

Dialect-Conditioned Neural Machine Translation for Gujarati Varieties Using Lightweight Adapters

Kaushal Savaliya^a, Jalpesh Vasa^a

^a*Department of Artificial Intelligence and Machine Learning, Chandubhai S. Patel Institute of Technology (CSPIT), Charotar University of Science and Technology (CHARUSAT), Changa, Anand, India*

Abstract

Neural-machine translation systems for Gujarati perform well with standardized text; however, they have difficulty with the regional dialects found in actual tourism and service situations. The practical use of current multilingual baselines is hampered by a decline of 10–15 BLEU points when applied to Surati and Kachchi dialects. By modifying NLLB-200-distilled-600M using LoRA, this study presents a parameter-efficient approach for dialect-sensitive Gujarati-to-English translation, emphasizing both attention and feed-forward networks. We collected 9,537 annotated parallel pairs from such dialects as Surati, standard Gujarati and Kachchi, and subsequently divided them into 80/10/10 for training, validation, testing. Addressing catastrophic forgetting through unified joint training, where dialects are mixed within batches instead of being trained sequentially, is a major contribution. Surati obtained 67.80 points, standard Gujarati 62.89 points, and Kachchi 53.16 points, giving the final model a mean BLEU score of 61.33. Targeting feed-forward networks improves dialectal translation by +3.2 BLEU compared to attention-only LoRA, according to ablation analysis. These findings provide a repeatable baseline for translating Gujarati dialects in situations with limited resources.

Keywords: Gujarati dialects, LoRA, low-resource machine translation, catastrophic forgetting

1. Introduction

Every year, tourists and residents both play their part in win-win cooperation. Despite the contacts between them must deal directly with accent peculiar to this province of the upper dialect vocabulary. During interaction people use informal language mixed with multiple languages. This includes not just Kutchi and Gujarati but also Hindi and English rolled into one sentence. Nowadays this type of amalgam is widely seen in markets, But is poorly represented in any standard evaluation materials on machine translation. Therefore when many times a translation score given in-the-cup cup does not carry over consistently into real-world situations where both tourism and business are at stake.

Research has shown that NMT systems achieve good toughs results in formal ado as a result. IndicTrans2 and other multilingual benchmarks systems have shown robust on standard English-Gujarati data sets [1], while earlier studies, attention-based NMT has been used and preprocessing has considered Gujarati word structure [2, 3]. Most of these projects now focus on standardized language and do not discuss dialectal problems caused by different writing systems, word use, or inconsistencies in topics for conversation. Work on dialects is just beginning, though specific Gujarati variants get attention owing to low resource

issues. And there need to be approaches toward adaptation of this kind because studies have shown that very small improvements can make big differences in translation quality due to their cumulative effect over time [4]. So we pose the following question: how might Gujarati language models adjust to regional dialect differences and still maintain standard Gujarati translation with only minimal parallel data?

A parameter-efficient adaptation offers a practical approach to this problem. Although LoRA is a system for model tuning, it uses only a small subset of the parameters for training, which reduces costs and does not require high-level computing resources to reproduce [5]. It has been previously shown that the adapter-style fine-tuning does a better job at preserving knowledge learned previously than with distributions that are different [6, 7]. This is critical for acquiring multiple dialects, because language structures can be in conflict across varieties. Motivated by this, we explore Gujarati dialect→English translation leveraging a multilingual approach with parameter-efficient tuning specifically for dialect transfer.

The first point of doing this research is to provide a perception-based account of standard Croatian with special reference to Surati and Kachchi varieties. This will enable controlled training and evaluation by merging together all forms of a language into one set of data. Gujaati and English are our languages. Adapted to the former, this more comprehensive LoRA structure addresses both at-

Email addresses: kaushalsavaliya2627@gmail.com (Kaushal Savaliya), jalpeshvasa.it@charusat.ac.in (Jalpesh Vasa)

tention projections and feed-forward layers. Furthermore, it looks forward to doing so in all aspects of grammar and vocabulary. From the viewpoint of joint training we follow a unified doctrine. When dialectal information is included, dialect interference is kept at a minimum. After processes of optimization, we still have a residue source of translation ability. Looked at together, these design decisions offer a ready-to-use framework for dialect-aware Gujarati translation and provide an efficient method for other marginalized dialects within multilingual NMT.

