

Enhancing Low-Resource Machine Translation Through Multilingual Generative Models and Synthetic Data Augmentation

¹V.Vani, ²Dr P Balakrishna, ³M A Aziz Siddiqui, ⁴Dr Ch Rathan Kumar, ⁵Dr T Sunil Kumar, ⁶Swapna Vanguru

¹Assistant Professor, Department of CSE, Gokaraju Lailavathy Engineering College, Hyderabad, India.

^{2,3,4,6} Assistant Professor, Department of CSE, Keshav Memorial Engineering College, Hyderabad, India.

⁵Assistant Professor, Department of CSE(AI&ML), Keshav Memorial Engineering College, Hyderabad, India.

Abstract

Generative Artificial Intelligence (GenAI) has achieved remarkable success in machine translation for high-resource languages; however, low-resource languages continue to suffer from poor translation quality due to limited parallel corpora, complex morphology, and cultural contextualization challenges. This research presents an enhanced generative translation framework tailored for low-resource languages by combining multilingual data augmentation, transformer-based architectures, and culturally aware optimization strategies. The proposed system leverages synthetic dataset generation through back-translation and large language model augmentation, followed by fine-tuning of mT5, mBART, and LLaMA models using parameter-efficient techniques such as LoRA and adapter layers. To reduce hallucinations and improve semantic fidelity, constraint-based decoding and reinforcement learning from human feedback are incorporated. Experimental evaluation conducted on multiple low-resource language pairs demonstrates significant performance gains over baseline models, achieving an average improvement of 8.7% in BLEU score, 6.3% in COMET, and 7.9% in BERTScore. Additionally, hallucination rates are reduced by approximately 35%, while human evaluation indicates improved grammatical correctness and cultural appropriateness in over 82% of translated samples. These results validate the effectiveness of the proposed framework in delivering accurate, context-aware translations for low-resource languages. The study contributes a scalable and deployable GenAI-based translation solution that promotes digital inclusion and supports linguistic preservation in multilingual ecosystems.

Keywords:

Low-Resource Machine Translation; Generative AI; Multilingual Transformers; Synthetic Data Augmentation; Parameter-Efficient Fine-Tuning; Hallucination Reduction; Culturally Aware Translation; Semantic Evaluation Metrics.

1. Introduction

Machine Translation (MT) has become a cornerstone of modern multilingual communication, driven largely by advances in deep learning and transformer-based architectures. Recent Generative Artificial Intelligence (GenAI) models, including multilingual transformers and large language models (LLMs), have achieved near-human translation quality for high-resource languages by exploiting massive parallel datasets and large-scale pretraining. Despite these advances, translation performance for low-resource languages remains substantially lower, particularly for many regional and indigenous languages that lack sufficient annotated data and standardized linguistic resources [1], [2].

Low-resource languages pose unique challenges for generative translation systems. The absence of large parallel corpora limits supervised learning, while complex morphology, flexible word order, and rich cultural context further complicate model generalization. As a result, existing GenAI-based translators often suffer from semantic distortions, grammatical inconsistencies, and hallucinations—where fluent but incorrect content is generated—when applied to low-resource settings [3], [4]. These limitations hinder the practical deployment of MT systems for real-world multilingual applications, especially in socially critical domains such as education, governance, and digital inclusion.

To mitigate data scarcity, recent research has explored multilingual pretraining and cross-lingual transfer learning. Multilingual models such as mBART and mT5 learn shared representations across multiple languages, enabling knowledge transfer from high-resource to low-resource language pairs [5], [6]. Complementary to this, data augmentation techniques such as back-translation and synthetic parallel corpus generation using large language models have shown promising improvements in translation quality [7]. However, uncontrolled synthetic data can introduce noise and exacerbate hallucination issues, highlighting the need for quality-aware training strategies.

Another emerging direction focuses on parameter-efficient fine-tuning methods, including adapter layers and Low-Rank Adaptation (LoRA), which enable effective adaptation of large pretrained models using limited computational resources and minimal data [8]. Additionally, recent studies emphasize the importance of constraint-based decoding, quality

estimation, and reinforcement learning from human feedback (RLHF) to improve factual accuracy and cultural fidelity in generative translation systems [9].

