RESEARCH ARTICLE                                                                      OPEN ACCESS

# Kanglish to English Machine Translation using Pretrained Multilingual Transformers

SANDHYA S*

*(MSc Data Science, NMKRV College for Women, Bengaluru-560011
Email: sandhyashanmugam459@gmail.com)

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

## Abstract:

Kannada is one of the Dravidian languages which is commonly spoken in Karnataka state of India. Kanglish is a code-mixed language which comprises of both Kannada and English languages. Over the course of time with advancement in technology and social media development, it is observed that many people who has a native language different from English tend to mix up their mother tongue with English in day-to-day life. Code-mixed languages play a vital role as most of the text data that are generated on online platforms are posted by people who use such code-mixed languages. Also, such mixed languages are used to browse information on search engines. However, machine translation for Kanglish to English is not explored widely. With the recent advancement in pretrained language transformers, our aim is to achieve the Machine Translation for Kanglish to English using pretrained language models and acquire the best accuracy.

*Keywords* — **Code-mixed language, Machine translation, Kanglish.**

\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*

## I. INTRODUCTION

WE SAY A PERSON IS BILINGUAL IF HE/SHE CAN UNDERSTAND AND COMMUNICATE IN 2 DIFFERENT LANGUAGES. A PERSON CAN DEVELOP TO BE A BILINGUAL NATURALLY WHEN HE/SHE LEARNS A DIFFERENT LANGUAGE APART FROM THEIR MOTHER TONGUE. MANY TEND TO MIX-UP THEIR NATIVE LANGUAGE WITH ENGLISH IN DAY-TO-DAY LIFE WITHOUT EVEN NOTICING IT. IN A SENSE IT HAS BEEN IMPLEMENTED AS A PART OF OUR INTERACTION WITH OTHER INDIVIDUALS WITH THE GROWING SENSE OF COMMUNITY AND TECHNOLOGY.

SINCE THE NUMBER OF USERS ARE ALSO INCREASED ON SOCIAL PLATFORM, THE TEXT DATA THAT ARE BEING GENERATED IN THE FORM OF TWEETS, COMMENTS, OR EVEN USER QUESTIONS ON WEB ARE IN THE FORM OF BILINGUAL. TO FULLY UNDERSTAND THE COMBINATION OF 2 DIFFERENT LANGUAGES IS BIT CHALLENGING AS THE SEMANTIC/GRAMMAR RULES OF EACH LANGUAGE DIFFERS.

AND IT IS TO BE NOTED THAT THE MACHINE TRANSLATION OF CODE-MIXED LANGUAGES LACKS EXPLORATION FOR INDIAN LANGUAGES. AND SINCE INDIA HAS DIVERSE LANGUAGES, IT IS AN ADDED PART OF THE CHALLENGE. BUT THERE ARE PRE-TRAINED MULTILINGUAL TRANSFORMERS THAT ARE AVAILABLE FOR THE INDIAN LANGUAGE IN WHICH WE CAN FINE-TUNE TO PERFORM OUR DESIRED MONOLINGUAL MACHINE TRANSLATION OF KANGLISH.

THE TERM CODE-MIXING MAY SOUND NEW, BUT THIS HAS BEEN IN PRACTICE AMONG US EVEN WITHOUT OUR OWN KNOWLEDGE. IN THE PRESENT WORLD WHERE PEOPLE LIVE IN VARIOUS PARTS OF THE SOCIETY APART FROM THEIR NATIVE, IT IS ESSENTIAL TO PICK UP A NEW LANGUAGE FOR EASY COMMUNICATION. AND SINCE THE HUMAN BRAIN CAN UNDERSTAND MORE THAN ONE LANGUAGE, IT IS QUITE COMMON TO MIX-UP THESE LANGUAGES IN DAILY LIFE. AND WITH THE RECENT DEVELOPMENT IN THE TECHNOLOGY, WHERE PEOPLE CONNECT TO OTHERS ALL OVER THE WORLD, THE TEXT DATA THAT GETS GENERATED ONLINE IS TREMENDOUS. AND THESE DATA NEED NOT BE MONOLINGUAL ALL THE TIMES. IT CAN BE OF BILINGUAL WHICH MIGHT BE CHALLENGING FOR PROCESSING.