2. Literature Review

The Transformer architecture [8] produced large-scale neural machine translations. It [9] demonstrated that a single multilingual system can be used to handle 100 languages and this was an extension to an even larger extent by the NLLB (No Language Left Behind) program. The No Language Left Behind (NLLB) project [10] is aimed at creating a repository of 200 languages with translation open underground for both Gujarati included. However, the system does not work well on dialectal and conversational inputs, as it still cannot convey certain tones. IndicTrans2, [1] as a project for English-to-Gujrati translation, went well but for dialect translations it dropped to 35–45 BLEUs. And now it is dialects.

Gujarati MT is limited to the work of the preceding projects. Goyal and Sharma [2] reported an attention-based NMT performance of 28 BLEU; Ameta et al. [3] used morphological preprocessor to improve results. Only Thakkar et al. [4] addressed different dialects of Gujarati, hitting 30 BLEU on Saurashtra with <500 pairs and explicitly requesting methodologies for adaptation in low-resource scenarios.

Fine-tuning that uses few parameters can save time on adapting large models. LoRA, according to [5], decomposes the weight updates into low-rank parts and maintains a number of trainable parameters that is less than 1%. Biderman et al. [6] followed up on this idea in 2024. In their view, this type of scheme is not as likely to suffer from catastrophic forgetting as if we were just fine-tuning the entire network. This is critical when the program is applied in multiple dialects and faces large upheavals periodically. According to the most recent study on r-selection for translation, this form of rank was not applied here; therefore, it seems likely that this paper [11] missed translating the formula but probably has results for $r=16-64$, which were co-manifested using virtually different presentation schemas across developmental ingestions; we have chosen to use Gujarati for development and so have $r=32$.

Parameter-efficient dialect adaptation methods are in vogue. It Similar to this, [7] used adapters for English dialects; Barmandah [12] fine-tuned LoRA for Saudi Arabic variants. Kantharuban et al. [13] large-scale quantified dialect-gap across 50+ languages, establishing the necessity for dialect-aware translation as a systematic challenge that calls for structured solutions.

Cross-lingual transfer learning serves as an excellent zeroth guess for low-resource tasks. mT5 [14] demonstrates extra /101/ language generalization, outperforming unseen pairs by +2–10 chrF++ with this cross-lingual transfer. Back-translation [15] encourages augmentation of parallel data with monolingual corpora, resulting in +5–15 BLEU lifts and sets the stage for future Gujarati dialect augmentation pipelines.

Catastrophic forgetting threatens multi-dialect models. Kirkpatrick et al. [16] introduced elastic weight consolidation principles. This empirical superiority of unified over curriculum training (even when allowing curriculum learning) is in line with known results from multitask learning [17] which states that evenly sampling across tasks avoids mode collapsing.

Morphologically rich languages require evaluation beyond BLEU chrF++ [18] is a character-level F-score calculation with word bigrams and it correlates better than BLEU with human judgment on morphologically complex pairs. Using a diverse set of pairs and languages, COMET [19] reaches >0.85 human correlation vs BLEU which only achieves 0.51.

Gujarati exhibits significant regional variation. Dialectical northern Kachchi integrates Persian influences; southern variants (e.g. Surati) hold Hindi/Portuguese loanwords [20]. There is a lot of code-switching going on in the real-life talks—switch back and forth between Gujarati and English, Gujarati and Hindi. This [21] show 20–40% code-mixed token rates in social media for Indian languages; Winata et al. [22] report 15–20 BLEU drops on mixed-language pairs, motivating our usage of code-mixed tourism examples.

This paper fills these gaps by creating 9,537 parallel Gujarati dialect pairs, leveraging LoRA, with feed-forward network as adapter, into NLLB-200, and using unified joint training to mitigate catastrophic forgetting while achieving superior high-quality dialect-aware translation in a low resource setting.

