

Towards Code-switched Classification Exploiting Constituent Language Resources

Kartikey Pant^{*} and Tanvi Dadu^{*}

International Institute of Information Technology, Hyderabad, India

Netaji Subhas Institute of Technology, New Delhi, India

kartikey.pant@research.iiit.ac.in

tanvid.co.16@nsit.net.in

Abstract

Code-switching is a commonly observed communicative phenomenon denoting a shift from one language to another within the same speech exchange. The analysis of code-switched data often becomes an assiduous task, owing to the limited availability of data. We propose converting code-switched data into its constituent high resource languages for exploiting both monolingual and cross-lingual settings in this work. This conversion allows us to utilize the higher resource availability for its constituent languages for multiple downstream tasks.

We perform experiments for two downstream tasks, sarcasm detection and hate speech detection, in the English-Hindi code-switched setting. These experiments show an increase in 22% and 42.5% in *F1-score* for sarcasm detection and hate speech detection, respectively, compared to the state-of-the-art.

1 Introduction

Code-switching is the juxtaposition of speech belonging to two different grammatical systems or subsystems within the same speech exchange. In other words, it denotes a shift from one language to another within the same utterance but not limited to beyond mere insertion of borrowed words, fillers, and phrases. It often involves morphological and grammatical mixing (Sitaram et al., 2019). Code-switching of speech is a common phenomenon in the multilingual communities (Hickey, 2012), especially among peers who have similar fluency in multiple languages. Code-switching among languages like Spanish-English, Hindi-English, Tamil-English is often seen to be used by bilingual speakers of these languages.

Code-switching is likely to be found in informal settings, including speech and social media,

and semi-formal such as navigation instructions. This phenomenon leads us to acknowledge code-switching as a legitimate communication form deserving careful linguistic analysis and processing. Processing code-switched communication would enhance user experience in various industrial settings, including advertising, healthcare, education, and entertainment (Sitaram et al., 2019).

Access to code-switched data is challenging and limited. This phenomenon makes the analysis and information extraction from code-switched languages a less explored and challenging task. Most code-switching studies focus on pairs of a high resource and a lower resource language, often making these studies data starved. Also, code-switched language is mostly used for informal and semi-formal settings, making it less archived, resulting in its limited availability in terms of context and volume.

Traditionally, most works in academia have limited themselves to language-independent methodologies like SVM, FastText, and CNN (Sitaram et al., 2019) for modeling code-switched languages. This can be attributed to the shortage of data in code-switched settings, limiting the exploitation of predictive performance achieved via transfer learning using pretrained contextualized word embeddings.

We propose a generic methodology for modeling Hindi-English code-switched data revolving around its translation and transliteration. We propose two approaches that exploit a) Transliteration to convert this low-resource task into a high-resource task using cross-lingual learning, and b) Translation to utilize monolingual resources of each of the constituent languages. The importance of this methodology is that we enable the exploitation of pretrained transformer-based embeddings that are available in the monolingual or cross-lingual setting but not code-mixed setting

^{*} Both authors contributed equally to the work.

due to the scarcity of data for pretraining. We then benchmark both the approaches for the tasks of hate speech detection and sarcasm detection using the current state-of-the-art and language-independent supervision learning as the baseline. We further provide insights into the proposed approaches and compare them analytically.

2 Related Works

The previous works on sarcasm and humor detection contain a myriad of methods employed over the years on multiple datasets, including statistical and N-gram analysis on spoken conversations from an American sitcom (Purandare and Litman, 2006), Word2Vec combined with KNN Human Centric Features on the *Pun of the Day* and *16000 One-Liner datasets* (Yang et al., 2015), and Convolutional Neural Networks on datasets with distinct joke types in both English and Chinese (Chen and Soo, 2018). Recently, transformer-based architectures have been used to detect humor on multiple datasets collected from Reddit, Kaggle, and Pun of the Day Dataset (Weller and Seppi, 2019).

The prior works exploring hate speech detection employ several machine learning-based methods like linguistic features with character n-grams (Waseem and Hovy, 2016), SVM with a set of features that includes n-grams, skip-grams, and clustering-based word representations (Malmasi and Zampieri, 2017). Neural networks like LSTM (Badjatiya et al., 2017), FastText (Badjatiya et al., 2017) and CNN (Gambäck and Sikdar, 2017) have also been used for hate speech detection. Recently, transformer-based architectures using BERT, ALBERT, and RoBERTa have been used to detect hate speech (Mozafari et al., 2019; Wiedemann et al., 2020).