MANY APPLICATIONS OF DATA SCIENCE RELY UPON THE TEXT DATA THAT IS BEING EASILY AVAILABLE ONLINE WHICH ARE IN TURN GENERATED BY REAL HUMANS. VARIOUS APPLICATIONS LIKE SENTIMENT ANALYSIS, PREDICTION OF SUCCESS OF A MOVIE OR AN ELECTION AND MANY MORE SUCH OPERATIONS ARE DONE BY ANALYSING THESE TEXTS. IN SUCH REQUIREMENTS, IT IS ESSENTIAL TO PREPROCESS THE DATA IN INITIAL STAGES BEFORE THE DEVELOPMENT OF THE MODELS.

NATURAL LANGUAGE PROCESSING (NLP) HAVE GROWN IN RECENT TIMES WHICH HELPS IN UNDERSTANDING AND PROCESSING OF HUMAN LANGUAGE. THOUGH IN INITIAL STAGES IT WAS WIDELY USED FOR FOREIGN LANGUAGES LIKE ENGLISH, GERMAN ETC., LATELY IT HAS BEEN IMPLEMENTED TO INDIAN LANGUAGE WITH THE ADVANCEMENT IN THE GENERATION OF DIVERSE DATA.

MACHINE TRANSLATION ALSO BEEN VERY USEFUL AND IS BEING IMPLEMENTED WIDELY IN RECENT DAYS AS IT IS SUPPORTING DIFFERENT LANGUAGES INCLUDING INDIAN LANGUAGES. MACHINE TRANSLATION HELPS WITH TRANSLATING ONE LANGUAGE SCRIPT TO THE OTHER LANGUAGE. ALSO, IT HELPS IN TRANSLITERATING A LANGUAGE IN OTHER LANGUAGE. THOUGH MACHINE TRANSLATION IS HELPING IN UNDERSTANDING THE LANGUAGE AND PERFORMING THE USER REQUIRED OPERATIONS BY UNDERSTANDING THE RULES AND

SEMANTICS OF EACH INDIVIDUAL LANGUAGE, IT IS QUITE CHALLENGING TO DEAL WITH THE CODE-MIXED LANGUAGE, WHERE THE SCRIPT FORMED WITH THE COMBINATION OF MORE THAN ONE LANGUAGE.

ONE SUCH CODE-MIXED LANGUAGE IS THE KANGLISH (KANNADA + ENGLISH). THOUGH DEALING WITH CODE-MIXED LANGUAGE HAVE BEEN IMPLEMENTED AND ARE IN USE FOR FOREIGN LANGUAGES, IT IS NOT EXPLORED FOR INDIAN LANGUAGES. THE KEY INSIGHT OF THIS PROJECT IS TO BUILD A MODEL WHICH CAN TRANSLATE THIS CODE-MIXED LANGUAGE KANGLISH TO A MONOLINGUAL LANGUAGE (ENGLISH IN OUR CASE), UTILIZING THE DIFFERENT SEQ2SEQ PRE-TRAINED MULTILINGUAL TRANSFORMER-BASED MODELS AND ACHIEVE THE BEST ACCURACY.

## II. RELATED WORK

IN THIS SECTION, WE WILL DISCUSS ABOUT THE PREVIOUS WORKS THAT WERE DONE IN THE FIELD OF CODE-MIXED MACHINE TRANSLATION BY SEVERAL RESEARCHERS AND DATA SCIENTISTS. WE CAN UNDERSTAND HOW THE WORK HAS BEEN PROGRESSING IN THIS AREA SINCE THE BEGINNING AND THE RESEARCH GAPS THAT NEEDS TO BE FILLED.

CODE-MIXING ARISES WHEN A PERSON MAKE USE OF 2 OR MORE LANGUAGES IN THE CONTEXT OF THE SAME COMMUNICATION. IT IS QUITE COMMON IN THE PRESENT SOCIETY WITH THE DEVELOPMENT IN SOCIAL MEDIA PLATFORMS AND MULTILINGUAL SOCIETY.