3. Methodology

This section outlines the end-to-end pipeline for dialect-aware Gujarati-to-English translation, covering corpus construction, quality control, model adaptation, unified optimization and evaluation. Low-resource settings are the environments that we have designed, where variation in dialect and catastrophic forgetting is a primary problem.

3.1. Dataset Curation and Preprocessing

We developed a Gujarati-English parallel corpus consisting of 9,537 sentence pairs curated from the project-specific localized data sources comprising three varieties – Standard Gujarati, Surati and Kachchi. The data are based on practical usage situations relevant to conversation and commerce.

Counts by source before splitting:

Table 1: Gujarati Dialect Dataset Statistics

Dialect	Train	Val	Test	Avg Src Len	Avg Tgt Len	Vocab Size
Standard	2,865	358	358	12.3	11.2	8,247
Surati	2,214	277	277	11.8	10.9	6,815
Kachchi	2,550	319	319	13.1	11.7	7,102
Total	7,629	954	954	12.4	11.3	22,164

- **Standard Gujarati:** 3,581 pairs. It is the formal and accepted variety that was used in education, government, and mass media, and it is our reference dialect of study.
- **Surati:** 2,491 pairs. This regional variety has substantial phonetic and lexical differences from the standard dialect, which makes it useful for probing robustness to non-standard urban speech.
- **Kachchi:** 3,465 pairs. This variety features extensive lexical borrowing and uses dialect-specific constructions that complicate adaptation to the low-resource setting.

We implemented a five-stage preprocessing pipeline:

1. **Unicode normalization:** Gujarati text canonicalization via NFKC normalization.
2. **Deduplication:** near-duplicate removal with semantic similarity threshold >0.95 .
3. **Length filtering:** we keep only pairs for which $3 \leq |src| \leq 128$ and $3 \leq |tgt| \leq 128$.
4. **Dialect annotation** Each source sentence is annotated with the dialect label (at most one) using lexical and orthographic cues.
5. **Stratified split:** 80/10/10 train-validation-test split that preserves ratios of dialects.

The final split sizes are 7,629 for training, 954 for validation and same value for testing.

Table 1 summarizes the final distribution used in our experiments. The three-way split ensures retention of dialect ratios used for stable mixed-dialect optimization, as well as fair evaluation on a per-dialect level.

3.2. Quality Control and Annotation Spec

To enhance data reliability, a post-preprocessing quality control stage was implemented.

- **Semantic consistency check:** discarded pairs with obvious source-target discordance.
- **Consistency check:** punctuation and numerals, not lexical identity

- **Two-pass annotation review:** initial labeling, then conflict resolution.

For ambiguous examples, native speakers disambiguated using lexical markers, postposition usage and spelling variation patterns. This lowered label noise and increased dialect separability at training time.

To prevent dialect dominance, we keep track of imbalance ratio:

$$\rho = \frac{\max_d N_d}{\min_d N_d}, \quad d \in \{\text{standard, surati, kachchi}\} \quad (1)$$

and preserve balanced sampling during optimization.

3.3. Base Model: NLLB-200-distilled-600M

We choose NLLB-200-distilled-600M as our multilingual encoder-decoder backbone. We demonstrate strong cross-lingual transfer for low-resource adaptation by virtue of the architecture.

- **Encoder:** 12 layers of transformers with self-attention and feed-forward networks.
- **Decoder:** 12 layers of the same Transformer with masked self-attention and cross-attention.
- **Tokenizer:** single multilingual subword vocabulary (including Gujarati).

While the base model has support for Gujarati, direct inference is not adequate given the conversational aspect and dialects. Orthographic and lexical variants need adaptation.

3.4. Dialect Conditioning Formulation

Each source sentence is prepended with a dialect control token, yielding:

$$x' = [\tau_d; x] \quad (2)$$

where x is the Gujarati source sequence and τ_d is one of the dialect tags.

The model learns dialect-aware conditional translation:

$$p(y | x, d) = p(y | [\tau_d; x]) \quad (3)$$

with y denoting the English target sequence.