There has been considerable work in code-switching, as is well documented by Sitaram et al. (2019). The use of code-switched data has been explored for the tasks of Part-of-Speech tagging (Vyas et al., 2014), Named Entity Recognition (Aguilar et al., 2018), Dependency Parsing (Pantunen et al., 2018), and Text Classification (Swami et al., 2018; Bohra et al., 2018).

Swami et al. (2018) introduced an English-Hindi code-switching dataset for sarcasm detection in social media text. Their dataset contains tweets annotated as sarcastic and non-sarcastic in the form of their constituent language-identified tokens. They further provide a feature-based model as a base-

line. We use their dataset for our sarcasm detection experiments.

Bohra et al. (2018) introduced a dataset of English-Hindi social media text, demonstrating code-switching for the task of hate speech detection. They also conduct benchmarking experiments using a feature-based model. Similar to Swami et al. (2018), their benchmarks only contain accuracies for their baseline systems and do not take F1-scores into account. We use their dataset for our hate speech detection experiments.

There have been attempts to exploit monolingual datasets for enhancing downstream tasks in a code-switched setting. Solorio and Liu (2008) explored the use of monolingual taggers for Spanish and English for Part-of-Speech Tagging in an English-Spanish code-switched dataset. They further show that their best results were obtained by using the output of monolingual POS taggers as input features. Gupta et al. (2018) explore using monolingual and bilingual resources for the task of Question Answering in an English-Hindi code-switched dataset.

For text classification involving code-switched text, traditional methodologies like SVM, FastText, CNN has been used (Sitaram et al., 2019). However, most works are unable to exploit the recent advances in pretrained models, including RoBERTa and XLM-RoBERTa (Liu et al., 2019; Ruder et al., 2019) due to data scarcity and lack of normalization in the code-switched setting. To overcome this, Srivastava and Vardhan (2020) attempt to use multilingual BERT for code-switched sentiment analysis of social media text using a transliteration based methodology. However, they failed to beat the traditional baselines, including fasttext.

3 Approach

This section outlines our approach for analyzing code-switched language by converting it into its respective high resource languages and fine-tuning monolingual and cross-lingual contextual word embeddings. We highlight the data preparation followed by the contextual embedding used for each of the variants. Figure 1 illustrates our proposed approach through an example code-switched sentence from one of the datasets.

