

CalBERT – Code-mixed Adaptive Language representations using BERT

Anonymous ACL submission

Abstract

A code-mixed language is a type of language that involves the combination of two or more language varieties in its script or speech. Code-mixed language has become increasingly prevalent in recent times, especially on social media. However, the exponential increase in the usage of code-mixed language, especially in a country like India which is linguistically diverse has led to various inconsistencies. Analysis of text is now becoming harder to tackle because the language present is not consistent and does not work with pre-defined existing models which are monolingual. We propose a novel approach to improve performance in Transformers by introducing an additional step called "Siamese Pre-Training", which allows pre-trained monolingual Transformers to adapt language representations for code-mixed languages with a few examples of code-mixed data. Our studies show that CalBERT is able to improve performance over existing pre-trained Transformer architectures on downstream tasks such as sentiment analysis. Multiple CalBERT architectures beat the state of the art F_1 -score on the Sentiment Analysis for Indian Languages (SAIL) dataset, with the highest possible improvement being 5.1 points. CalBERT also achieves the state-of-the-art accuracy on the IndicGLUE Product Reviews dataset by beating the benchmark by 0.4 points.

1 Introduction

Code-mixed language is a form of language wherein syntactic elements from one language are inserted into another language in such a way that the semantics of the resultant language remain the same. Code-mixed language is more prevalent in a multilingual society like India, where most of the population is at least bi-lingual. Code-mixed language is most prevalent on social media platforms such as Facebook and Twitter. Given the interest of

enterprises in determining business insights from social media posts, generating a mechanism for the proper analysis, and understanding of code-mixed language gains even more importance.

Language representations are used in Natural Language Understanding tasks such as sentiment analysis and human-like conversation systems. The current state of the art methods for learning representations use Transformer architectures pre-trained on vast amounts of natural language data. However, almost all of these learnt representations are monolingual and have been created using a single language only, or pre-trained on multiple languages individually. These representations struggle when a language might be code-mixed and hence display low performance on code-mixed natural language tasks due to inherent inconsistencies and large variability in data.

Our novel approach generates Code-mixed Adaptive Language representations using Bidirectional Encoder Representations from Transformers (CalBERT) by adding an additional pre-training step called "Siamese Pre-Training" where a pre-trained monolingual Transformer (Vaswani et al., 2017) is provided with different representations of the same sentence or phrase and learns to minimise the distance between their embeddings. We evaluate CalBERT on sentiment analysis of the Hinglish language using the benchmark SAIL 2017 dataset (Das and Gambäck, 2014) and obtain an F_1 -score of 62%, thus obtaining an improvement of 5.1 points, or 8.9%. We also evaluate CalBERT on the sentiment analysis task released by IIT Patna using the IndicGLUE Product Reviews dataset (Kakwani et al., 2020), and obtain an accuracy of 79.37%, a 0.5% increase over the existing benchmark.

2 Background

BERT (Devlin et al., 2018) and similar transformer architectures derived from BERT (such as RoBERTa (Liu et al., 2019), DistilBERT (Sanh et al., 2019) and others) are used to extract language representations from text corpora. These language representations are learned using a bidirectional attention mechanism (Bahdanau et al., 2014), and incorporate contextual information about the individual tokens in the corpus, and can be fine-tuned for a variety of tasks. The large majority of existing models are trained on monolingual corpora, and as such, produce representations that are more attuned for high performance in tasks involving a single language. As such, these representations suffer from poor performance when applied to tasks involving code-mixed language that contains multiple language varieties in a single script language.

The methodology proposed in this work seeks to adapt existing representations for monolingual text into representations that can be fine-tuned for code-mixed tasks. This allows for representations of code-mixed language to be generated without having to pre-train a Transformer model from scratch on large quantities of code-mixed data, which is a time consuming and computationally intensive process.

3 Previous Work

Mikolov, in his paper titled “Efficient Estimation of Word Representations in Vector Space” (Mikolov et al., 2013), proposes an novel approach to compute word embeddings without the added complexity of a hidden layer while performing better than neural network language model. These word vectors outperform the SemEval 2012 task-2 benchmark however perform poorly on out of vocabulary words and morphologically rich languages.

