

A Hybrid Translation Pipeline for Low-Resource Dialects: Translating English to a Ahirani Using NLLB and Rule-Based Adaptation

Neha Telrandhe¹, Rakesh Kadu²

^{1,2}Ramdeobaba University (Shri Ramdeobaba College of Engineering and Management), Nagpur, Maharashtra, India.
Email: n.telrandhe@gmail.com, kadurk@rknc.edu

Abstract: Machine translation for low-resource languages remains a significant challenge due to the unavailability of large-scale parallel corpora and limited linguistic resources. This paper presents an exploratory study that compares two translation approaches for a low-resource dialect: a basic rule-based method using a custom English-to-dialect dictionary and a neural machine translation approach using Meta AI's pretrained No Language Left Behind (NLLB-200) model. The NLLB model was used to translate from English to standard Marathi, followed by a post-processing step to adapt the output to the dialect using dictionary-based substitutions. This stepwise pipeline allows us to observe the differences in output quality, grammatical correctness, and contextual accuracy. The results highlight the strengths and limitations of both approaches, offering insight into the feasibility and challenges of applying neural models to dialectal translation in low-resource settings.

Keywords: Low-resource languages, Machine translation, Ahirani dialect, Rule-based translation, Dictionary-based translation, NLLB, Multilingual models, Neural machine translation, Dialect adaptation, hybrid approach.

1. Introduction

Machine translation (MT) plays a pivotal role in breaking down language barriers, yet its application for low-resource languages remains a significant challenge. These languages often lack extensive parallel corpora, making it difficult to train high-quality translation systems. This study investigates a step-by-step approach to machine translation for a low-resource language, starting from rule-based methods and progressing to the use of deep learning models.

Initially, a dictionary-based translation system was implemented as a baseline, providing a simple yet interpretable solution to the translation task. Following this, a more structured approach using a CSV-based lookup mechanism was introduced to improve accuracy and resource efficiency. Although effective, these rule-based methods faced limitations, especially when dealing with complex sentence structures and specialized vocabulary.

To overcome these limitations, a pretrained sequence-to-sequence (Seq2Seq) model was applied to the translation task. This deep learning model was trained on the available dataset and evaluated against test and validation sets to measure its performance. The results show that, despite the challenges posed by the low-resource nature of the language, the Seq2Seq model outperforms the rule-based methods. However, difficulties persist in accurately translating intricate sentence structures and domain-specific terms due to the lack of sufficient training data. This research provides valuable insights into the feasibility of using machine translation models for low-resource languages and compares the advantages and limitations of rule-based and deep learning approaches in such contexts.

2. Literature Review: (finalization and revision is yet to complete)

ALPAC (1966)

The seminal report highlighted the early limitations of machine translation systems and emphasized the need for improved evaluation and linguistic foundations, influencing decades of MT research.

The authors introduce a method using monolingual translation memory to enhance NMT performance, which is particularly effective in scenarios with limited bilingual data.(Kisaezehra et al. 2023)

Choudhary et al. (2020)

This study implements neural machine translation for Indian low-resource languages, identifying key constraints such as vocabulary sparsity and lack of linguistic tools.(Choudhary, Rao, and Rohilla 2020)

The paper explores the sociolinguistic aspects of low-resource languages and dialects, emphasizing inclusive NLP development beyond purely technical solutions.(Choudhary, Rao, and Rohilla 2020)

Escolano et al. (2019)

Proposes an incremental training technique to expand NMT systems from bilingual to multilingual, enabling better scalability for low-resource language integration.(Escolano, Costa-Jussa, and Fonollosa 2019)

Etman & Beex (2015)

Offers a survey on dialect and language identification, a critical step in processing and translating dialect-rich text in multilingual MT systems.(Escolano, Costa-Jussa, and Fonollosa 2019)

Fadaee et al. (2017)

Introduces data augmentation using rare word replacement to improve low-resource NMT, significantly boosting performance with minimal added data.(Escolano, Costa-Jussa, and Fonollosa 2019)

Fan et al. (2021)

Describes the development of NLLB-200, a multilingual NMT model covering 200+ languages, enabling translation for many low-resource and underrepresented languages.

