marcos Zampieri and Binyam Gebrekidan Gebre. 2012. [Automatic identification of language varieties: The case of portuguese](#). In *Proceedings of KONVENS*.
- Marcos Zampieri, Shervin Malmasi, Nikola Ljubešić, Preslav Nakov, Ahmed Ali, Jörg Tiedemann, Yves Scherrer, and Noëmi Aepli. 2017. [Findings of the vardial evaluation campaign 2017](#). In *Proceedings of the Fourth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial)*, pages 1–15. Association for Computational Linguistics.
- Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Ahmed Ali, Suwon Shon, James Glass, Yves Scherrer, Tanja Samardžić, Nikola Ljubešić, Jörg Tiedemann, Chris van der Lee, Stefan Grondelaers, Nelleke Oostdijk, Dirk Speelman, Antal van den Bosch, Ritesh Kumar, Bornini Lahiri, and Mayank Jain. 2018. [Language identification and morphosyntactic tagging: The second vardial evaluation campaign](#). In *Proceedings of the Fifth Workshop on NLP for Similar Languages, Varieties and Dialects (VarDial 2018)*, pages 1–17. Association for Computational Linguistics.
- Chunting Zhou, Chonglin Sun, Zhiyuan Liu, and Francis C. M. Lau. 2015. [A C-LSTM neural network for text classification](#). *CoRR*, abs/1511.08630.