Reimers et al, in their work titled “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks” (Reimers and Gurevych, 2019) modify the existing BERT architecture to use a Siamese network with shared weights and cosine similarity as the loss function to predict the semantic textual similarity. The model was able to train in linear time yields a score of 84.88% in sentiment prediction of movie reviews.

“A Passage to India: Pre-trained Word Embeddings for Indian Languages” (Kumar et al., 2020) by Kumar et al, shows that models trained on sub word representations perform better as Indian lan-

guages are morphologically rich. Their multilingual embeddings when evaluated on next sentence prediction using pre-trained BERT gives 67.9% accuracy. Another paper titled “Towards Sub-Word Level Compositions for Sentiment Analysis of Hindi-English Code Mixed Text” (Joshi et al., 2016) by Joshi et al, discusses the advantage of using sub word representations compared to character level tokens to deal with inconsistencies in code-mixed text. This method was evaluated on a custom dataset and yields an F_1 -score of 65.8%.

Choudhary et al (Choudhary et al.), in their paper “Sentiment Analysis of Code-Mixed Languages leveraging Resource Rich Languages” propose a novel method which uses contrastive learning. They use a Siamese model to map code-mixed and monolingual language to the same space by clustering based on similarity between their skip-gram vectors. This works on the hypothesis that similar words have similar contexts. The model achieves 75.9% F-score on the HECM dataset, an improvement over existing approaches.

4 Proposed Approach

We propose a novel approach to adapt Transformer representations for a code-mixed language from existing representations that exist in another language by introducing an additional pre-training step using a Siamese network architecture (Bromley et al., 1993). We attempt to minimise the distance between the two semantic spaces of the corresponding languages and obtain a joint or shared representation for equivalent sentences in both languages. This additionally enables the generation of code-mixed language representations from an existing language’s representations, without the need to pre-train a Transformer from scratch.

Our novel approach is implemented by using a shared pre-trained Transformer layer in each branch of the Siamese network and use an appropriate loss function to bring the embeddings closer between each pair of sentences. To Siamese-pretrain CalBERT, the same sentence needs to be obtained in both the language that the Transformer was pre-trained in, along with the language for which the Transformer is trying to adapt representations.

5 Methodology

5.1 Terminologies

- *Base Language*: The single language that was used to pre-train the Transformer architec-

ture. Could be any task like Masked Language Modelling (Devlin et al., 2018), Next Sentence Prediction and so on. For example, in the case of the Hinglish language the base language is English.

- *Target Language*: The code-mixed language for which the Transformer architecture is trying to adapt representations. This language is a super-set of the base language, since it may contain sentences entirely in the base language. For example, Hinglish.
- *Siamese Pre-training*: The novel additional pre-training step proposed to adapt representations.

5.2 Siamese Pre-training

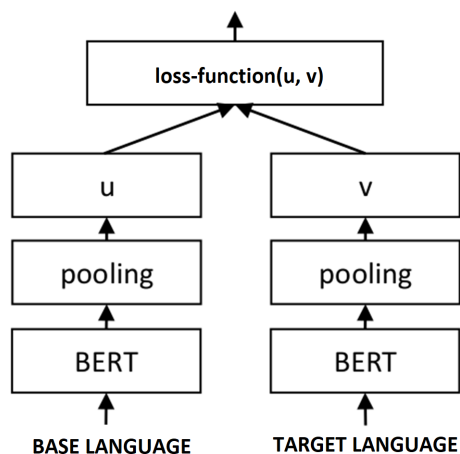


Figure 1: SBERT Architecture used for Siamese Pre-training

A Siamese network is a network consisting of two or more branches, wherein each branch receives an input. The network learns to distinguish between the examples provided through each branch. Siamese networks have been used extensively in computer vision for image classification.

