

The Best of both Worlds: Dual Channel Language modeling for Hope Speech Detection in low-resourced Kannada

Adeep Hande¹, Siddhanth U Hegde², Sivanesan Sangeetha³,
Ruba Priyadharshini⁵, Bharathi Raja Chakravarthi⁴

¹Indian Institute of Information Technology Tiruchirappalli

²University Visvesvaraya College of Engineering, Bangalore University

³National Institute of Technology Trichy ⁴ULTRA Arts and Science College

⁵National University of Ireland Galway

adeeph18c@iiitt.ac.in, siddhanthhegde227@gmail.com

sangeetha@nitt.edu, rubapriyadharshini.a@gmail.com

bharathi.raja@insight-centre.org

Abstract

In recent years, various methods have been developed to control the spread of negativity by removing profane, aggressive, and offensive comments from social media platforms. There is, however, a scarcity of research focusing on embracing positivity and reinforcing supportive and reassuring content in online forums. As a result, we concentrate our research on developing systems to detect hope speech in code-mixed Kannada. As a result, we present DC-LM, a dual-channel language model that sees hope speech by using the English translations of the code-mixed dataset for additional training. The approach is jointly modelled on both English and code-mixed Kannada to enable effective cross-lingual transfer between the languages. With a weighted F1-score of 0.756, the method outperforms other models. We aim to initiate research in Kannada while encouraging researchers to take a pragmatic approach to inspire positive and supportive online content.

1 Introduction

The last decade has seen a drastic increase in social media users, owing primarily to easier access to the internet as a result of global modernization (Johnson, 2021). As a result of the surge, several minority groups have turned to social media for support and reassurance. This, however, poses a serious risk to adolescents and young adults who are avid internet users. Social media apps like Facebook, Twitter, and YouTube have become an essential part of their daily lives (Kietzmann et al., 2011). Certain ethnic groups or individuals are victims of social media manipulation to foster destructive or disruptive behaviour, which is a common

scenario in cyberbullying (Abaido, 2020). However, these systems ignore potential biases in the dataset on which they are trained and may harm a specific group of social media users, frequently leading to gender/racial discrimination among its users (Davidson et al., 2019).

As a result, there is a need to detect hope speech in social media. Several marginalised groups seek comfort and assistance from social media content that they can relate to and empathise with others' situations (Chakravarthi, 2020). This type of speech is essential for everyone because it encourages people to improve their quality of life by taking action. Hope speech aims to inspire people suffering from depression, loneliness, and stress by providing assurance, reassurance, suggestions, and support (Herrestad and Biong, 2010). Because most social media in multilingual communities still revolve around English, the phenomenon of code-mixing is common. According to studies, code-mixing is an essential component of social media in multilingual countries (Jose et al., 2020).

Kannada (ISO 639-3:kan) is one of India's low-resource Dravidian languages. Dravidian languages are spoken by over 200 million people, mostly in southern India and northern Sri Lanka (Steever, 1998). The language is primarily spoken by people in Karnataka, India, and it is also recognised as an official language of the state (Hande et al., 2020). Kannada script, also known as Catanese, is an alphasyllabary of Brahmic scripts that evolved into the Kadamba script (Chakravarthi et al., 2019). Kannada has over 43 million speakers¹. However, as previously stated, Kannada is a language with limited resources due to a lack of

¹<https://www.ethnologue.com/language/kan>

language technologies.

Our work aims to detect hope speech in low-resourced code-mixed languages. We develop models on hope speech detection in low-resourced kannada. we propose that a language model would learn effectively with the help of the parent translations. We make use of translations with Google Translate API and experiment with several multilingual language models to find the best performing model. We define Dual channel language model as a model that uses two translations, namely, code-mixed Kannada and English. We present DC-LM, (Dual-Channel Language Model) based on the architecture of BERT that uses the translation of the dataset as additional input for training, performing better in contrast to the typical fine-tuned multilingual BERT. We perform a comprehensive analysis of our models on the dataset along with a thorough error analysis on its predictions on the dataset.

2 Related Work

Researchers have worked on extracting data from social media, particularly from user comments on YouTube, Facebook, and Twitter (Chakravarthi et al., 2020; Severyn et al., 2014). Most information extracted from social media does not adhere to grammatical rules and is written in code-mixed, or non-native scripts, as is common among users from multilingual countries (Jose et al., 2020; Bali et al., 2014). People can communicate on social media without face-to-face interaction, but they are prone to misunderstandings because they do not consider the perspectives of others. There have been few previous efforts on hope speech identification, with the only dataset contribution being (Chakravarthi, 2020), a large multilingual corpus manually annotated for English, Tamil, and Malayalam, with around 28K, 20K, and 10K comments, respectively.

