

A Sequence-to-Sequence Approach for Mizo Spelling Correction using Pre-trained Transformer Models

Abstract

The proliferation of digital communication has highlighted the need for robust language processing tools for low-resource languages. Mizo, a Tibeto-Burman language spoken by over a million people, currently lacks effective, publicly available spelling correction systems. This paper addresses this gap by presenting a novel approach to building a Mizo spelling corrector. We leverage the power of transfer learning by fine-tuning a pre-trained multilingual sequence-to-sequence model, mT5-small. Due to the absence of a large, annotated Mizo spelling error corpus, we developed a systematic methodology for generating a synthetic dataset by introducing linguistically-motivated errors, including character-level noise, tonal diacritic manipulation, and common orthographic substitutions specific to Mizo. Our model achieves a BLEU score of 94.17, a Word Error Rate (WER) of 2.96%, and a Character Error Rate (CER) of 0.79% on a held-out test set. This work establishes a strong baseline for Mizo spelling correction and demonstrates that our synthetic data generation strategy is a viable and effective method for bootstrapping NLP tools in low-resource settings.

1. Introduction

The increasing integration of technology into daily life has made digital communication a primary mode of interaction globally. For speakers of low-resource languages like Mizo, this digital shift presents both opportunities and challenges. While it allows for greater connectivity, the lack of fundamental natural language processing (NLP) tools, such as spelling and grammar correctors, acts as a significant barrier to effective communication and digital inclusion (Joshi et al., 2020).

Mizo features unique linguistic characteristics, such as its tonal nature marked by diacritics (e.g., â, ê, î, ô, û) and specific characters like 'ṭ', which are common sources of orthographic errors in digital text. The absence of a standardized, large-scale corpus of

Mizo text, let alone a corpus of annotated spelling errors, makes traditional data-driven approaches to building spelling correctors infeasible.

To overcome these challenges, this paper explores the effectiveness of fine-tuning pre-trained multilingual transformer models. Specifically, we investigate the use of mT5 (Xue et al., 2021), a sequence-to-sequence model trained on a vast corpus of over 100 languages. Our core contribution is a detailed methodology for creating a high-quality synthetic training corpus by systematically introducing Mizo-specific errors into clean text.

The main contributions of this work are threefold:

1. We present a novel, linguistically-informed methodology for generating a synthetic spelling correction dataset for the Mizo language.
2. We train and evaluate an mT5-small model, establishing the first strong public baseline for Mizo automated spelling correction.
3. We provide a detailed analysis of the model's performance, highlighting its strengths in correcting specific error types and discussing its limitations, thereby paving the way for future research.

2. Related Work

2.1. Spelling Correction Systems

Traditional spelling correction systems relied on dictionary lookups and rule-based methods based on minimum edit distance (Damerau, 1964). Later, statistical approaches using noisy channel models became prominent, treating a misspelled word as a "noisy" version of the correct word (Brill & Moore, 2000).

With the advent of deep learning, neural network architectures, particularly sequence-to-sequence (Seq2Seq) models with LSTM or GRU encoders and decoders, demonstrated superior performance by capturing contextual information (Sak et al., 2014). The introduction of the Transformer architecture (Vaswani et al., 2017) and pre-trained models like BERT (Devlin et al., 2019) and T5 (Raffel et al., 2020) revolutionized the field, framing spelling correction as a text-to-text translation task.

2.2. NLP for Low-Resource Languages

A primary challenge for low-resource languages is data scarcity. A significant body of research focuses on overcoming this through techniques like transfer learning, where a model pre-trained on high-resource languages is fine-tuned on a smaller, task-specific dataset (Agić, 2017). Multilingual models like mBERT and mT5 have been particularly successful, demonstrating the ability to transfer grammatical and semantic knowledge to languages unseen or underrepresented during pre-training.

Furthermore, synthetic data generation, or data augmentation, has emerged as a crucial technique. This involves creating "noisy" data from clean text to train models for tasks like machine translation (Sennrich et al., 2016) and grammatical error correction (Rei et al., 2017), a strategy we adapt for Mizo spelling correction.

3. Methodology

Our approach consists of three main stages: (1) creating a clean Mizo text corpus, (2) generating a synthetic parallel corpus of incorrect-correct sentence pairs, and (3) fine-tuning a pre-trained sequence-to-sequence model.