A pre-trained monolingual Transformer being trained on only a single language is capable of generating language representations for that language only, thus performing poorly when the language is code-mixed. Pre-training a Transformer from scratch on code-mixed data is a difficult task, since it is computationally expensive and also time consuming. Since the Transformer has learnt language representations for one of the languages, it is far

more optimal to adapt these existing representations for new and similar words belonging to the other language.

A pre-trained monolingual Transformer is provided with different representations of the same sentence or phrase and learns to minimize the distance between their embeddings (Fig. 1). The two representations in this case are the transliterated version of the code-mixed sentence in the target language, and the translated version in the base language. Since the two versions have the same semantic meaning, their similarity should ideally be as high as possible and thus the distance between their representations should be minimum. Since the Transformer already knows representations for the base language, it only needs to adapt representations to map the target language to the same space.

The Siamese pre-training is called thus because of its use of the Siamese network architecture, wherein the two arms of the network are fed with two semantically equivalent input sentences in the base and the target language respectively. The network learns the embeddings for both sentences by comparing their similarity (normally, this involves the use of a labelled dataset with sentence pairs and the corresponding similarities, but here all the sentence pairs have maximum similarity). The loss function (Eqn. 1 and ??) used here is used to bring the representations from the two branches closer. In our work we use the cosine distance as the loss function, although similar loss functions such as contrastive loss can also be used. Minimizing the cosine distance implies that the similarity is maximized between the two sentence embeddings, which is the desired outcome.

The Siamese pre-training is performed as the last pre-training step (after the usual pre-training strategies used in a Transformer such as masked language modelling or next sentence prediction), since it needs existing base language representations to learn the target language effectively. Additionally, our work shows that a significant amount of data is not required to effectively perform Siamese pre-training and that the model is able to learn with a fraction of the data size that was used for pre-training, thus showing that our approach does not require vast computational resources to improve performance.

A variety of BERT-based architectures were pre-trained and fine-tuned for the purposes of this

work. The models that were evaluated were BERT, RoBERTa and DistilBERT. These three architectures are used either pre-trained on an English corpus (as available publicly on the HuggingFace Hub), or pre-trained on the corpus of code-mixed data that was collected as part of this experiment.

6 Workflow

We focus our efforts on Transformer architectures based on the Bidirectional Encoder Representations from Transformers (BERT) architecture, since they are bidirectional models capable of learning representations from a given sequence using both forward as well as backward context. We demonstrate our novel approach on the Hindi-English (Hinglish) code-mixed language and data for the same was obtained from social media and news articles to maintain a good balance of well structured as well as informal code-mixed language usage.

6.1 Equations

The objective is to minimize the distance between the representation in the base language and the representation in the target language. The loss function, hence, needs to reflect this minimization of the distance between the two representations.

The cosine distance loss function in (1) represents the angle between two vectors. Minimizing the cosine distance implies that the two vectors are highly similar as a smaller angle between the two vectors creates a smaller cosine distance.

$$l(\vec{a}, \vec{b}) = \frac{\vec{a} \cdot \vec{b}}{\|\vec{a}\| \|\vec{b}\|} \quad (1)$$

6.2 Data Collection

Since code-mixed Hinglish data is abundantly present on platforms like social media, we choose to use social media as our source of data. Platforms like Twitter and Facebook were scraped and compiled together. There are also several online archives available which host code-mixed data for several tasks. Additionally, Hinglish code-mixed data is already available on IndicCorp, which is one of the largest publicly-available corpora for Indian languages.

After compiling all the sources of data together, they need to be converted into a suitable pairwise format to be used along with CalBERT.

6.3 Data Preprocessing

Although a large amount of data is obtained (Table 1), most of it is not in a format that can be directly used to train CalBERT. Since Hinglish code-mixed language can exist in both the Hindi (Devanagiri) script as well as Roman script, we need to transliterate all of them to the same single script language. For our work, we choose the script language to be Roman since most popular Transformer architectures have been pre-trained in this script language, and hence have representations for the base language, in our case English.