Several researchers have worked to promote positivity on social media by developing and analysing systems that filter out malignancy on social media by focusing on very specific events such as crisis and war (Palakodety et al., 2020), inter-country social media dynamics (Sarkar et al., 2020), and protests (Sohn and Lee, 2019). The authors conducted a shared task on hope speech detection for comments scraped from YouTube in these languages to encourage more research into hope speech for English, Malayalam, and Tamil (Chakravarthi and Muralidaran, 2021). The organ-

isers of the collaborative task used the HopeEDI (Chakravarthi, 2020) Multilingual hope speech dataset. In Malayalam (Hossain et al., 2021), fine-tuning a pretrained XLM-RoBERTa model resulted in the best-weighted F1-score of 0.854. In Tamil (Sharma and Arora, 2021), an ensemble of synthetically generated code-mixed data for training ULM-FiT, baseline-KNN, and a fine-tuned RoBERTa achieved the best score of 0.61. The authors fed the combination of pretrained XLM-R and Tf-Idf Vectors as inputs to an inception block, leading to a weighted F1-Score of 0.93 (Huang and Bai, 2021).

3 Dataset

We use the code-mixed Kannada Hope speech dataset (Hande et al., 2021b). The dataset has two labels, namely Hope and Not-Hope. Table 1 refers to the dataset statistics. Some examples of Hope speech and Not-hope speech classes are shown in Fig 1. For a person, *Hope* can be defined as an inspiration to people battling depression, loneliness, and stress by assuring promise, reassurance, suggestions, and support (Chakravarthi, 2020). Dataset is annotated based on the following guidelines:

Hope speech:

- The comment comprises an inspiration provided to participants by their peers and others, offering reassurance and insight.
- Comment talks about equality, diversity, and inclusion
- Comment talks about the survival story of people from marginalised groups.

Non-hope speech

- The comment produces hatred towards a person or a marginalised group.
- The comment is very discriminatory and attacks people without thinking of the consequences.
- The comment comprises racially, ethnically, sexually, or nationally motivated slurs.
- The comments do not inspire Hope in the readers' mind.

3.1 Pre-Processing

As the data is extracted from the comments section of YouTube, preprocessing would be imperative. To better adapt algorithms to the dataset, we follow the steps for preprocessing comments as listed below.

1. URLs and other links are replaced by the word, ‘URL’.
2. The emojis are replaced by the words that the emoji represents, like happy, sad, among other emotions depicted by emojis. As emojis mainly depict a user’s intention, it would be imperative to replace them with their meanings to pick up their cues. As most models are pretrained only on unlabelled text, we feel that it would be necessary.
3. Multiple spaces in a sentence and other special characters are removed as they do not contribute significantly to the overall intention.

Language Pair	Kannada-English
Vocabulary Size	18,807
Number of Posts	6,176
Number of Sentences	6,871
Tokens per post	9
Sentences per post	1

Table 1: Dataset Statistics

Class	Non-hope Speech	Hope Speech
Training	3,265	1,675
Development	391	227
Test	408	210
Total	4,064	2,112

Table 2: Class-wise distribution of Train-Development-Test Data

We use *nlTK*² for tokenizing words and sentences and calculating the corpus statistics as shown in Table 2. We observe that the vocabulary size is significant due to code-mixed data in a morphologically rich language (Hande et al., 2021a).

We find that non-hope speech makes up the majority of the dataset. The dataset had 7,572 comments after annotation, with *Not-Kannada* having

²<https://www.nltk.org/>

a distribution of 1,396 out of 7,572 comments. We removed the comments labelled as *Not-Kannada*, resulting in a dataset of 6,176 comments. The dataset is divided into three sections: train, development, and test. The training set accounts for 80% of the distribution, while the development set accounts for 10%, which is equal to the distribution of the test set. Table 2 shows the class-wise distribution of data for the train, development, and testing phases. The classes are not evenly distributed across the dataset, with Non-hope speech accounting for 65.81 percent and Hope speech accounting for 34.19 percent. The difference in the distribution after removing the sentences with the *Not-Kannada* label is shown in Table 2.