3.1. Corpus Creation

A clean Mizo corpus was aggregated from various digital sources, including the Mizo Bible, online news articles from publications such as 'Vanglaini', and publicly available educational materials. The collected text was normalized by standardizing punctuation and removing duplicate sentences. The final clean corpus consists of [e.g., 55,000] unique sentences, which serves as the ground truth for our synthetic data generation.

3.2. Synthetic Error Generation

To create a realistic training dataset, we designed a noise injection pipeline that simulates common errors made by Mizo speakers. Each clean sentence from the corpus was used to generate multiple corrupted versions by applying the following error types with a set probability:

- **Character-level Noise:**
 - **Substitution:** A character is replaced with a random character, with a higher probability of being a keyboard-adjacent key.
 - **Insertion:** A random character is inserted at a random position.
 - **Deletion:** A random character is deleted.
 - **Transposition:** Two adjacent characters are swapped.

- **Mizo-Specific Orthographic Errors:**

- **Tonal Diacritic Manipulation:** As Mizo is a tonal language, diacritics are a frequent source of error. Our script introduces errors by (a) randomly removing a circumflex from a vowel (e.g., kân → kan) and (b) substituting one accented vowel for another.
- **ṭ to t Substitution:** The retroflex stop ṭ is often incorrectly typed as t. We systematically replace ṭ with t to simulate this common error.
- **Phonetically Similar Substitutions:** Common phonetic confusions (e.g., ph to f) are simulated.

```
def introduce_mizo_specific_error(word, diacritic_error_prob=0.2,
special_char_error_prob=0.1):

    # Diacritic errors

    if random.random() < diacritic_error_prob:

        for original_char, alternatives in DIACRITIC_MAP.items():

            if original_char in word:

                word = word.replace(original_char,
random.choice(alternatives), 1)

                break # Apply only one diacritic error per word for
simplicity

    # Special char errors (e.g., 'ṭ' to 't')

    if random.random() < special_char_error_prob:

        for original_char, replacement_char in MIZO_SPECIAL_CHARS.items():

            if original_char in word:

                word = word.replace(original_char, replacement_char, 1)

                break # Apply only one special char error per word

    return word
```

```

def corrupt_sentence(sentence, char_error_prob_word=0.2,
diacritic_error_prob=0.3, special_char_error_prob=0.2):

    corrupted_words = []

    # Split by spaces, but keep punctuation attached to words initially

    words = sentence.split()

    for word in words:

        original_word = word

        # Apply Mizo specific errors first

        word = introduce_mizo_specific_error(word, diacritic_error_prob,
special_char_error_prob)

        # Then apply generic character errors to the (potentially already
modified) word

        if random.random() < char_error_prob_word:

            word = introduce_char_error(word, error_prob=1.0) # Ensure an
error is attempted if condition met

        corrupted_words.append(word)

    return " ".join(corrupted_words)

```

This process yielded a parallel corpus of [e.g.,1000] (incorrect, correct) sentence pairs.

3.3. Model Architecture

We selected the mT5-small model, a multilingual variant of the T5 (Text-to-Text Transfer Transformer) model. T5 frames all NLP tasks as a text-to-text problem, making it inherently suitable for spelling correction where the input is a noisy sentence and the output is a clean one. To provide the model with task-specific context, every input sentence was prefixed with the string "correct mizo: ".

4. Experimental Setup

4.1. Dataset

The generated parallel corpus was split into training, validation, and testing sets in an 80:10:10 ratio. The test set was held out and used exclusively for the final evaluation to ensure an unbiased assessment of the model's generalization capabilities.

4.2. Training

The mT5-small model was fine-tuned using the Hugging Face Transformers library (Wolf et al., 2020) on a single NVIDIA T4 GPU. The training was configured with the hyperparameters detailed in Table 1. We employed the `load_best_model_at_end` strategy, which saves the model checkpoint that achieves the lowest validation loss.

Hyperparameter	Value
Base Model	google/mt5-small
Optimizer	AdamW
Learning Rate	5e-5 (for initial), 1e-5 (for continued)
Batch Size	8
No. of Epochs	1->10->50->500
Weight Decay	0.01

Table 1: Training Hyperparameters

4.3. Evaluation Metrics

We evaluate the model's performance using three standard metrics:

- **BLEU (Bilingual Evaluation Understudy):** Measures the n-gram overlap between the model's output and the reference sentence. While traditionally used for machine translation, it is effective for measuring text generation quality.
- **Word Error Rate (WER):** The percentage of words that need to be substituted, deleted, or inserted to transform the model's output into the reference sentence. Lower is better.
- **Character Error Rate (CER):** Similar to WER but at the character level. It is particularly useful for spelling correction as it penalizes minor typos. Lower is better.