For code-mixed data that already exists in the Roman script, we need to obtain the translated version of the same for the base language representation. This has been performed by using software automation tools such as Selenium combined with Google Translate. Alternatively, for code-mixed data that exists in the Devanagiri script, we employ the use of a transliteration library and convert them into the ITRANS ASCII transliteration scheme.

The input to CalBERT consists of sentence pairs (Table 2). Each pair consists of the transliteration of the code-mixed sentence in the target language, and the translation of the same code-mixed sentence in the base language. Both sentences are represented in the Roman script in our work. Since both sentences have the same semantic meaning, the objective is to reduce the distance between the corresponding sentence representations. The shared Transformer layer in CalBERT thus allows it to learn joint representations for both languages and adapts base language representations to the target language.

Due to computational limitations, we did not Siamese pre-train our models on the data obtained from IndicNLP (Kakwani et al., 2020). However, our experiments show that the small portion of data obtained via scraping was able to boost performance significantly.

Source	Number of sentences
Social Media	147731
IndicNLP	8466307

Table 1: CalBERT Dataset Metrics

6.4 Evaluation Metric

Since CalBERT is trained in a task-agnostic manner, it can be fine-tuned on any downstream natural language understanding task like sentiment analy-

Base Language	Target Language
in reply, pakistan got off to a solid start	jisake javaba mem pakistan ne achchi shuruata ki thi.
by this only, we can build our country and make it great again	isake jariye hi hama desha ka nirmana karemge aura use mahana bana payemge
people started gathering	logom ki bhida jama hone lagi
obtain loans from a bank or an individual	kisi bank ya vyakti se rrina prapta kara sakati hai
he was later taken to a hospital and treated	bada mem use aspatala le jaya gaya aura ilaja kiya gaya
it's our cultural heritage	yaha hamari samskritika virasata hai

Table 2: Base and target language pairs used to train CalBERT

sis and question answering. We evaluate CalBERT on the Sentiment Analysis for Indian Languages (SAIL) 2017 dataset, which consists of Hinglish code-mixed data in the Roman script.

The dataset serves as a good evaluation strategy, since it has been compiled with data from various sources ranging from news articles to even social media and consists of various code-mixed forms for the same word and is often considered a difficult dataset to perform well on. It also serves as a benchmark dataset for sentiment analysis on Hindi-English and Bengali-English data, with the highest possible F_1 -score achieved on the dataset being 56.9%.

We also evaluate CalBERT on the IndicGLUE Product Reviews dataset released by IIT Patna, which serves as another benchmark dataset for sentiment analysis of Hinglish text, but in the Hindi (Devanagiri) script. All models evaluated on this script have either been trained on English scripted text, or on Devanagiri scripted text. We however propose to evaluate CalBERT on a code-mixed version of the dataset by transliterating the text from Devanagiri to Roman.

The dataset is popularly used for evaluating the performance of Transformers built for Indian languages such as IndicBERT and consists of reviews for products across various categories. The highest possible F_1 -score obtained on this dataset by

IndicBERT is 71.32%.

6.5 Siamese Pre-training CalBERT

To Siamese pre-train a Transformer (Table 3), we first initialise a Siamese network with a shared layer containing the Transformer whose representations we intend to adapt (BERT, RoBERTa or DistilBERT). This can be done effectively using a sentence-transformer architecture. We add a single pooling layer to each branch, and then finally combine the pieces together by adding a suitable loss function to reduce the distance between the representations. In our preliminary experiments, we observe that the contrastive and cosine distance loss functions perform nearly the same and hence use the cosine distance loss function for all the experiments performed. This shared Transformer layer can then be extracted and fine-tuned on other downstream tasks.

Hyperparameter	Value
Number of epochs	10
Number of warm-up steps	100
Learning Rate	2.5×10^{-5}
Weight decay	0.01
Optimizer	AdamW
Maximum gradient normalization	1.0

Table 3: Hyperparameters for Siamese Pre-training

7 Evaluating CalBERT

CalBERT is meant to be fine-tuned for downstream tasks involving code-mixed data. Due to the abundance of code-mixed Hinglish data available on social media and the much need for code-mixed Transformer architectures, we evaluate performance on the popular downstream task of sentiment analysis.