- T_1 : ತುಂಬು ಹೃದಯದ ಶುಭಾಶಯಗಳು ಕನ್ನಡ ಚಿತ್ರರಂಗದ ಅಭಿಮಾನಿಗಳಿಂದ
Transliteration: Tumbu hrdayada shubhasayagalu kannada citrarangada abhiniganigalinda.
Translation: Best wishes to the Kannada Cinema Industry from the bottom of my heart.
Label: Hope
This comment is classified as hope, as the speaker motivates and inspires the reader by his/her/their greetings to the Kannada Cinema Industry; Hence the comment instigates hope to its readers.
- T_2 : ಸಾರ್ ನಿಮ್ಮ ತಂದೆ ನಿಮಗೆ ಕಲಿಸಿದ ಸಂಸ್ಕೃತ ಸಂಸ್ಕೃತಿ ನಮಗೆ ತುಂಬಾ ಇಷ್ಟ ಆಯ್ತು ಮತ್ತು ನೀವು ಅವರು ತೋರಿಸಿದ ಮಾರ್ಗದರ್ಶನದಲ್ಲಿ ನಡೀತಾ ಇರೋದು
Transliteration: Sir nimma tande nimage kalisida sanskara sansthe namage thumba ishta aytu mattu neevu avaru toresida margadhharshanadalli nadita erodu
Translation: Sir I like the culture your father had taught you, I hope you follow the path he guides you in.
Label: Hope
The sentence is classified as hope, due to the nature of the comment, appreciating the cultures and the behavioural knowledge interpreted by the son from his father.
- T_3 : Yaru tension agbede yakandre dislike madiravru mindrika kadeyavru
Translation: No one needs to worry as the people who disliked this are fans of Mandrika
Label: Not-hope
This sentence is classified as Not-hope. Despite the comment consoling someone because their opinion was disliked, the comment spreads hate to the person named Mandrika.
- T_4 : ಟ್ರೋಲ್ ಅಂದ್ರೇ, ಬ್ರೋ ನಾನು ಟಿಕ ಟಾಕ್ ಗೆ ಅಡಿಕ್ಟ್ ಆಗಿದೆ ಬಟ್ ನಮ್ ದೇಶಕ್ಕಿಂತ ದೊಡ್ಡಲ್ಲ, ಈ ಟಿಕ ಟಾಕ್ ಅಷ್ಟೇ ಇನ್ನೊಂದ್ ವಿಷ್ಣು ನಮ್ ದೇಶದ್ ರೊಪೊಸೋ ಡೌನ್‌ಲೋಡ್ ಮಾಡಿ ಓಪನ್ ಮಾಡಿ ನೊಡುದು
Transliteration: Troll andre, bro naanu tiktok ge addict agide but namma deshakkinta doddadalla, ee tiktok ashte ennond namm deshada rofoso download madi nodu.
Translation: For Troll, bro, I am addicted to TikTok, but it is not bigger than our nation; download our own Indian app Rofoso.
Label: Not-hope
This comment can be classified as Not-hope. Even though the comment states that TikTok is not more significant than the nation, expressing patriotism, the comment may or may not be factually correct. Hence, the comment spews unnecessary hatred towards TikTok.

Figure 1: Examples of Hope speech and Not-hope speech classes.

4 Methodology

We perform extensive analysis on the Kannada hopespeech dataset using a variety of classifiers, ranging from simple machine learning algorithms

to complex deep learning algorithms. To tabulate our results, we employ the scikit-learn library (Buitinck et al., 2013). We conduct our experiments in the manner described below. We ran an average of 5 runs on each model to tabulate the results. We avoid using stopwords or other lemmatisation techniques because Kannada is a morphologically rich language. For machine learning algorithms, we used the scikit-learn library. We used the Pytorch implementation of the pretrained language models available on Huggingface Transformers³. We fine-tuned the models on Google Colaboratory⁴ for its easier access to GPU resources and User Interface.

4.1 Machine Learning Algorithms

For our experiments, we used Logistic Regression (LR). The input features are Term Frequency Inverse Document Frequency (TF-IDF) values ranging from 1 to 5-grams, with the inverse regularisation parameter, C, set to 0.1. It is a control variable that, by being positioned inversely to the lambda regulator, retains the strength modification of regularisation. We applied uniform weights to KNN for classification with 3, 4, 5, and 7 neighbours. We use *Minkowski* as the distance metric, with the distance metric’s power parameter (p) set to 2 and uniform weights for the neighbours. The maximum depth for decision trees and random forests was 500, and the minimum sample splits were 5, with *emphGini* as the criterion. We test a Naive Bayes classifier for multinomially distributed data, with ($\alpha = 1$) for Laplace smoothing to avoid zero probabilities.