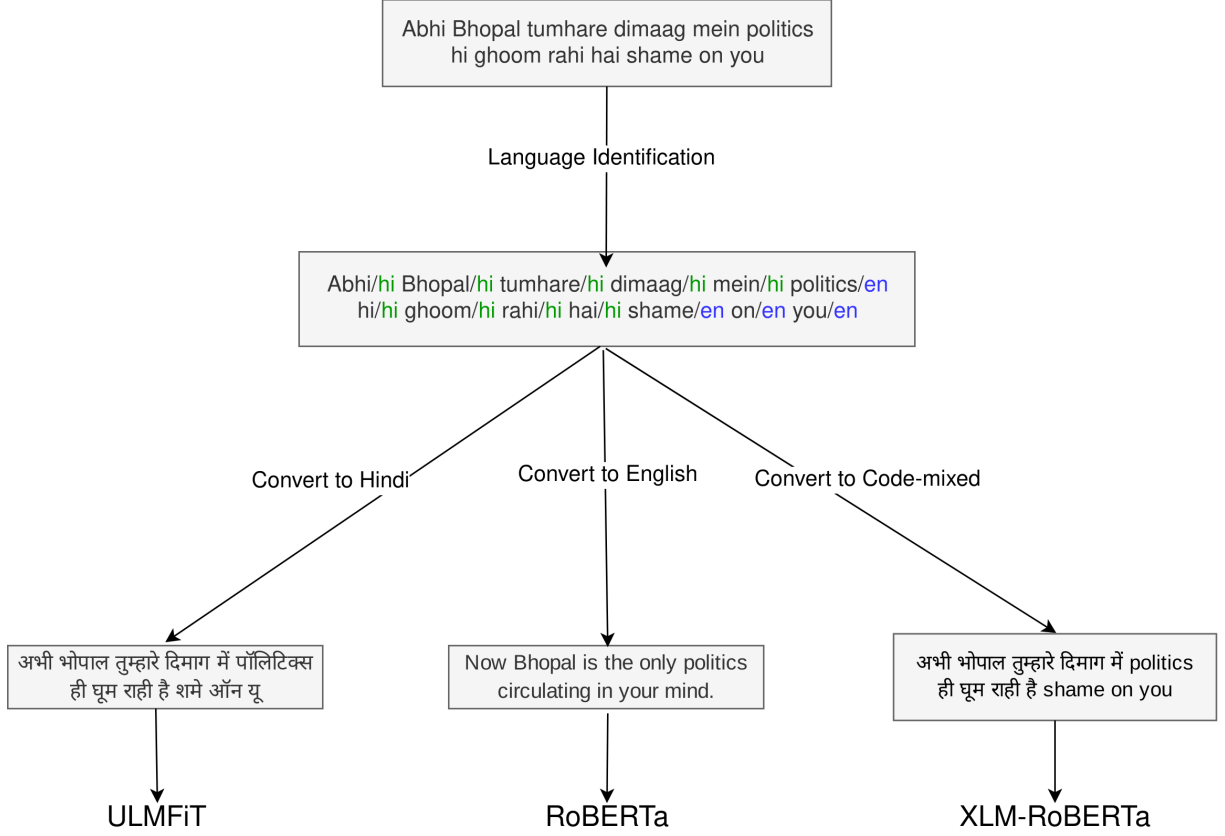


Figure 1: An example sentence demonstrating the proposed approach.

3.1 High Resource Language A - Hindi

3.1.1 Data Preparation

This section outlines our approach for analyzing code-switched language by converting it into its respective high resource languages and fine-tuning monolingual and cross-lingual contextual word embeddings. We highlight the data preparation followed by the contextual embedding used for each of the variants.

This section presents the steps performed for converting the code-switched data into one of the high resource languages, Hindi:

1. **Language Identification:** Each token was identified as either containing Hindi, English, or neither of them.
2. **Transliteration into Devanagari:** Each language-marked token, containing Hindi and English words, was transliterated from the Latin script to the Devanagari script. We use Microsoft Translate Text¹ service’s Transliteration API for the transliteration process.

¹<https://www.microsoft.com/en-us/translator/business/translator-api/>

Language Identification for the Hate Speech Detection task was done using a supervised FastText model trained on the Sarcasm Detection corpus (Swami et al., 2018) containing tokens with their respective languages- English, Hindi, and Other (includes Hashtags, Mentions, and Punctuation marks).

Dataset	Percentage of Hindi Tokens
Sarcasm Detection	88.259%
Hate Speech	86.437%

Table 1: Percentage of Hindi Tokens in comparison to Language-marked Tokens.

Table 1 shows that the majority of the words present in the dataset are Hindi, allowing us to assume that code-switched language contains Hindi and English words borrowed into Hindi. This phenomenon allows us to take advantage of converting all words from Latin to Devanagari via transliteration.

3.1.2 Classification Model

For the classification of the text translated to Hindi, we use Universal Language Model Fine-tuning for Text Classification, ULMFiT (Howard and Ruder,

2018). We use a pretrained language model trained on 172K Wikipedia articles in Hindi, attaining a perplexity of 34.06². We fine-tune our language model using discriminative fine-tuning, using different learning rates for different layers of the language model and slanted triangular learning rates. We then fine-tune the final model with the text translated in Hindi, after augmenting the trained language model with two additional linear blocks for classification. These blocks use batch normalization and dropout, with ReLU activations for the intermediate layer and a softmax activation for the last layer.

3.2 High Resource Language B - English

3.2.1 Data Preparation

This section outlines the conversion process of the code-switched data into the other high resource language, English:

1. **Language Identification:** Similar to the previous subsection, each token was identified to be either containing Hindi, English, or neither of them.
2. **Transliteration into Devanagari:** We used the same model as in the previous subsection for transliterating each language-marked token from the Latin script to the Devanagari script.
3. **Translation into English:** Each full sentence from the transliterated text was translated into English using the Translate API of the Microsoft Translate Text service.