7.1 SAIL 2017

The Sentiment Analysis of Indian Languages (SAIL) 2017 dataset is a collection of sentences in two popular Indian code-mixed languages – Hindi-English and Bengali-English. The datasets are composed of sentences from various sources like news articles as well as social media and are in the Roman script. There is also a high degree of variability in the data, with multiple forms existing for the same word and different styles of writing. The sentences are classified into 3 polarities – positive,

neutral and negative. The SAIL 2017 task is a challenging benchmark, and is widely regarded as one of the benchmark datasets for sentiment analysis of the Hinglish language. The highest documented F_1 -score on the benchmark is 56.9%.

Dataset Split	Number of examples
Train set	9945
Test set	1238
Validation set	1240

Table 4: SAIL 2017 Dataset Metrics

The SAIL 2017 dataset is already partitioned into train, test and validation splits (Table 4). All comparisons are made using the F_1 score on the validation split. For our experiments, we fine-tune existing pre-trained models with and without CalBERT’s additional Siamese pretraining step and compare the score obtained by each model. The experiments are performed multiple times using the same set of hyperparameters and the highest F_1 -score over 10 runs is recorded.

Hyperparameter	Value
Learning Rate	2×10^{-6}
Training Batch Size	4
Evaluation Batch Size	4
Number of epochs	15
Weight Decay	0.08

Table 5: Hyperparameters for SAIL 2017 Fine-Tuning

CalBERT outperforms the SAIL 2017 benchmark (Table 6) F_1 -score by 5.1 points, or 8.9% with the CalBERT-XLM-RoBERTa model (XLM-RoBERTa with CalBERT’s Siamese pre-training). We also improve upon on the benchmark precision by 3.5% and the recall by 10.3%. Additionally, all other CalBERT architectures also outperform the benchmark F_1 -score, with the minimum improvement obtained by CalBERT-DistilBERT being 3%.

We observe that Transformer architectures also obtained improved performance metrics on the SAIL 2017 benchmark. For this reason we evaluate CalBERT’s Siamese pre-training against a native Transformer that was used to Siamese pre-train CalBERT (Table 7). Our experiments show that CalBERT has improved the F_1 -score on all native Transformer architectures as well, thus showing that additional pre-training can improve performance on the given code-mixed task.

Model Type	F_1 -Score	Precision	Recall
CalBERT-XLM-RoBERTa	0.620	0.618	0.618
CalBERT-RoBERTa	0.612	0.615	0.614
CalBERT-BERT	0.588	0.581	0.583
CalBERT-DistilBERT	0.586	0.587	0.586
CalBERT-IndicBERT	0.530	0.529	0.531
SAIL-2017 Benchmark	0.569	0.597	0.56

Table 6: Comparison of F_1 -scores on SAIL 2017 benchmark

Model Type	CalBERT	F_1 -Score
RoBERTa	✓	0.613
RoBERTa		0.608
BERT	✓	0.588
BERT		0.585
DistilBERT	✓	0.584
DistilBERT		0.580
XLM-RoBERTa	✓	0.620
XLM-RoBERTa		0.608
IndicBERT	✓	0.530
IndicBERT		0.544
SAIL-2017 Benchmark		0.569

Table 7: Influence of CalBERT on model F_1 -scores

Pre-training a Transformer from scratch is computationally expensive as well as time-consuming. Additionally, it requires a vast amount of data to effectively pre-train a Transformer to learn usable language representations. Since there are no code-mixed Transformer architectures that exist at the time of writing this paper, we experiment by pre-training some popular Transformers (Table 8) using Masked Language Modelling on a subset of our code-mixed data. The size of the subset taken is the same as that was used for CalBERT experiments to compare between the two approaches. Our experiments show that the models pre-trained from scratch do not outperform the benchmark, and are significantly lower than all CalBERT architectures, thus proving that Transformers do need enormous amounts of data to provide good results. Additionally, it reinforces our hypothesis that it is far more optimal as well as effective to apply CalBERT’s Siamese pre-training to adapt a pre-trained Transformer’s representations for another code-mixed language.