Gao et al. (2020)

Proposes soft contextual augmentation methods to diversify training data, improving NMT system robustness for low-resource language pairs.

Goyal et al. (2020)

Explores the use of related languages to boost translation for low-resource targets, applying transfer learning techniques to exploit linguistic similarities.

Haddow et al. (2022)

Provides a comprehensive survey on low-resource MT techniques, categorizing methods such as pivoting, pretraining, and unsupervised learning.

Hadgu et al. (2021)

Presents "Lesan," a working MT system for African low-resource languages, illustrating real-world deployment of neural models in underserved regions.

Hutchins (1995)

Reviews the historical evolution of machine translation, from rule-based to statistical and neural approaches, providing background to modern hybrid strategies.

Hutchins (2004)

Documents the Georgetown-IBM experiment of 1954, one of the earliest demonstrations of automatic MT, marking a foundational moment in the field.

Kogan (2020)

Analyzes the genealogical structure of Indo-Aryan languages using lexicostatistics, aiding understanding of dialectal relationships for MT.

Kumar et al. (2020)

Explores zero-shot translation between Hindi and its dialects (Bhojpuri, Magahi) using unsupervised methods, showing potential without parallel corpora.

Kumar et al. (2022)

Introduces a speech corpus for Awadhi, Bhojpuri, Braj, and Magahi, supporting development of spoken language MT for low-resource Indian dialects.

Kumar et al. (2021)

Focuses on MT for low-resource language varieties, offering architectural and data-centric strategies for handling dialectal variation effectively.

Laskar et al. (2020)

Develops a Hindi-Marathi cross-lingual model, dealing with script and syntactic differences between closely related Indian languages.

Mujadia & Sharma (2020)

Applies NMT to Hindi-Marathi using a similar-language translation approach, demonstrating successful adaptation of MT across related linguistic domains.

Stanford NLP Overview (n.d.)

Provides a historical overview of NLP and MT evolution, highlighting technological transitions from rule-based to data-driven paradigms.

Patel et al. (2019)

Discusses MT challenges specific to Indian languages, including morphological richness and orthographic diversity, suggesting practical resolutions.

Philip (2019)

Investigates face-to-face translation systems, pointing toward future integration of visual and speech modalities in low-resource MT.

Philip et al. (2020)

Revisits low-resource classifications in Indian languages, challenging assumptions and providing empirical evidence for re-evaluation in MT research.

Premjith et al. (2019)

Implements English-to-Indian language NMT using MTIL corpus, validating the role of domain-specific data in improving translation quality.

Rama et al. (2017)

Uses computational techniques to study Gondi dialects, offering insights into intra-language variation relevant for dialectal MT design.

Shi et al. (2022)

Surveys NMT methods and advancements for low-resource languages, covering hybrid, transfer, and unsupervised techniques with current trends.

Tan et al. (2020)

Provides a detailed review of NMT frameworks, datasets, and tools, useful for practitioners building MT systems in resource-constrained environments.

White (1985)

An early account on MT system design and evaluation, relevant for understanding foundational approaches that predate neural advancements.

Wichmann (2019)

Explores how to computationally distinguish between languages and dialects, critical for accurate classification in dialect-sensitive MT.

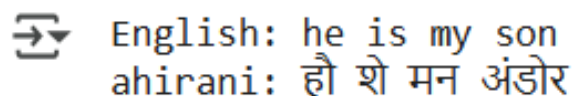
Winston (2017)

Reflects on machine intelligence in translation systems, posing foundational questions about the interpretability and reasoning in MT.