3.2.2 Classification Model

For classification of the generated English text, we use a strong model for monolingual English classification, *RoBERTa_{Large}*. RoBERTa (Liu et al., 2019) is a contextualized word representation model, pre-trained using a bidirectional Transformer-based encoder. It is trained using significantly larger data with carefully tuned hyperparameters and performs competitively in benchmarking texts like GLUE (Wang et al., 2018) for text classification and SQuAD (Rajpurkar et al., 2016) for question answering. Since the model is pretrained on a large generic dataset, it can be fine-tuned for a specific task in a simpler yet efficient manner without making significant changes

in its architecture. We fine-tune its large variant, *RoBERTa_{Large}*, for both the tasks.

3.3 Cross-lingual Classification

3.3.1 Data Preparation

Each language-marked token needs to have its correct script to utilize cross-lingual word representations. We, therefore, employ the following steps to prepare the dataset for cross-lingual classification:

1. **Language Identification:** Similar to the previous subsections, each token was identified to be either containing Hindi, English, or neither of them.
2. **Transliteration of Hindi words into Devanagari:** We used the same model as in the previous subsections for transliterating each token marked to be Hindi from the Latin script to the Devanagari script.

3.3.2 Classification Model

For cross-lingual classification, we use XLM-RoBERTa, a state-of-the-art contextual word representation released by Facebook (Ruder et al., 2019). XLM-R is a cross-lingual unsupervised contextual word representation pretrained on a 100-language sized dataset. We exploit its advantage of massive cross-lingual transfer learning by fine-tuning it for both the tasks. Due to its high scalability and ease of use in fine-tuning, the model enables us to achieve competitive predictive performance while using minimum training resources. We use the large variant, XLM-R, for both the tasks.

4 Experimental Results

4.1 Downstream Tasks

This section introduces the datasets used for the downstream tasks of Sarcasm Detection and Hate Speech Detection.

For the task of sarcasm detection, we use English-Hindi code-switched dataset released by Swami et al. (2018). The dataset contains 5,250 tweets sampled from Twitter with sarcastic or non-sarcastic tags. It also provides language tags (English/Hindi/Others) for all the tokens present in tweets. The dataset has 0.790 Cohen’s Kappa score for measuring inter-annotator agreement, implying a high quality of the annotation schema.

For the task of hate speech detection, we use Bohra et al. (2018)’s dataset containing tweets showing English-Hindi code-switching. From their

²<https://github.com/goru001/nlp-for-hindi>

released dataset, we were able to get 4,575 tweets along with their binary label denoting whether the tweet contains hate speech or not. We preprocess the dataset and remove hyperlinks embedded in the tweets after language identification of each token. The authors measure the inter-annotator agreement through Kohen’s Kappa score, which is reasonably high, with a value of 0.982.

4.2 Baselines

We use the following neural architectures as baselines for our experiments:

1. FastText (Joulin et al., 2016): FastText is a library released by Facebook which uses bag of words and bag of n-grams as features for text classification. It relies on capturing partial information about the local word order efficiently.
2. Convolutional Neural Networks (CNN) (Kim, 2014): Convolutional neural networks are multistage trainable neural network architectures developed for classification tasks employing convolution and pooling for extracting features from the text.

4.3 Experimental Settings

In this subsection, we outline the experimental settings used for each of the models used in the experiment. We evaluated our model on a held-out test dataset for all experiments, consisting of 10% of the total dataset. For validation purposes, we further split our training dataset using a 90 : 10 train-validation split. We evaluate all the models on the following metrics: *F1*, *Precision*, *Recall*, *Accuracy*.

We use FastText’s recently open-sourced automatic hyperparameter optimization functionality and run 100 trials of optimization. We use a two-dense-layered architecture for Convolutional Neural Network with 3 and 4 convolution layers for sarcasm and hate speech detection, respectively. We further used a sequence length of 100 and an embedding size of 300 with a dropout rate of 0.1 and 0.2.