Model Type	F_1 -Score	Precision	Recall
SAIL-2017 Benchmark	0.569	0.597	0.56
DistilBERT-big	0.553	0.551	0.557
DistilBERT-small	0.551	0.552	0.556
BERT-small	0.551	0.550	0.554
BERT-big	0.543	0.542	0.547
RoBERTa-small	0.533	0.531	0.537
RoBERTa-big	0.524	0.523	0.529

Table 8: Comparison of F_1 -scores on code-mixed pre-trained Transformers with limited data

Table 9 showcases some predictions made by the various CalBERT model types. The outputs are in the form of integers, where a value of 0 indicates a negative sentiment, 1 indicates a neutral sentiment and 2 indicates a positive sentiment.

7.2 IndicGLUE Product Reviews

The IIT Patna product review dataset was released by the IndicNLP organization as part of the IndicGLUE collection of datasets that are meant to

Sentence	True Label	CalBERT-XLM-RoBERTa
sab ka Bhai meri Jan Salman khan	POSITIVE	POSITIVE
hahaha ! gazab imagination hai teri !	POSITIVE	POSITIVE
lagta hai aaj bhi bating nahi milega #indvssa	NEGATIVE	NEGATIVE
Or mastery kaisi chal ri hai apki	NEUTRAL	NEUTRAL

Table 9: CalBERT Results for SAIL 2017 Dataset

be used for evaluation of models trained for NLU tasks on Indian languages. The dataset consists of product reviews in Hindi taken from a popular e-commerce website. Like the SAIL dataset, there is a high degree of variability in this data too, with multiple forms existing for the same word and different styles of writing. The sentences are classified into 3 polarities – positive, neutral and negative.

Dataset Split	Number of examples
Train set	4182
Test set	523
Validation set	523

Table 10: IndicGLUE Dataset Metrics

Hyperparameter	Value
Learning Rate	2×10^{-6}
Training Batch Size	16
Evaluation Batch Size	16
Number of epochs	50
Weight Decay	0.02

Table 11: Hyperparameters for IndicGLUE Product Reviews Fine-Tuning

We observe that CalBERT is able to beat the state-of-the-art accuracy on the dataset by 0.4 points, or 0.5% with the CalBERT-XLM-RoBERTa model. We also improve on the score set by the IndicBERT model, by achieving an improvement of 8.05 points, or 11.2%. However, we observe that the other Transformer architectures do not perform well on this dataset, thus being consistent with previous attempts made by other authors (Kakwani et al., 2020).

Model Type	Accuracy
CalBERT-XLM-RoBERTa	0.794
IndicGLUE Benchmark	0.789
CalBERT-RoBERTa	0.639
CalBERT-BERT	0.612
CalBERT-IndicBERT	0.602
CalBERT-DistilBERT	0.564

Table 12: Comparison of accuracy on IndicGLUE Product Review Dataset

We also experiment with our custom pre-trained Transformers as used in the SAIL-2017 experiment on this dataset.(Table 13). As seen previously, the models again do not outperform the benchmark and perform very poorly. However, we observe that in certain cases, CalBERT out-performs the comparable from-scratch trained Transformer model, hence showing the effectiveness of Siamese pre-training as used in CalBERT, to learn language representations of code-mixed language.

Model Type	Accuracy
IndicGLUE Benchmark	0.789
IndicBERT	0.713
DistilBERT-big	0.671
DistilBERT-small	0.625
BERT-small	0.659
BERT-big	0.656
RoBERTa-small	0.627
RoBERTa-big	0.627

Table 13: Comparison of accuracy on code-mixed pre-trained Transformers with limited data

The marginal improvement in Table 7 can be attributed to the vast difference in the size of the dataset that was used to pre-train and Siamese pre-train. Due to computational limitations, CalBERT was only Siamese pre-trained on a 15% sample of the entire available data. However, the results are promising and they display the potential of CalBERT’s Siamese pretraining if it had been possible to pre-train on all the data available to us.