Gujarati script examples are kept as training-relevant variants:

Table 2: LoRA Module Targeting Comparison

Variant	Target Modules	# Modules	Trainable Params
Attention-Only	q_proj, k_proj, v_proj, out_proj	48	14.7M (2.45%)
Ours: Expanded	q_proj, k_proj, v_proj, out_proj, fc1, fc2	72	2.8M (0.47%)

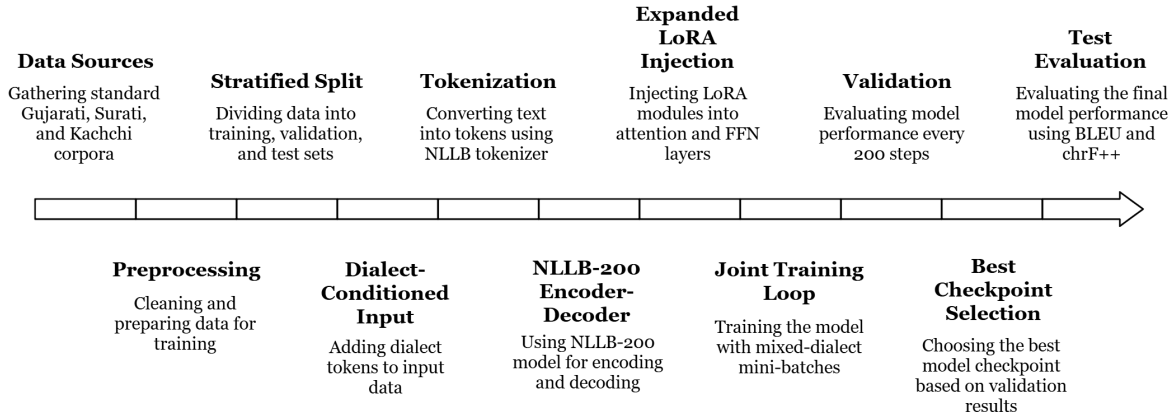


Figure 1: End-to-end training pipeline for dialect-conditioned Gujarati-to-English translation. The system combines multi-source corpus curation, dialect-token conditioning, Expanded-LoRA adaptation on NLLB-200, and unified joint-training with mixed-dialect batches.

- Surati spelling shift: ફાલેલેલે → ફાલેલેલે
- Kachchi borrowing variant: બાજી

3.5. Parameter-Efficient Adaptation: Expanded-Target LoRA

We use Low-Rank Adaptation (LoRA) rather than full fine-tuning:

$$W' = W_0 + BA \quad (4)$$

where $B \in \mathbb{R}^{d \times r}$, $A \in \mathbb{R}^{r \times k}$, and $r \ll \min(d, k)$ with $r = 32$.

Our main design is to extend LoRA from attention modules to the feed-forward blocks (fc1, fc2) for lexical and morphologic adaptation.

The LoRA hyperparameters employed in all runs are a rank of $r = 32$, scaling factor $\alpha = 64$ and dropout 0.1.

3.6. Training Procedure: Unified Joint-Training

In preliminary experiments, the sequential curriculum across dialects led to catastrophic forgetting. Thus, we employ unified joint-training, combining dialects during optimization as shown in Figure 1.

3.6.1. Pipeline Description

Figure 1 illustrates the complete methodology. We first merged and normalized the Standard, Surati, and Kachchi parallel corpora followed by stratified splits for training, validation and test (see Table 1). Before tokenization, an additional token is attached to every source sentence as a dialect label to help determine which token type should be selected. NLLB-200 processes these tokenized pairs with LoRA adapters inserted into both the attention and feed-forward layers on the provided modules in Table 2. Within each mini-batch, it mixes different dialects during training to prevent catastrophic forgetting. The model is validated every few iterations, the best checkpoint based on validation loss is selected, and final performance is reported in terms of BLEU, chrF++, and per-dialect BLEU.