We set the maximum depth for the decision tree classifier to 500 and the minimum sample splits to 5, using Gini as the criterion. We looked at random forest classifiers with the same parameters as decision trees. Furthermore, we evaluate a Naive Bayes classifier for multinomially distributed data, with $\alpha = 1$ for Laplace smoothing to avoid zero probabilities.

4.2 Fine-tuning pretrained Language Models

The success of the transformer architecture (Vaswani et al., 2017) has resulted in the researchers adapting to transformer-based models from conventional recurrent neural networks (RNN). We have fine-tuned four pretrained language models for hope speech detection, all of

which are based on the primary architecture of BERT. Because all models were pre-trained on unlabeled monolingual or multilingual data, the models may struggle to classify code-mixed sentences. Because this is a binary classification task, we use Binary Crossentropy as the loss function. By decoupling weight decay from gradient update, we use the Adam optimizer (AdamW) available on Huggingface Transformers (Loshchilov and Hutter, 2019). The corpus is first tokenized to cleave

Hyper-parameters	Characteristics
Optimizer	AdamW
Batch Size	[32, 64, 128]
Dropout	0.1
Loss	Binary cross-entropy
Learning rate	2e-5
Max length	128
Epochs	10

Table 3: Hyper-parameters used for fine-tuning BERT-based language models

the word into tokens. During tokenization, the special tokens needed for sentence classification, the [CLS] token at the start of a sentence and the [SEP] token at the end. Post the addition of the special tokens, the tokens are replaced by ids (*input_ids*), and *attention_masks* for training. During fine-tuning, we extract the pooled output of the [CLS] token and feed the output through an activation layer (Sigmoid) to compute the output prediction probabilities for the given sentence (Hande et al., 2021c).

We used two language models that are part of the pretrained architecture of the BERT (Devlin et al., 2019). We use **bert-base-uncased**, a monolingual language model with a 12-layer, 768-hidden dimension, 12-heads, and 110 million parameters that has been pretrained only on lower cased English text. (Pires et al., 2019), a multilingual version of BERT, is pretrained on publicly available Wikipedia dumps of the top 100 languages. We use **bert-base-multilingual-cased**⁵, which is pretrained on cased text from the top 104 languages and has 12 layers, 768 hidden dimensions, 12 heads, and 179 million parameters. Both models use the same parent architecture, with the only difference being the corpora used during pretraining.

³<https://huggingface.co/transformers/>

⁴<https://colab.research.google.com/>

⁵<https://github.com/google-research/bert/blob/master/multilingual.md>

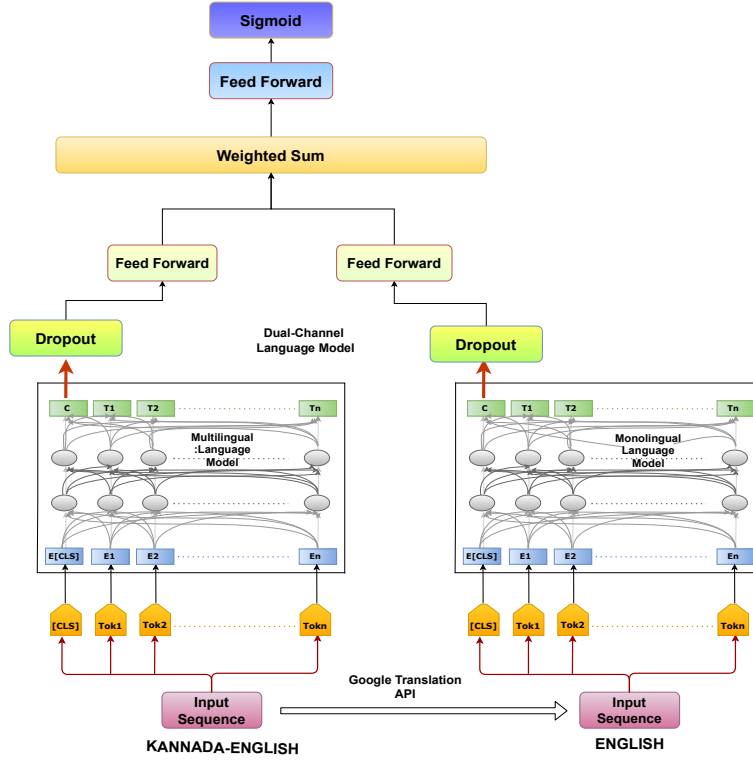


Figure 2: Dual-Channel BERT-based Language Model [DC-LM]