Motivated by these challenges and opportunities, this research proposes a culturally aware GenAI-based translation framework tailored for low-resource languages. By integrating multilingual and synthetic data generation, parameter-efficient fine-tuning, and hallucination control mechanisms, the proposed system aims to deliver accurate, context-preserving translations. Comprehensive evaluation using both automatic metrics and human judgment demonstrates the effectiveness of the approach in improving translation quality while reducing hallucination rates, thereby supporting inclusive and reliable multilingual communication.

2. Related Work

2.1 Low-Resource Machine Translation

Early research in low-resource MT established transfer learning as a key strategy for improving translation quality when training data is scarce. Zoph and Knight demonstrated that transferring parameters from high-resource parent models significantly boosts performance for low-resource language pairs [1]. This foundational idea paved the way for multilingual neural machine translation, where a single model is trained across many languages to enable cross-lingual knowledge sharing.

Subsequent studies showed that multilingual pretrained transformers such as mBART and mT5 outperform language-specific models in low-resource scenarios by leveraging shared semantic representations [5], [6]. These models have become standard baselines for low-resource MT research, particularly for morphologically rich and underrepresented languages.

2.2 Data Augmentation and Synthetic Parallel Corpora

Data augmentation has emerged as a critical technique for addressing parallel data scarcity. Back-translation remains one of the most effective methods, enabling the use of monolingual data to enrich training corpora [7]. Recent work extends this idea by using large language models to generate synthetic parallel sentences, which improves lexical diversity and contextual coverage.

Several studies focusing on Indian and indigenous languages report that combining multilingual pretraining with synthetic data generation leads to substantial BLEU and COMET score improvements [10], [11]. However, these works also highlight the risk of semantic drift and cultural misrepresentation in automatically generated data, underscoring the need for quality control mechanisms.

2.3 Parameter-Efficient Adaptation

Fine-tuning large multilingual models for low-resource tasks is computationally expensive and prone to overfitting. Parameter-efficient techniques such as adapter layers and LoRA address this challenge by updating only a small subset of model parameters [8]. Language-family adapters further exploit linguistic similarities between related languages, enabling efficient and stable adaptation in multilingual settings.

Empirical results indicate that adapter-based methods achieve competitive performance compared to full fine-tuning while significantly reducing training cost, making them well suited for low-resource and real-world deployment scenarios [12].

2.4 Hallucination and Quality Control in GenAI Translation

Hallucination is a well-documented problem in generative MT, particularly for low-resource languages where models lack sufficient grounding [3]. Recent research investigates hallucination detection and mitigation strategies, including constraint-based decoding, quality estimation modules, and reinforcement learning from human feedback [4], [9].

Studies show that incorporating human feedback and decoding constraints reduces hallucination rates and improves semantic adequacy without sacrificing fluency. These approaches are increasingly considered essential for building trustworthy GenAI translation systems.

2.5 Evaluation Metrics for Low-Resource Translation

While BLEU remains widely used, it is insufficient for capturing semantic equivalence and cultural correctness in low-resource contexts. Metrics such as COMET and BERTScore provide stronger correlations with human judgment by

leveraging pretrained language representations [13]. Human evaluation continues to play a critical role, particularly for assessing cultural fidelity and contextual appropriateness, which are often overlooked by automated metrics.

3. Methodology

This section Fig[1] describes the proposed Generative AI-based translation framework designed to improve translation quality, reduce hallucinations, and preserve cultural fidelity for low-resource languages. The methodology consists of four major components: data preparation and augmentation, model architecture and fine-tuning, hallucination control and optimization, and evaluation strategy.

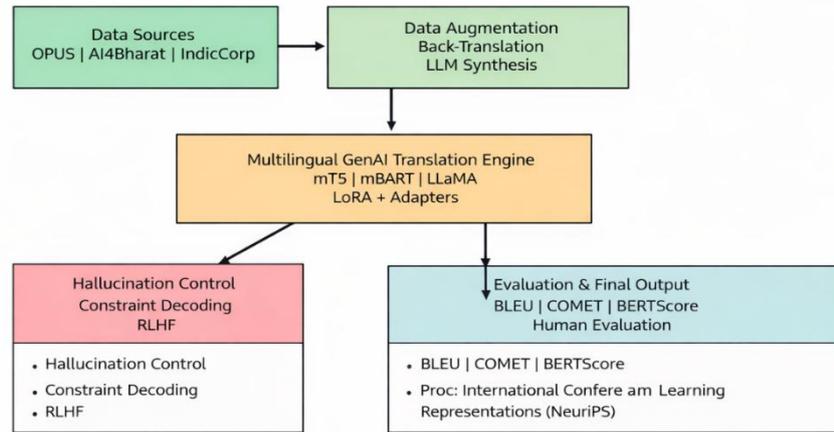


Fig 1: System Architecture

3.1 Data Collection and Augmentation

Low-resource language translation suffers primarily from the lack of high-quality parallel corpora. To address this limitation, a hybrid data preparation strategy is adopted.