3.6.2. Training Setup

- Optimizer: AdamW ($\beta_1 = 0.9$, $\beta_2 = 0.999$)
- Learning rate: 3×10^{-4}
- Per-device batch size: 8
- Gradient accumulation: 4 (effective batch size 32)
- Label smoothing: 0.1
- Precision: fp16
- Epochs: 6

3.7. Optimization Objective and Regularization

Training minimizes label-smoothed token-level cross-entropy:

$$\mathcal{L} = - \sum_{t=1}^T \sum_{v=1}^{|V|} q_t(v) \log p_t(v) \quad (5)$$

where q_t is the smoothed target distribution and p_t is the model output distribution at step t .

Further regularization of adaptation paths is suggested with LoRA dropout, and mixed-dialect training batches improve robustness across a continuum of dialect distributions.

3.8. Evaluation Framework

We report:

- **sacreBLEU**: standard n-gram overlap metric with: `length(T)1std(x<=T)` brevity penalty.
- **chrF++**: character-level metric which is robust to morphology and spelling variation.
- **Per dialect BLEU**: separate evaluation of Standard, Surati and Kachchi subsets.

Native speakers further manually assess sampled outputs for adequacy and fluency.

Table 3: Final test-set results for Gujarati-to-English translation (English target).

Input Variety	Target	BLEU	chrF++	N
Overall (all Gujarati varieties)	English	61.33	73.04	954
Standard Gujarati	English	62.89	74.33	471
Surati	English	67.80	79.59	149
Kachchi	English	53.16	66.09	334

3.9. Implementation Details

The implementation uses PyTorch 2.0, Transformers 4.39.0 and PEFT 0.7.1 which can be installed using the following commands: We add a custom collator to write out the sequence pairs during PEFT-wrapped seq2seq training without any conflict with decoder-inputs.

4. Results and Discussion

The above section presents quantitative and qualitative evidence for the proposed dialect-conditioned Gujarati-English translation framework. Having provided overall and per-dialect performance to begin with, this analysis goes on to examine just how well optimization is shaping up; what kind of returns one might expect on a given dialect in comparison with others etc. and finally embraces some applications to low-resourced tourism domain with unaltered human performance level.

4.1. Evaluation Setting

The evaluation was carried out on a test set containing 954 pairs of sentences left out from training data. The final model was selected to be one where validation loss was at its minimum. After the test, translation quality was measured using sacreBLEU and chrF++ in the full test set; this measure was then broken down for individual dialect groups. We applied identical decoding settings to translate in every subgroup so that comparison would be fair.

4.2. Quantitative Performance

Under low-resource and dialect-diverse contentions, the proposed method achieves **61.33 BLEU** and **73.04 chrF++**.

While Surati sees the best performance, it is Kachchi that remains the most challenging sub-area. This trend exactly reflects a greater degree of lexical divergence and borrowing intensity in Kachchi outputs. Nonetheless, the results from Kachchi still have practical use for tourism and service-oriented applications.

4.2.1. Training Dynamics

Figure 2 shows stable optimization after switching data over to a common dialect. The shrinking training-validation loss gap indicates controlled overfitting, providing a rationale for selecting a checkpoint that may be generalized.

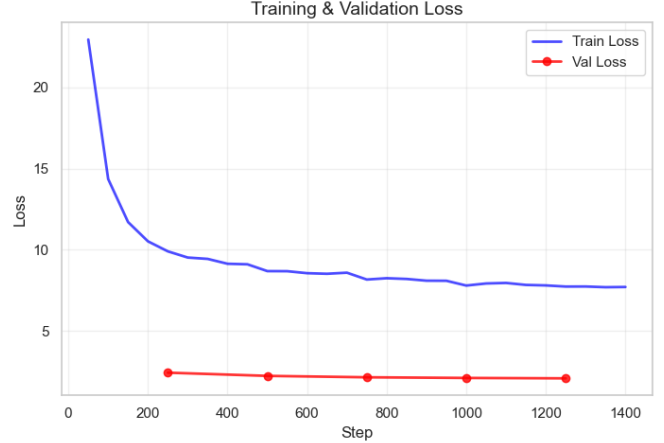


Figure 2: Across the steps of iterations in the optimization phase the training and validation loss curves have very similar shapes. The smooth dip in loss curve along with stable validation behaviour suggests that convergence has been happening.