For ULMFiT, we use a pretrained language model trained on 172K Wikipedia articles in Hindi. We use the SentencePiece tokenizer for tokenizing the texts. For language modeling, we used a batch size of 16 and BPTT of 70. We use AWD-LSTM architecture with a dropout rate of 0.5 as a classifier. For *RoBERTa_{Large}* and XLM-R, we fine-tune

with a learning rate of $1 * 10^{-5}$ for 3 epochs, each with a maximum sequence length of 50.

4.4 Results

Table 2 and Table 3 show the experimental results obtained for the sarcasm detection dataset and hate speech dataset, respectively. We observe that our proposed approach of converting code-switched data into respective high resource language outperforms previously used approaches like FastText and CNN significantly. For both sarcasm detection and hate speech detection, preparing code-switched data for cross-lingual classification and using it to fine-tune XLM-R gives the best results. This is reflected in the model obtaining an *F1-score* of 0.850 for sarcasm detection and 0.724 for hate speech detection. For both tasks, respectively, it is 22% and 42.5% higher than CNN, the best performing baseline.

Model	Accuracy	Precision	Recall	F1
FastText	76.22%	0.245	0.951	0.390
CNN	95.71%	0.806	0.610	0.694
ULMFiT	94.85%	0.800	0.733	0.765
<i>RoBERTa_{Large}</i>	95.99%	0.791	0.883	0.835
XLM-R	96.37%	0.806	0.900	0.850

Table 2: Experimental results for the sarcasm detection task.

Model	Accuracy	Precision	Recall	F1
FastText	68.78%	0.609	0.417	0.495
CNN	58.08%	0.411	0.664	0.508
ULMFiT	68.78%	0.545	0.506	0.525
<i>RoBERTa_{Large}</i>	71.62%	0.707	0.716	0.711
XLM-R	71.83%	0.730	0.718	0.724

Table 3: Experimental results for the hate speech detection task.

Further we observe that our proposed approach performs significantly well for all the metrics taken under consideration with stark increase in *F1-score*. In sarcasm detection, it outperforms *F1-score*, *Precision* where as in hate speech detection, it outperforms *F1-score*, *Precision*, *Recall* and *Accuracy* by a significant margin. In sarcasm detection, a slight increase in *Accuracy* can be attributed to the imbalance distribution of classes, which is reinforced by a significant increase in *F1-score*.

The experimental results obtained for both sarcasm and hate speech detection show our approach’s effectiveness in leveraging the pretrained contextualized word embeddings for the code-switched settings. It also successfully tackles data

shortage when dealing with code-switched data in pretrained settings by converting code-switched language into its constituent languages.

5 Conclusion

In this work, we show that converting code-switched data into its constituent high resource language enables us to exploit the performance of the pretrained models for those languages on downstream tasks. We perform experiments on the English-Hindi code-switching pair and benchmark our approach for sarcasm detection and hate speech detection. The experiments show an improvement of up to 42.5% in *F1-score* compared to strong baselines like CNN. Our findings pave the way for further research to utilize monolingual resources for code-switched data for multiple downstream tasks and extend this methodology for other code-switched language pairs.

Acknowledgments

We would like to thank Sai Krishna Rallabandi for reviewing the drafts and the mentoring. We would also like to thank the anonymous reviewers for providing critical suggestions.