4.2.1 RoBERTa

In contrast to BERT, RoBERTa (Liu et al., 2019), disregards the Next Sentence Prediction (NSP) loss from its pretraining because the authors found no improvement regardless of the loss function. RoBERTa tokenizes using byte-pair encoding (BPE) rather than BERT’s WordPiece tokenization. Textbfrobert-base is a monolingual language model pretrained on 160GB of unlabeled English texts, with 12 layers, 768 hidden dimensions, 12 heads, and 125 million parameters.

4.2.2 XLM-RoBERTa

XLM-RoBERTa is based on large-scale unsupervised cross-lingual learning. **xlm-robert-base**, the smaller version of the model, has 270 million parameters, 12-layers, 768 hidden states, and 8 heads, and was trained on 2.5 TB of newly created clean Common Crawl data in 100 languages.

4.2.3 Dual-Channel Language Model

We propose a Dual-Channel LM (DC-LM), as shown in Fig 2, by fine-tuning a language model based on the transformer architecture on the code-mixed data and its translation in English. We use the Googletrans API ⁶ to translate the code-mixed KanHope to English. This API employs

the GoogleTrans Ajax API⁷ to make calls to detect methods and translate. We invoke the *Translator* function and set the destination language to English, as the *Translator* attempts to identify the language’s source on its own. The use of two channels of pretrained language models is dependent on the advancements of English language models. We obtain more training data for hope speech in English by translating the sentences to English. We believe that when using Dual Channel language model, one model for the code-mixed Kannada-English texts - a multilingual language model - and the other model for the translated English texts - a monolingual language model (pretrained on English), learn better from two languages rather than one. The weighted sum will be the weighted sum of two pooled outputs obtained from the [CLS] token. To fine-tune the code-mixed sentences, we tokenized them with a pretrained multilingual tokenizer and the translated English sentences with a monolingual tokenizer pretrained on English. The first channel (RoBERTa, BERT, or XLNeT) received the translated text, whereas the multilingual language model received the usual raw text (mBERT or XLM-RoBERTa). The pooled output was extracted from the [CLS] token of both models, as

⁶<https://pypi.org/project/googletrans/>

⁷<https://translate.google.com/>

Model	Not-Hope			Hope						
	P	R	F1	P	R	F1	Acc	W(P)	W(R)	W(F1)
Logistic Regression	0.681	0.964	0.798	0.788	0.228	0.354	0.693	0.721	0.693	0.634
KNN	0.705	0.890	0.787	0.659	0.364	0.469	0.696	0.688	0.696	0.670
Decision Tree	0.732	0.797	0.763	0.591	0.500	0.542	0.688	0.680	0.688	0.681
Random Forest	0.736	0.867	0.796	0.673	0.469	0.553	0.720	0.713	0.720	0.706
Naive Bayes	0.719	0.885	0.793	0.674	0.408	0.508	0.709	0.702	0.709	0.688
mBERT	0.757	0.854	0.802	0.680	0.531	0.596	0.735	0.728	0.735	0.726
BERT	0.758	0.780	0.769	0.604	0.575	0.589	0.704	0.701	0.704	0.702
DC-LM(bert-mbert)	0.771	0.836	0.802	0.672	0.575	0.619	0.740	0.734	0.740	0.735
DC-LM(roberta-mbert)	0.788	0.838	0.812	0.690	0.614	0.650	0.756	0.752	0.756	0.752
DC-LM(roberta-xlmr)	0.777	0.779	0.778	0.621	0.618	0.620	0.720	0.720	0.720	0.720
DC-LM(bert-xlmr)	0.727	0.735	0.731	0.589	0.587	0.591	0.650	0.655	0.647	0.651
DC-LM(xlnet-mbert)	0.757	0.759	0.758	0.601	0.598	0.600	0.700	0.700	0.701	0.726
DC-LM(xlnet-xlmr)	0.798	0.851	0.829	0.702	0.635	0.639	0.770	0.758	0.767	0.766

Table 4: Class-wise Precision (P), Recall (R), and F1-Scores for both the classes of the dataset. DC-LM(model1-model2): model1: Monolingual, model2: Multilingual

shown in Fig 2, and a layer took the weighted sum of both pooled outputs. The overall output was then fed into a feed-forward network, which was then activated with a sigmoid function.