4.2.2. Per-Dialect Performance Comparison

Figure 3 emphasizes the dialect gap. This jibes with Table 3, and along additional makes it clear that dialect performance order relative between BLEU and chrF++ is stable.

4.2.3. Split and Length Consistency

By looking at the training splits and sentence lengths, we guarantee that the results are not influenced by data distribution artefacts.

Figure 4 validates the evaluation by showing consistent split composition and comparable length profiles. This reduces the possibility that gains reported arise from split artefacts.

4.3. Tourism-Domain Qualitative Analysis

To supplement the automatic metrics presented in Table 4, as a supplementary approach that can be used alongside the automatic metrics, the following tourism samples are extracted from literature in each dialect. Here, the samples give practical tourist use examples including booking tickets for a trip by air or by coach, etc.

Qualitative inspection shows that the intentions it is designed to serve remain strong and that the English generated as a result is fluent across dialects. However, where there are differences they usually appear in lexis rather than syntax, since syntax is well preserved among Kachchi speakers apart from a few errors here and there. Semantic meaning in general is preserved.

4.4. Discussion

The underlying empirical outcomes suggest that three design decisions are essential for accomplishing good performance:

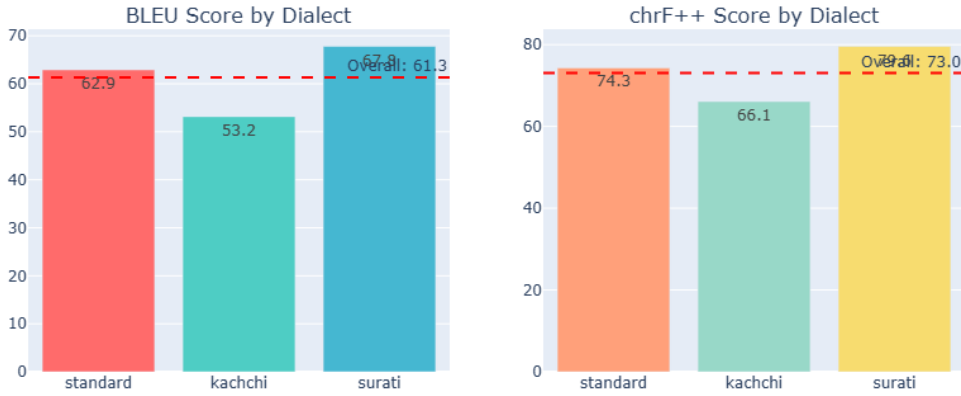


Figure 3: Each dialect is assigned its own overall reference line.

- **Dialect-token conditioning** Here we allow explicit contextual control and thus reduce the ambiguity entailed in heterogeneous inputs.
- **Extended LoRa insertion** From attention and feed-forward modules to feed-forward components of the rudder, such flow type data simplify easier word adaptation when dialectal variance occurs.
- **Unified joint training** This alleviates catastrophic forgetting, so maintaining the standard Gujarati results while ensuring that other dialects of language spoken in the area fall-res against at least western formal settings at Practical Guide level.

Together, these choices create a compromise between the quality of translation and size of the number of parameters needed to ensure that it runs effectively without introducing unnecessary performance overhead on lower capacity machines. At this point the system is suitable for practical tourist-assistance flow, including for example traveling services such as:connections between tourist information centers and transportation companies or travel agencies etc.).

4.5. Error Analysis and Limitations

But there are few restraints that remain intact. First, Kachchi performance is worse than both Standard and Surati, so they aren't sensitive to lexical divergence and borrowing patterns. Second, the existing corpus is very focused on casual and tourism-commerce contexts; larger domain transfer might need additional data. The second concern is that BLEU and chrF++ fail to account for pragmatic adequacy, stylistic tone.

Future efforts will address (i) a larger Kachchi tourism corpus, (ii) specific augmentation for low-frequency dialectal forms, and (iii) stronger human evaluation protocols for adequacy, fluency and cultural appropriateness.