MANY WORKSHOPS (DIAB ET AL.,2014 [7], 2016 [6][8]; AGUILAR ET AL., 2018B [10]) HAVE BEEN CONDUCTED IN NOTABLE CONFERENCES TO ADVANCE IN THE CODE-MIXED MACHINE TRANSLATION FIELD. BUT THESE WORKSHOPS WERE DIVERSELY DEALING WITH RELATED TASKS LIKE LANGUAGE IDENTIFICATION (SOLORIO ET AL., 2014 [7][9]; MOLINA ET AL., 2016 [8]), NER (AGUILAR ET AL., 2018A [10]; RAO AND DEVI, 2016 [11]), PoS TAGGING (JAMATIA ET AL., 2018 [13]), SENTIMENT ANALYSIS, AND QUESTION ANSWERING (CHANDU ET AL., 2018 [12]).

THOUGH THESE WORKSHOPS WERE STEPPING STONES TO ADVANCE AND EXPLORE CODE-MIXED TRANSLATION, THE INADEQUACY OF DATASETS WERE OF A MAJOR HINDRANCE.

GANESH JAWAHAR, EL MOATEZ BILLAH NAGOUDI, MUHAMMAD ABDUL-MAGEED AND LAKS V.S. LAKSHMANAN PUBLISHED A PAPER ON "EXPLORING TEXT-TO-TEXT TRANSFORMERS FOR ENGLISH TO HINGLISH MACHINE TRANSLATION WITH SYNTHETIC CODE-MIXING". IN THIS PAPER THEY USED THE TRANSFORMER-BASED ENCODER-DECODER MODELS LIKE mT5 AND mBART ON MODELS TO TRANSLATE MONOLINGUAL ENGLISH TEXT INTO HINGLISH CODE-MIXED LANGUAGE. ADDED TO WHICH THEY ALSO PROPOSED A METHOD FOR GENERATING CODE-MIXED TEXTS FROM BILINGUAL REPRESENTATIONS. AND THE RESULT SHOWED THAT THEIR mT5 MODEL GIVES BETTER PERFORMANCE WITH 12.67 BLEU SCORE.

SATYAM DUTTA, HIMANSHI AGARWAL AND PRADEEP KUMAR ROY FROM INDIAN INSTITUTE OF INFORMATION TECHNOLOGY, SURAT, GUJARAT, INDIA PUBLISHED A PAPER ON "SENTIMENT ANALYSIS ON MULTILINGUAL CODE-MIXED KANNADA LANGUAGE" WHICH PRESENTS A MODEL THAT HELPS IN SENTIMENT ANALYSIS OF DRAVIDIAN CODE-MIXED KANNADA COMMENTS. USING THE BERT MODEL, IT ACHIEVED A F1-SCORE OF 0.66 ON VALIDATION DATASET, AND F1-SCORE OF 0.619 ON TEST DATASET.

HINGLISH TO ENGLISH MACHINE TRANSLATION USING MULTILINGUAL TRANSFORMERS BY VIBHAV AGARWAL, POOJA RAO S B, DINESH BABU JAYAGOPI FROM INTERNATIONAL INSTITUTE OF INFORMATION TECHNOLOGY BANGALORE, INDIA AND UNIVERSITY OF LAUSANNE, SWITZERLAND PROPOSED MACHINE TRANSLATION OF CODE-MIXED HINGLISH TO ENGLISH LANGUAGE USING THE PRETRAINED TRANSFORMERS LIKE mBART AND mT5 ON THE PHINC DATASET. ALSO, THEY REPORT A SIGNIFICANT JUMP FROM 15.3 TO 29.5 IN BLEU SCORES WHICH WAS A 92.8% IMPROVEMENT OVER THE SAME DATASET [19].