DC-LM (model1-model2) is a dual-channel model that uses *model1* for translated text and *model2* for code-mixed texts. *model1* is trained on translated text using two language models based on BERT and RoBERTa. We use two multilingual models for the *model2*, mBERT and XLM-RoBERTa.

DC(bert-mbert): This model employs *bert-base-uncased* for the English text and *bert-base-multilingual-cased* for the code-mixed Kannada-English. The same method is used for all other Dual-Channel language models.

5 Results and Discussion

The results of experiments carried out for classifying hope speech with various models are listed in Table 4 in terms of precision and recall for the individual classes, as well as overall accuracy, weighted averages of Precision, Recall, and F1-score. In our test set, there are 390 instances of *not-hope speech* and 228 samples of *hope speech*. Our experiments’ code is available⁸.

We use four language models for the dual-channel LM, listed in Table 4. We fine-tune multilingual BERT and the uncased base version of BERT separately to assess the significance of improving performance in DC-LM if any. Out of the two BERT models, multilingual BERT performs

better than the BERT model that was pretrained only on English, with a minor increase of 2.1%. However, the performance between the machine learning algorithms and pretrained language models differ by around 7.8%. We trained three dual-channel language models based on the possible combinations between the monolingual and multilingual models. *DC-LM (bert-mbert)* used the monolingual BERT (only English) for the translated text, while the multilingual BERT for the code-mixed Kannada-English texts. DC-LM(bert-mbert) achieves a weighted F1-Score of 0.740, an improvement of 0.5% from mBERT and 3.6% from monolingual BERT. When *XLNet* is used for the translated texts and *XLM-RoBERTa* for the code-mixed texts, it achieves the best performance of all the models, having an F1-Score of 0.766. The principal reason for this increase comes down to the better hyper-parameter tuning and pretraining strategy used in XLM-RoBERTa and XLNet.

DC-LM (roberta-xlmr) has also been fine-tuned to evaluate if there is cross-lingual transfer between the models. Despite being pre-trained on 2.5 TB of data and using an unsupervised cross-lingual learning scale, we find that this model performs worse than DC-LM (bert-mbert). One of the causes for XLM-poor R’s performance, we feel, is its tokenizations. Despite the fact that the developers of XLM-R claim that the model’s performance is unaffected by the type of encoding used in tokenizations, it is discovered that Byte-Pair Encoding (BPE) has a lower morphological alignment with the actual code-mixed text (Jain et al., 2020). In

⁸<https://github.com/adeepH/DC-LM>

Label	Texts	Predictions
Not-Hope	Text: Finally, sonu gowda b day dhinane tiktok ban aythu Translation: Finally, TikTok got banned on Sonu Gowda’s Birthday	Hope
Not-Hope	Text: Found 806 rashmika mangannas Translation: Found 806 Rashmika monkeys	Hope
Hope	Text: Guru ee desha uddhara agatte indian youth volle ide Translation: Brother this country will develop as Indian youth are fantastic	Not-Hope
Hope	Text: thogari tippa supar Translation: Thogari Tippa Super	Not-Hope

Table 5: Predictions on the Test Set

contrast to BERT’s WordPiece tokenization, XLM-R employs the BPE tokenizer, which results in more subwords. We believe XLM-RoBERTa performs worse than multilingual BERT since Kannada is a semantically rich language (Tanwar and Majumder, 2020).

Surprisingly, the monolingual BERT (only English) performed worse than some machine learning algorithms in terms of precision, recall, and F1 scores. We believe this is due to the dataset’s characteristics.

5.1 Error Analysis

We observe that the model predict 331 out of 390 samples correctly for the *Not-hope* label, while the model predicts 145 out of 228 samples correctly for the other class. We observe that several texts have been misclassified for reasons beyond the scope of the model. We have tabulated some predictions in Table 5

Text: “Thogari Tippa“ super

Thogari Tippa is the name of a popular movie that talks about equality. The model identifies it as “Not-Hope Speech“, whereas the dataset classified it as *Hope speech*. The lack of knowledge about the movie is likely the reason why the model predicted incorrectly.