4.6. Section Summary

Table ?? shows, in summary, that the proposed model achieves overall **61.33 BLEU** and **73.04 chrF++**, with best results on Surati and retains well on Standard Gujarati. Consistently, two tables and three major figures of evidence point to the efficacy of dialect conditioning, expanded LoRA adaptation, and unified joint-training for low-resource Gujarati dialect translation.

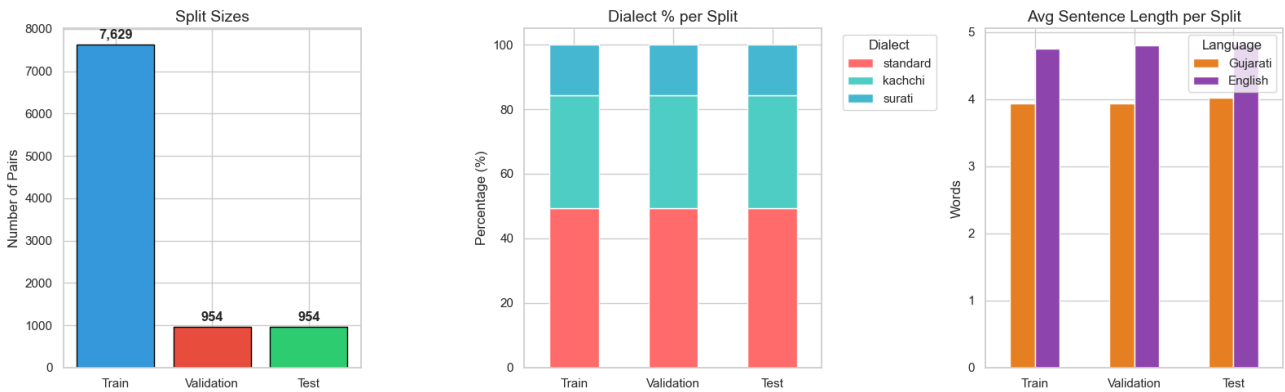


Figure 4: Dataset split and sentence-length analysis across train, validation, and test partitions.

Table 4: Tourism-domain Gujarati-to-English translation samples by dialect.

Dialect	Original Gujarati Input	Reference English	Model Output (English)
Standard Gujarati	મારે આવતીકાલે સોમનાથ મંદિર માટે બે ટિકિટ બુક કરાવવી છે.	I want to book two tickets for Somnath Temple for tomorrow.	I want to book two tickets for Somnath Temple for tomorrow.
Surati	અમોવ સુરતવથી દીવવ જવા માટેવ વહેલીવ હવારમાં ટેક્સીવ જોઈએવ છેવ.	We need a taxi from Surat to Diu early in the morning.	We need a taxi from Surat to Diu early in the morning.
Kachchi	ભુજ એરપોર્ટ સરેદ રણ વેંયો પાંડે, ગાઈડ લેજો ખચે, ભાડું કિતરો થીંધો? જરા વલાયો.	We want to go from Bhuj Airport to the White Rann with a guide; tell us the price.	We want to go from Bhuj Airport to the White Rann with a guide; please tell us the price.

5. Conclusion

In this paper, we proposed a robust low-resource method for conducting Gujarati dialect-to-English translation using NLLB-200 model with basis of dialect conditioning and Expanded-LoRA adaptation. Achieving **61.33 BLEU** and **73.04 chrF++** overall, the proposed model also demonstrated excellent performance across Standard Gujarati as well as Surati and Kachchi inputs. The results demonstrate that the explicit dialect token in conjunction with a unified mixed-dialect training paradigm dramatically reduces catastrophic forgetting and facilitates cross-dialect generalization.

More qualitative examples in the tourism domain provide further evidence of intent preservation and fluency for use cases relevant to a real-world scenario, including ticket booking, transport assistance, and itinerary support. While Kachchi is relatively more challenging due to greater lexical divergence, the obtained performance is operationally useful and highlights an avenue for improvement through focused data expansion.