## III. PROPOSED WORK

IN THIS SECTION, WE EXPLAIN ABOUT THE mBART AND mT5 BASED MACHINE TRANSLATION MODEL AND THE DATASET THAT WILL BE USED TO PROCEED THE EXPERIMENT.
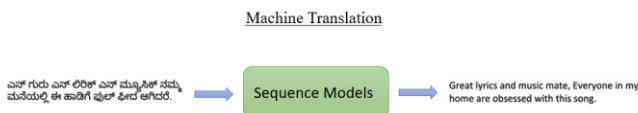
### A. Dataset

For the dataset we create a csv file with the Kanglish sentences along with their English translations. And the Kanglish sentences will be of the kind that are fully Kannada, Kannada mixed with English and Kannada transliterated in English wordings. The dataset is created in this way so that the model will perform well to deal with the different styles of input. The data is collected from internet source such as public corpora or social media comments etc.

TABLE I
A SAMPLE FROM THE COLLECTION OF DATA

| | text | translate |
|---|---|---|
| 0 | ಒಂದು ದೇಶದ ಮುಂದುವರಿಯುವುದು ಅದರ ಆರ್ಥಿಕ ಸ್ಥಿತಿಯನ್ನ... | The progress of a country does not depend on i... |
| 1 | ಕನ್ನಡದಲ್ಲಿ ಡೈಲಿ ಟೆಕ್ ಅಪ್ಡೇಟ್ಸ್ ಪಡೆಯಲು ಸಬ್ಸ್ಕ್... | Subscribe to our channel for Daily Tech update... |
| 2 | Super sar song | Superb song\n |
| 3 | Tiktokers present situation... ನನೋಡುವವರು ಯಾರು ... | Tik Chockers' Parent Status: Who's watching ou... |
| 4 | Super ಸಾಂಗ್ ಪೆರಿ ನೈಸ್.... | The song is very peppy.\n |

### B. Machine Translation Model

Sequence to Sequence (often abbreviated to seq2seq) models is a special class of Recurrent Neural Network architectures that we typically use (but not restricted) to solve complex Language problems like Machine Translation, Question Answering, creating Chatbots, Text Summarization, etc.

Machine Translation



mBART is a seq2seq multilingual bidirectional auto-encoder pretrained model which is an improved model of BART (Lewis et al., 2020 [4][5]) but on large-scale monolingual corpora of 35 languages. It is an architecture that consists of 12 encoder and decoder layers each with 16 attention heads and model dimensions being 1024 resulting in roughly 680 million parameters.

mT5 is a T5 model's multilingual variant pretrained on 101 languages. It has quite similar architecture with 2 encoder and decoder layers each, model dimensions being 1024 and 12 attention heads resulting in approximately 770 million parameters [19].

We propose our model such that the transformers and their associated weights are from HuggingFace's Transformer (Wolf et al., 2020 [1]) package. We planned to use mbart-cc-25 model weight for mBART and mT5-base model weight for mT5. We also plan to finetune the mBART and mT5 models on our created dataset and acquire good performance by altering the hyperparameters and the number of epochs.

Finally, to check the measure the performance accuracy of the model we tend to use the BLEU score for each of the model respectively and compare the results.

## IV. CONCLUSIONS

OUR KEY INSIGHT OF THE PROJECT IS TO USE THE LARGE MULTILINGUAL TRANSFORMERS AND DEMONSTRATE THEIR PERFORMANCE ON TRANSLATING KANGLISH CODE-MIXED LANGUAGE INTO ENGLISH LANGUAGE.

THERE ARE DEEP LEARNING MODELS BY WHICH WE CAN ACHIEVE THIS PROJECT. BUT FINE TUNING THE ALREADY EXISTING I.E., PRETRAINED MODELS FOR OUR REQUIREMENTS WILL BE HELPFUL IN TRANSLATING THE CODE-MIXED LANGUAGES.

AS A PART OF OUR FUTURE WORK, WE WOULD LIKE TO IMPROVE OUR MODEL TO TRANSLATE ENGLISH TO CODE-MIXED LANGUAGE TRANSLATION WHICH IN TURN WILL HELP TO IMPROVE THE CODE-MIXED TO ENGLISH TRANSLATION. ALSO, AS AN EXTENSION TO IMPLEMENT CODE-MIXED TO ENGLISH TRANSLATION TO THE DIVERSE KANNADA SLANGS WHICH VARIATES AS PER DIFFERENT PLACES OF KARNATAKA.