Text: “Guru ee desha uddhara agatte bedu bhai indian youth tumba volle ide“

The text praises the Indian youth, suggesting that India will develop because of them. The model identifies it as *Not-Hope Speech*, even though it should have classified it as *Hope Speech*.

6 Conclusion

A surge in the active users on social media has inadvertently increased the amount of online content available on social media platforms. There is a need to motivate positivity and hope speech in platforms

to instigate compassion and assert reassurance. In this paper, we work on KanHope, a manually annotated code-mixed data of hope speech detection in an under-resourced language, Kannada, consisting of 6,176 comments crawled from YouTube and propose DC-LM, a Dual-Channel BERT-based model that uses the best of both worlds: Code-mixed Kannada-English and Translated English texts. Several pretrained multilingual and monolingual language models were analysed to find the best approach that yields a tremendous weighted F1-Score. We have also trained the dataset on preliminary machine learning algorithms to baseline for future work on the dataset. We believe that this dataset will expand further research into facilitating positivity and optimism on social media. We have developed several models to serve as a benchmark for this dataset. We aim to promote research in Kannada.

7 Acknowledgments

The author Bharathi Raja Chakravarthi was supported in part by a research grant from Science Foundation Ireland (SFI) under Grant Number SFI/12/RC/2289_P2 (Insight_2), co-funded by the European Regional Development Fund as well as by the EU H2020 programme under grant agreement 825182 (Prêt-à-LLOD), and Irish Research Council grant IRCLA/2017/129 (CARDAMOM-Comparative Deep Models of Language for Minority and Historical Languages) for his postdoctoral period at National University of Ireland Galway.

References

- Ghada M. Abaido. 2020. [Cyberbullying on social media platforms among university students in the united arab emirates](#). *International Journal of Adolescence and Youth*, 25(1):407–420.

- Kalika Bali, Jatin Sharma, Monojit Choudhury, and Yogarshi Vyas. 2014. “I am borrowing ya mixing ?” an analysis of English-Hindi code mixing in Facebook. In *Proceedings of the First Workshop on Computational Approaches to Code Switching*, pages 116–126, Doha, Qatar. Association for Computational Linguistics.
- Lars Buitinck, Gilles Louppe, Mathieu Blondel, Fabian Pedregosa, Andreas Mueller, Olivier Grisel, Vlad Niculae, Peter Prettenhofer, Alexandre Gramfort, Jaques Grobler, Robert Layton, Jake VanderPlas, Arnaud Joly, Brian Holt, and Gaël Varoquaux. 2013. API design for machine learning software: experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, pages 108–122.
- Bharathi Raja Chakravarthi. 2020. [HopeEDI: A multilingual hope speech detection dataset for equality, diversity, and inclusion](#). In *Proceedings of the Third Workshop on Computational Modeling of People’s Opinions, Personality, and Emotion’s in Social Media*, pages 41–53, Barcelona, Spain (Online). Association for Computational Linguistics.
- Bharathi Raja Chakravarthi, Mihael Arcan, and John P. McCrae. 2019. [WordNet gloss translation for under-resourced languages using multilingual neural machine translation](#). In *Proceedings of the Second Workshop on Multilingualism at the Intersection of Knowledge Bases and Machine Translation*, pages 1–7, Dublin, Ireland. European Association for Machine Translation.
- Bharathi Raja Chakravarthi, Navya Jose, Shardul Suryawanshi, Elizabeth Sherly, and John Philip McCrae. 2020. [A sentiment analysis dataset for code-mixed Malayalam-English](#). In *Proceedings of the 1st Joint Workshop on Spoken Language Technologies for Under-resourced languages (SLTU) and Collaboration and Computing for Under-Resourced Languages (CCURL)*, pages 177–184, Marseille, France. European Language Resources association.
- Bharathi Raja Chakravarthi and Vigneshwaran Muralidaran. 2021. [Findings of the shared task on hope speech detection for equality, diversity, and inclusion](#). In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 61–72, Kyiv. Association for Computational Linguistics.
- Thomas Davidson, Debasmita Bhattacharya, and Ingmar Weber. 2019. [Racial bias in hate speech and abusive language detection datasets](#). In *Proceedings of the Third Workshop on Abusive Language Online*, pages 25–35, Florence, Italy. Association for Computational Linguistics.
- 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Adeep Hande, Siddhanth U Hegde, Ruba Priyadarshini, Rahul Ponnusamy, Prasanna Kumar Kumareshan, Sajeetha Thavareesan, and Bharathi Raja Chakravarthi. 2021a. [Benchmarking multi-task learning for sentiment analysis and offensive language identification in under-resourced dravidian languages](#). *arXiv preprint arXiv:2108.03867*.
- Adeep Hande, Ruba Priyadarshini, and Bharathi Raja Chakravarthi. 2020. [KanCMD: Kannada CodeMixed dataset for sentiment analysis and offensive language detection](#). In *Proceedings of the Third Workshop on Computational Modeling of People’s Opinions, Personality, and Emotion’s in Social Media*, pages 54–63, Barcelona, Spain (Online). Association for Computational Linguistics.
- Adeep Hande, Ruba Priyadarshini, Anbukkarasi Sampath, Kingston Pal Thamburaj, Prabakaran Chandran, and Bharathi Raja Chakravarthi. 2021b. [Hope speech detection in under-resourced kannada language](#). *arXiv preprint arXiv:2108.04616*.
- Adeep Hande, Karthik Puranik, Konthala Yasaswini, Ruba Priyadarshini, Sajeetha Thavareesan, Anbukkarasi Sampath, Kogilavani Shanmugavadeivel, Durairaj Thenmozhi, and Bharathi Raja Chakravarthi. 2021c. [Offensive language identification in low-resourced code-mixed dravidian languages using pseudo-labeling](#). *arXiv preprint arXiv:2108.12177*.
- Henning Herrestad and S. Biong. 2010. Relational hopes: A study of the lived experience of hope in some patients hospitalized for intentional self-harm. *International Journal of Qualitative Studies on Health and Well-being*, 5.
- Eftekhari Hossain, Omar Sharif, and Mohammed Moshuiul Hoque. 2021. [NLP-CUET@LT-EDI-EACL2021: Multilingual code-mixed hope speech detection using cross-lingual representation learner](#). In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 168–174, Kyiv. Association for Computational Linguistics.
- Bo Huang and Yang Bai. 2021. [TEAM HUB@LT-EDI-EACL2021: Hope speech detection based on pre-trained language model](#). In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 122–127, Kyiv. Association for Computational Linguistics.
- Kushal Jain, Adwait Deshpande, Kumar Shridhar, Felix Laumann, and Ayushman Dash. 2020. [Indic-transformers: An analysis of transformer language models for indian languages](#). *arXiv preprint arXiv:2011.02323*.
- Joseph Johnson. 2021. [Number of internet users worldwide](#).