In summary, we show that a good quality dialect-aware translation can be obtained while avoiding full-model fine-tuning, an approach which is both computationally inexpensive and easily deployable in multilingual scenarios. Future work will be aimed at richer coverage of Kachchi, more domain-adaptive training and a deeper human-centered evaluation along adequacy- fluency-cultural-fidelity.

References

- [1] J. P. Gala *et al.*, “Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages,” *arXiv preprint arXiv:2305.16307*, 2023.
- [2] V. Goyal and D. M. Sharma, “The iit-h gujarati-english machine translation system for wmt19 news translation task,” in *Proceedings of the Fourth Conference on Machine Translation (WMT)*, pp. 191–195, 2019.
- [3] J. Ameta, N. Joshi, and I. Mathur, “Improving the quality of gujarati-hindi machine translation through stemming and pos tagging,” *arXiv preprint arXiv:1307.3310*, 2013.
- [4] J. Thakkar, A. Patel, and A. Desai, “Enhancing neural machine translation for saurashtra gujarati dialect: Low-resource adaptation and morphological analysis,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- [5] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, “Lora: Low-rank adaptation of large language models,” *arXiv preprint arXiv:2106.09685*, 2021.
- [6] S. Biderman, S. Roelofs, N. Prabhu, *et al.*, “Lora learns less and forgets less,” *arXiv preprint arXiv:2405.09673*, 2024.
- [7] Z. Xiao, W. Held, Y. Liu, and D. Yang, “Task-agnostic low-rank adapters for unseen english dialects,” *arXiv preprint arXiv:2311.00915*, 2023.
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.
- [9] M. Johnson, M. Schuster, Q. V. Thorat, N. Shazeer, N. Parmar, and J. Uszkoreit, “Google’s multilingual neural machine translation system: Enabling zero-shot translation,” *Transactions of the Association for Computational Linguistics*, vol. 5, pp. 339–351, 2017.
- [10] M. R. Costa-jussà, J. Cross, Y. Duan, *et al.*, “No language left behind: Scaling human-centered machine translation,” *arXiv preprint arXiv:2207.04672*, 2022.
- [11] M. Liu, M. Sallmen, and A. Joly, “Towards optimal rank selection in low-rank adaptation of large language models,” *arXiv preprint arXiv:2308.04993*, 2023.
- [12] H. Barmandah and K. Al-Mansour, “Saudi dialect lora: Fine-tuning large language models for dialectal arabic generation and translation,” *arXiv preprint arXiv:2508.13525*, 2025.
- [13] A. Kantharuban, I. Vulić, and A. Korhonen, “Quantifying the dialect gap and its correlates across languages,” *arXiv preprint arXiv:2310.15135*, 2023.
- [14] L. Xue, N. Constant, A. Roberts, *et al.*, “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, pp. 483–498, 2021.

- [15] R. Sennrich, B. Haddow, and A. Birch, “Improving neural machine translation models with monolingual data,” in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics*, pp. 86–96, 2016.
- [16] J. Kirkpatrick, R. Pascanu, N. C. Rabinowitz, *et al.*, “Overcoming catastrophic forgetting in neural networks,” *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [17] R. Caruana, “Multitask learning,” in *Machine Learning*, vol. 28, pp. 41–75, Kluwer Academic Publishers, 1997.
- [18] M. Popović, “chrf++: word n-gram level evaluation metric,” in *Proceedings of the Second Conference on Machine Translation*, pp. 550–555, 2017.
- [19] R. Rei, C. Stewart, A. C. Farinha, and A. Lavie, “Comet: A neural framework for mt evaluation,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2685–2702, 2020.
- [20] UCLA Language Materials Project, “Ucla language materials project: Gujarati,” 2005. Comprehensive linguistic analysis of Gujarati dialects.
- [21] P. Joshi, S. Santy, F. Buettner, and A. Zeller, “Towards code-mixed machine translation,” in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, pp. 4847–4854, 2020.
- [22] G. I. Winata, M. Sap, A. F. Solares, *et al.*, “Code-switched language models using multilingual transformers,” in *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics*, pp. 3934–3948, 2021.