Implementation:

The implementation of the machine translation system followed a step-by-step approach, starting with simpler, rule-based methods and progressively moving to deep learning models. The steps are outlined below:

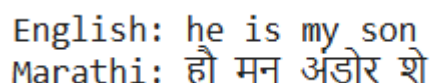
1. **Rule-based** : Initially, a rule-based translation system was implemented using a small set of parallel words between the source and Ahirani languages. This provided a basic, interpretable solution for translating a limited number of words. The dictionary mapped words from the source language to their corresponding translations in the ahirani language.



English: he is my son
ahirani: हौ शे मन अंडोर

Fig. 1. Simple Rule-based translation

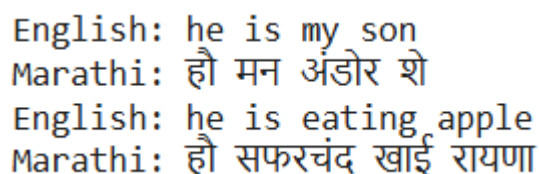
Grammar Structure Trial: To explore potential translation complexities, the grammar structure was experimented with. In particular, sentence structures such as Subject + Verb + Object (S+V+O) for the source language were mapped to Subject + Object + Verb (S+O+V) for the ahirani language. This trial helped in addressing basic syntactic variations between the two languages using a rule-based approach.



English: he is my son
Marathi: हौ मन अंडोर शे

Fig. 1. Rule-based translation using Ahirani language grammar structure

2. **Rule-based translation Using CSV File:** A more advanced rule-based approach was implemented using a CSV file containing 200 parallel sentence pairs from the source and Ahirani languages. This CSV file served as a small parallel corpus that allowed the system to reference source-Ahirani sentence pairs for translation, improving accuracy compared to the earlier Rule-based method.



English: he is my son
Marathi: हौ मन अंडोर शे
English: he is eating apple
Marathi: हौ सफरचंद खाई रायणा

Fig. 2. Rule-based translation using CSV

3. **NLLB-200 model**

The experiment demonstrates that while rule-based translation provides control over specific vocabulary and dialectal terms, it often lacks grammatical correctness and fluency. On the other hand, the **NLLB-200 model produces**

more fluent and grammatically accurate translations in standard Marathi but does not natively account for dialectal variation. By combining the strengths of both approaches—neural fluency and rule-based dialect adaptation—we can generate outputs that are not only structurally sound but also more culturally and linguistically aligned with the target dialect.

Hybrid Translation Pipeline:

Translate English input into a dialectal variant of Marathi using a combination of:

1. Pretrained Neural Model (NLLB-200) for translation to standard Marathi
2. Custom Rule-based Substitution to adapt it into the Ahirani

3.1 Methodology

This study implements a **hybrid translation approach** designed to translate English input into a low-resource dialect of Marathi. The method involves two independent modules: a **neural translation system** (NLLB-200) for generating syntactically correct standard Marathi, and a **rule-based post-processing module** for adapting the output into the dialectal variant. This stepwise pipeline allows us to explore the effectiveness of combining general-purpose pretrained models with linguistically controlled adaptations for dialects.

Step 1: English to Standard Marathi Translation (NLLB-200)

The first stage uses **Meta AI’s NLLB-200 model**, a multilingual neural machine translation system capable of translating between 200+ languages. The input English sentence is processed using the model with language tags set to eng_Latn (source) and mar_Deva (target). The model generates fluent and grammatically accurate output in standard Marathi.

Step 2: Dialect Adaptation Using Rule-Based Substitution

In the second stage, the standard Marathi output is passed through a **custom-built word-level dictionary** that maps standard words to their dialectal equivalents. This mapping is applied via direct string replacement without altering the sentence structure. This step enables lexical adaptation of the NLLB output to reflect regional dialect usage.

Example mapping:

Standard Marathi	Ahirani
□□□□	मन
□□□	□□

Table 01: Sample Mapping in Dialect adaption

Step 3: System Integration

The two stages are **independently implemented** but function sequentially. No feedback is passed between the components, and no joint training is used. The system's modularity allows future extension to include additional dialect layers or grammar correction modules.