## REFERENCES

[1] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.

[2] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward ¨ Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 8024–8035.

[3] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4- 9, 2017, Long Beach, CA, USA, pages 5998–6008.

[4] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pretraining for natural

language generation, translation, and comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871–7880, Online. Association for Computational Linguistics.

[5] Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. Transactions of the Association for Computational Linguistics, 8:726–742.

[6] Mona Diab, Pascale Fung, Mahmoud Ghoneim, Julia Hirschberg, and Thamar Solorio, editors. 2016. Proceedings of the Second Workshop on Computational Approaches to Code Switching. Association for Computational Linguistics, Austin, Texas.

[7] Mona Diab, Julia Hirschberg, Pascale Fung, and Thamar Solorio, editors. 2014. Proceedings of the First Workshop on Computational Approaches to Code Switching. Association for Computational Linguistics, Doha, Qatar.

[8] Giovanni Molina, Fahad AlGhamdi, Mahmoud Ghoneim, Abdelati Hawwari, Nicolas Rey-Villamizar, Mona Diab, and Thamar Solorio. 2016. Overview for the second shared task on language identification in code-switched data. In Proceedings of the Second Workshop on Computational Approaches to Code Switching, pages 40–49, Austin, Texas. Association for Computational Linguistics.

[9] Thamar Solorio, Elizabeth Blair, Suraj Maharjan, Steven Bethard, Mona Diab, Mahmoud Ghoneim, Abdelati Hawwari, Fahad AlGhamdi, Julia Hirschberg, Alison Chang, and Pascale Fung. 2014. Overview for the first shared task on language identification in code-switched data. In Proceedings of the First Workshop on Computational Approaches to Code Switching, pages 62–72, Doha, Qatar. Association for Computational Linguistics.

[10] Gustavo Aguilar, Fahad AlGhamdi, Victor Soto, Thamar Solorio, Mona Diab, and Julia Hirschberg, editors. 2018. Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-switching. Association for Computational Linguistics, Melbourne, Australia.

[11] Pattabhi R. K. Rao and S. Devi. 2016. Cmee-il: Code mix entity extraction in indian languages from social media text @ fire 2016 - an overview. In FIRE.

[12] Khyathi Chandu, Ekaterina Loginova, Vishal Gupta, Josef van Genabith, Gunter Neumann, Manoj Chin- ¨ nakotla, Eric Nyberg, and Alan W. Black. 2018. Code-mixed question answering challenge: Crowd-sourcing data and

techniques. In Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching, pages 29–38, Melbourne, Australia. Association for Computational Linguistics.

[13] Anupam Jamatia, Bjorn Gamb ¨ ack, and Amitava Das. ¨ 2018. Collecting and Annotating Indian social media Code-Mixed Corpora. In Computational Linguistics and Intelligent Text Processing, pages 406– 417, Cham. Springer International Publishing.

[14] Somnath banerjee, Kunal Chakma, Sudip Kumar Naskar, Amitava Das, Paolo Rosso, Sivaji Bandyopadhyay, and Monojit Choudhury. 2016. Overview of the Mixed Script Information Retrieval (MSIR). In Proceedings of FIRE 2016. FIRE.

[15] 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.

[16] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzman, Edouard Grave, Myle Ott, Luke Zettle- ´ moyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 8440– 8451, Online. Association for Computational Linguistics.

[17] Alexis Conneau and Guillaume Lample. 2019. Cross lingual language model pretraining. In Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pages 7057–7067.

[18] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 483–498, Online. Association for Computational Linguistics.

[19] Vibhav Agarwal, Pooja Rao S B, Dinesh Babu Jayagopi, International Institute of Information Technology

Bangalore, India; and University of Lausanne, Switzerland. "Hinglish to English Machine Translation using Multilingual Transformers", Proceedings of the Student Research Workshop associated with RANLP-2021, pages 16-21, held online, Sep 1-3, 2021