- Navya Jose, B. R. Chakravarthi, Shardul Suryawanshi, Elizabeth Sherly, and John P. McCrae. 2020. A survey of current datasets for code-switching research. *2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS)*, pages 136–141.
- Jan H. Kietzmann, Kristopher Hermkens, Ian P. McCarthy, and Bruno S. Silvestre. 2011. [Social media? get serious! understanding the functional building blocks of social media](#). *Business Horizons*, 54(3):241–251. SPECIAL ISSUE: SOCIAL MEDIA.
- 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.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Shriphani Palakodety, Ashiqur R. KhudaBukhsh, and Jaime G. Carbonell. 2020. [Hope speech detection: A computational analysis of the voice of peace](#).
- Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. [How multilingual is multilingual BERT?](#) In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4996–5001, Florence, Italy. Association for Computational Linguistics.
- Rupak Sarkar, HIRAK SARKAR, Sayantan Mahinder, and Ashiqur R. KhudaBukhsh. 2020. [Social media attributions in the context of water crisis](#).
- Aliaksei Severyn, Alessandro Moschitti, Olga Uryupina, Barbara Plank, and Katja Filippova. 2014. [Opinion mining on YouTube](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1252–1261, Baltimore, Maryland. Association for Computational Linguistics.
- Megha Sharma and Gaurav Arora. 2021. [Spartans@LT-EDI-EACL2021: Inclusive speech detection using pretrained language models](#). In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 188–192, Kyiv. Association for Computational Linguistics.
- Hajung Sohn and Hyunju Lee. 2019. [Mc-bert4hate: Hate speech detection using multi-channel bert for different languages and translations](#). In *2019 International Conference on Data Mining Workshops (ICDMW)*, pages 551–559.
- Sanford B Steever. 1998. Introduction to the dravidian languages. *The Dravidian languages*, 1:39.
- Ashwani Tanwar and Prasenjit Majumder. 2020. [Translating morphologically rich indian languages under zero-resource conditions](#). *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 19(6).
- 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 *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.