5. Results

The model's performance on the held-out test set is presented in Table 2.

Model	BLEU ↑	WER ↓	CER ↓
Baseline:Lavenshtein	85.90%	7.20%	2.47%
Our Model (mT5-small)	94.17%	2.96%	0.79%

Table 2: Quantitative Results on the Test Set

Our fine-tuned mT5 model significantly outperforms the baseline, demonstrating the effectiveness of the transfer learning approach. A high BLEU score coupled with low WER and CER indicates that the model generates text that is both structurally sound and orthographically accurate.

Qualitative Analysis

To better understand the model's behavior, we analyzed specific examples from the test set, as shown in Table 3.

Input (Incorrect)	Model Output	Ground Truth	Analysis
Mizo tawng thiam ka duh	Mizo ṭawng thiam ka duh	Mizo ṭawng thiam ka duh	Success: Correctly identified and fixed the common t → ṭ error.
Kan ram dinhmun hi kan vawn nun zl a ngai	Kan ram dinhmun hi kan vawn nun zel a ngai	Kan ram dinhmun hi kan vawn nun zel a ngai	Success: Fixed a character substitution (zl → zel).

Pathian thu chu kan nun kawng engtua ni.	Pathian thu chu kan nun kawng engtua ni.	Pathian thu chu kan nun kawng engtu a ni.	Failure: Failed to identify a subtle semantic error (if one existed).
Vawnah chanchin tha ka hria	Vawinah chanchin tha ka hria ka hria.	Vawinah chanchin tha ka hria ka hria.	Success: Fixed missing 'i' characters.

Table 3: Qualitative Error Analysis

6. Discussion

The quantitative results strongly suggest that fine-tuning pre-trained multilingual models is a highly effective strategy for developing NLP tools for Mizo. The model's low CER (e.g., 0.79%) is particularly encouraging, as it indicates a strong capability to fix the exact types of character-level errors we simulated during training.

Our qualitative analysis reveals that the model excels at correcting common, predictable errors, such as the t/ṭ substitution and single-character typos. Its ability to correct tonal diacritics suggests that the synthetic manipulation of these characters was a successful training strategy.

Limitations:

Despite its strong performance, our model has limitations. First, its effectiveness is intrinsically tied to the quality and diversity of our synthetic data. It may struggle with complex, real-world errors that were not represented in our generation script, such as word-order issues or semantic mistakes. Second, the base corpus, while substantial, may not cover all domains of the Mizo language, potentially limiting the model's vocabulary. Finally, the use of mT5-small might limit the model's capacity compared to larger variants.

7. Conclusion and Future Work

In this paper, we presented a successful methodology for building a Mizo spelling correction system by fine-tuning an mT5-small model on a synthetically generated dataset. Our work establishes a strong baseline and demonstrates the viability of this low-resource NLP development paradigm.

For future work, we propose several directions. First, collecting and annotating a corpus of real-world Mizo spelling errors would be invaluable for creating a more robust

evaluation benchmark and for further fine-tuning. Second, experimenting with larger pre-trained models, such as mT5-base or mT5-large, could yield further performance improvements. Finally, this work could be extended to the more complex task of Grammatical Error Correction (GEC) for Mizo, building upon the foundation laid here.

References

[Use a citation manager like Zotero or BibTeX. Format according to the conference/journal style (e.g., ACL, APA).]

Agić, Ž. (2017). Cross-lingual dependency parsing for closely related languages: A survey. *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*.

Brill, E., & Moore, R. C. (2000). An improved error model for noisy channel spelling correction. *Proceedings of the 38th Annual Meeting on Association for Computational Linguistics*.

Damerau, F. J. (1964). A technique for computer detection and correction of spelling errors. *Communications of the ACM*.

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*.

Joshi, P., Santy, S., Budhiraja, A., Bali, K., & Choudhury, M. (2020). The state and fate of linguistic diversity and inclusion in the NLP world. *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*.

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., ... & Hugging Face. (2020). Transformers: State-of-the-art natural language processing. *Proceedings of the*

2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations.

Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., ... & Raffel, C. (2021). mT5: A massively multilingual pre-trained text-to-text transformer. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics.*