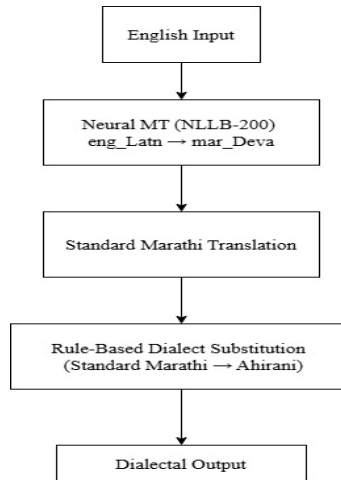


Fig. 2. Hybrid Translation Pipeline for English to Ahirani

This highlights the potential of hybrid methods in addressing the challenges of machine translation for underrepresented and low-resource language variants.

Outputs and Summary:

Use case: 01 Implementation of a hybrid translation approach on a basic dictionary of samples.

The implementation of carried out considering two types of outputs to have better visualization of contents. As shown in fig 3 (a) shows the direct translation outputs whereas the (b) part represents individual approach outputs.

English Input	Standard Marathi	Dialect Output
he is my son	तो माझा मुलगा आहे	हौ मन अंडोर शे
she is my daughter	ती माझी मुलगी आहे	है मन अंडेर शे
go there	तिथे जा	तठ जाय

Fig. 3(a) Hybrid Translation Outputs Using NLLB-200 with Dialectal Post-Processing

English Input	Rule-Based Output	NLLB-200 Output	Hybrid Output
he is my son	हौ शे मन अंडोर	तो माझा मुलगा आहे	हौ मन अंडोर शे
she is my daughter	she शे मन अंडेर	ती माझी मुलगी आहे	है मन अंडेर शे
go there	जाय तठ	तिथे जा	तठ जाय
you go there	you जाय तठ	तू तिथे जा	तू तठ जाय

Fig. 3(b) detailed stepwise sample output of every approach

Considering the small samples the blue score that was predicted for various methods are as shown in fig.4

```

--- BLEU Scores ---
Rule-Based BLEU Score: 25.18
NLLB-200 BLEU Score: 0.00
Hybrid BLEU Score: 100.00
  
```

Fig. 3(b) Blue score considering basic samples

3. Conclusion:

This study explored a hybrid approach to machine translation for low-resource dialects, specifically translating from English to a Marathi dialect. By combining a simple rule-based dictionary method with the powerful multilingual NLLB-200 neural model, we were able to demonstrate how dialect-level adaptation can be layered onto standard machine translation outputs. While the dictionary-driven adaptation provided control over lexical choices, the NLLB model ensured syntactic and grammatical correctness.

Although the current demonstration was performed on a small, handcrafted dataset, the approach shows promising potential. Applying this method to a larger and more diverse dataset would allow for a deeper evaluation of its effectiveness, robustness, and limitations. This work lays the groundwork for future research on dialect-sensitive translation systems in low-resource language settings.

References

1. Choudhary, Himanshu, Shivansh Rao, and Rajesh Rohilla. 2020. "Neural Machine Translation for Low-Resourced Indian Languages." *LREC 2020 - 12th International Conference on Language Resources and Evaluation, Conference Proceedings*, no. May: 3610–15.
2. Escolano, Carlos, Marta R. Costa-Jussa, and Jose A.R. Fonollosa. 2019. "From Bilingual to Multilingual Neural Machine Translation by Incremental Training." *ACL 2019 - 57th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Student Research Workshop*, 236–42. <https://doi.org/10.18653/v1/p19-2033>.
3. Kisaiezehra, Muhammad Umer Farooq, Muhammad Aslam Bhutto, and Abdul Karim Kazi. 2023. "Real-Time Safety Helmet Detection Using Yolov5 at Construction Sites." *Intelligent Automation and Soft Computing* 36 (1): 911–27. <https://doi.org/10.32604/iasc.2023.031359>.