8 Conclusion

We demonstrate the use of BERT and BERT-based architectures in learning code-mixed language representations for Hindi-English code-mixed lan-

Sentence	True Label	CalBERT-XLM-RoBERTa
agar 3G network chahiye to baahar se dongle lagane padegi	NEGATIVE	NEUTRAL
ek haath se istemaal nahin kar sakte hain	NEGATIVE	POSITIVE
micromax canvas series ke nae tablet me 3000 mah ki battery lagi hain , jo kaphi lambe samay tak chalti hain	POSITIVE	NEUTRAL
is philm se vaha ummid dhul gai	NEGATIVE	NEUTRAL
isme HDMI out put hai jisse aap store ki gai philmo aur picture ko apani tv par dekh sakte hain	NEUTRAL	NEUTRAL

Table 14: CalBERT Results for IndicGLUE Product Review Dataset

guage, and evaluate the performance of the learned embeddings on a benchmark sentiment analysis task. We present a task and language-agnostic approach to generating cross-language representations for sentences, which can further be fine-tuned on any specific downstream task.

We show a 8.9% improvement in the F_1 score achieved by the novel Siamese pre-training method over the existing benchmark score. We also show that CalBERT also outperforms the native Transformer architectures which were used to Siamese pre-train CalBERT, thus showing that Siamese pre-training can help existing Transformers adapt to a code-mixed language.

Owing to computational limitations at our end, we are yet to find out the extent of possible improvement that CalBERT can show on the benchmark dataset, but we postulate that training on more examples may result in a more significant increase in the performance of the model.

References

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

800	Jane Bromley, James W Bentz, Léon Bottou, Isabelle Guyon, Yann LeCun, Cliff Moore, Eduard Säckinger, and Roopak Shah. 1993. Signature verification using a “siamese” time delay neural network. <i>International Journal of Pattern Recognition and Artificial Intelligence</i> , 7(04):669–688.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in neural information processing systems</i> , pages 5998–6008.	850
801			851
802			852
803			853
804			854
805			855
806	Narendra Choudhary, Rajat Singh, Ishita Bindlish, and Manish Shrivastava. Sentiment analysis of code-mixed languages leveraging resource rich languages.		856
807			857
808			858
809	Amitava Das and Björn Gambäck. 2014. Identifying languages at the word level in code-mixed Indian social media text . In <i>Proceedings of the 11th International Conference on Natural Language Processing</i> , pages 378–387, Goa, India. NLP Association of India.		859
810			860
811			861
812			862
813			863
814	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. <i>arXiv preprint arXiv:1810.04805</i> .		864
815			865
816			866
817			867
818	Aditya Joshi, Ameya Prabhu, Manish Shrivastava, and Vasudeva Varma. 2016. Towards sub-word level compositions for sentiment analysis of hindi-english code mixed text. In <i>Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers</i> , pages 2482–2491.		868
819			869
820			870
821			871
822			872
823	Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, Gokul N.C., Avik Bhattacharyya, Mitesh M. Khapra, and Pratyush Kumar. 2020. IndicNLP Suite: Monolingual Corpora, Evaluation Benchmarks and Pre-trained Multilingual Language Models for Indian Languages. In <i>Findings of EMNLP</i> .		873
824			874
825			875
826			876
827			877
828			878
829	Saurav Kumar, Saunack Kumar, Diptesh Kanojia, and Pushpak Bhattacharyya. 2020. “a passage to india”: Pre-trained word embeddings for indian languages. In <i>Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)</i> , pages 352–357.		879
830			880
831			881
832			882
833			883
834			884
835			885
836	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. <i>arXiv preprint arXiv:1907.11692</i> .		886
837			887
838			888
839			889
840	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. <i>arXiv preprint arXiv:1301.3781</i> .		890
841			891
842			892
843			893
844	Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .		894
845			895
846			896
847	Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. <i>arXiv preprint arXiv:1910.01108</i> .		897
848			898
849			899