References

- Gustavo Aguilar, Fahad AlGhamdi, Victor Soto, Mona Diab, Julia Hirschberg, and Tamar Solorio. 2018. [Named entity recognition on code-switched data: Overview of the CALCS 2018 shared task](#). In *Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching*, pages 138–147, Melbourne, Australia. Association for Computational Linguistics.
- Pinkesh Badjatiya, Shashank Gupta, Manish Gupta, and Vasudeva Varma. 2017. [Deep learning for hate speech detection in tweets](#). In *Proceedings of the 26th International Conference on World Wide Web Companion*, WWW '17 Companion, page 759–760, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
- Aditya Bohra, Deepanshu Vijay, Vinay Singh, Syed Sarfaraz Akhtar, and Manish Shrivastava. 2018. [A dataset of Hindi-English code-mixed social media text for hate speech detection](#). In *Proceedings of the Second Workshop on Computational Modeling of People's Opinions, Personality, and Emotions in Social Media*, pages 36–41, New Orleans, Louisiana, USA. Association for Computational Linguistics.
- Peng-Yu Chen and Von-Wun Soo. 2018. [Humor recognition using deep learning](#). In *NAACL-HLT (2)*. Association for Computational Linguistics.
- Björn Gambäck and Utpal Kumar Sikdar. 2017. [Using convolutional neural networks to classify hate-speech](#). In *Proceedings of the First Workshop on Abusive Language Online*, pages 85–90, Vancouver, BC, Canada. Association for Computational Linguistics.
- Vishal Gupta, Manoj Chinnakotla, and Manish Shrivastava. 2018. [Transliteration better than translation? answering code-mixed questions over a knowledge base](#). In *Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching*, pages 39–50, Melbourne, Australia. Association for Computational Linguistics.
- Raymond Hickey. 2012. *The Handbook of Language Contact*, volume 1. 'John Wiley & Sons'.
- Jeremy Howard and Sebastian Ruder. 2018. [Universal language model fine-tuning for text classification](#). In *ACL*. Association for Computational Linguistics.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. [Bag of tricks for efficient text classification](#). Cite arxiv:1607.01759.
- 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, Doha, Qatar. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.
- Shervin Malmasi and Marcos Zampieri. 2017. [Detecting hate speech in social media](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 467–472, Varna, Bulgaria. INCOMA Ltd.
- Marzieh Mozafari, Reza Farahbakhsh, and Noel Crespi. 2019. A bert-based transfer learning approach for hate speech detection in online social media. In *COMPLEX NETWORKS*.
- Niko Partanen, Kyungtae Lim, Michael Rießler, and Thierry Poibeau. 2018. [Dependency parsing of code-switching data with cross-lingual feature representations](#). In *Proceedings of the Fourth International Workshop on Computational Linguistics of Uralic Languages*, pages 1–17, Helsinki, Finland. Association for Computational Linguistics.
- Amruta Purandare and Diane Litman. 2006. [Humor: Prosody analysis and automatic recognition for F*R*I*E*N*D*S*](#). In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*, pages 208–215, Sydney, Australia. Association for Computational Linguistics.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Sebastian Ruder, Anders Søgaard, and Ivan Vulić. 2019. [Unsupervised cross-lingual representation learning](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts*, pages 31–38, Florence, Italy. Association for Computational Linguistics.
- Sunayana Sitaram, Khyathi Raghavi Chandu, Sai Krishna Rallabandi, and Alan W. Black. 2019. [A survey of code-switched speech and language processing](#). *CoRR*, abs/1904.00784.
- Thamar Solorio and Yang Liu. 2008. Part-of-speech tagging for english-spanish code-switched text. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '08*, page 1051–1060, USA. Association for Computational Linguistics.
- Aditya Srivastava and V. Harsha Vardhan. 2020. [Hcms at semeval-2020 task 9: A neural approach to sentiment analysis for code-mixed texts](#).
- Sahil Swami, Ankush Khandelwal, Vinay Singh, Syed Sarfaraz Akhtar, and Manish Shrivastava. 2018. A corpus of english-hindi code-mixed tweets for sarcasm detection. *ArXiv*, abs/1805.11869.
- Yogarshi Vyas, Spandana Gella, Jatin Sharma, Kalika Bali, and Monojit Choudhury. 2014. [POS tagging of English-Hindi code-mixed social media content](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 974–979, Doha, Qatar. Association for Computational Linguistics.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Zeera Waseem and Dirk Hovy. 2016. [Hateful symbols or hateful people? predictive features for hate speech detection on twitter](#). In *Proceedings of the NAACL Student Research Workshop*, pages 88–93, San Diego, California. Association for Computational Linguistics.
- Orion Weller and Kevin D. Seppi. 2019. Humor detection: A transformer gets the last laugh. In *EMNLP/IJCNLP*.
- Gregor Wiedemann, Seid Muhie Yimam, and Chris Biemann. 2020. [UHH-LT & LT2 at semeval-2020 task 12: Fine-tuning of pre-trained transformer networks for offensive language detection](#). *CoRR*, abs/2004.11493.
- Diyi Yang, Alon Lavie, Chris Dyer, and Eduard H. Hovy. 2015. [Humor recognition and humor anchor extraction](#). In *EMNLP*. The Association for Computational Linguistics.