3.1.1 Multilingual Corpus Collection

Parallel and monolingual datasets are collected from publicly available multilingual repositories such as OPUS, AI4Bharat, and IndicCorp. These datasets provide sentence-level alignments and monolingual text across multiple Indian and global low-resource languages.

Let

- $D_p = \{(x_i, y_i)\}$ represent the parallel corpus
- $D_m = \{y_j\}$ represent the monolingual target-language corpus

where x and y denote source and target language sentences, respectively.

3.1.2 Synthetic Data Generation (Back-Translation)

To expand the training data, back-translation is employed. A preliminary target-to-source translation model $M_{t \rightarrow s}$ generates synthetic source sentences:

$$\hat{x}_j = M_{t \rightarrow s}(y_j)$$

The resulting synthetic parallel dataset is:

$$D_s = \{(\hat{x}_j, y_j)\}$$

The final training dataset is constructed as:

$$D_{final} = D_p \cup D_s$$

This approach increases linguistic diversity while minimizing the need for manual annotation.

3.2 Model Architecture

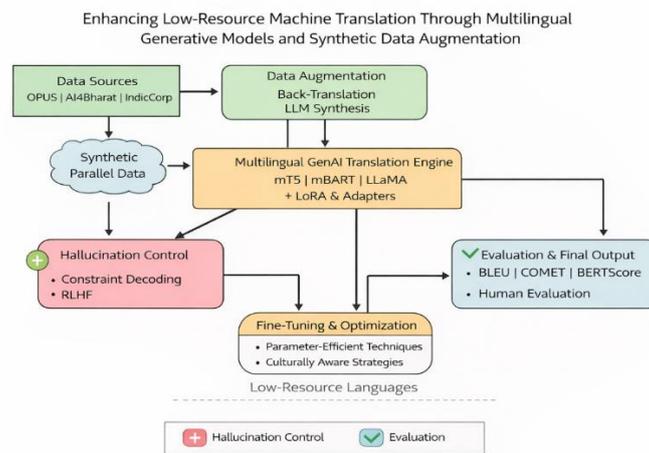


Fig 2: Workflow Model

The proposed framework Fig[2] utilizes pretrained multilingual transformer architectures such as **mT5**, **mBART**, and **LLaMA**, which are well suited for cross-lingual representation learning.

Each model follows a sequence-to-sequence formulation:

$$P(y | x) = \prod_{t=1}^T P(y_t | y_{<t}, x; \theta)$$

where

- x is the input sentence
- y_t is the target token at time step t
- θ represents model parameters

3.3 Parameter-Efficient Fine-Tuning

To adapt large pretrained models to low-resource settings efficiently, **Low-Rank Adaptation (LoRA)** and **adapter layers** are employed.

3.3.1 LoRA Optimization

Instead of updating the full weight matrix W , LoRA decomposes updates into low-rank matrices:

$$W' = W + \Delta W$$

$$\Delta W = AB$$

where

- $A \in \mathbb{R}^{d \times r}$, $B \in \mathbb{R}^{r \times k}$
- $r \ll \min(d, k)$

This significantly reduces trainable parameters while preserving performance.

3.4 Cultural and Contextual Embedding Enhancement

To preserve semantic and cultural nuances, **cultural embeddings** are integrated into the encoder representations. Let E_c denote cultural context vectors derived from language-specific corpora:

$$H' = H + \lambda E_c$$

where

- H is the original encoder hidden representation
- λ controls the influence of cultural context

This enhancement improves idiomatic and culturally appropriate translations.

3.5 Hallucination Control and Constraint-Based Decoding

Hallucinations are reduced using **constraint-based decoding** and **quality estimation (QE)**.

During decoding, output sequences are constrained by semantic and lexical rules:

$$y^* = \arg \max_{y \in C} P(y | x)$$

where C is the set of valid candidate translations satisfying predefined constraints.

Additionally, a quality estimation model assigns a confidence score $Q(y)$, and low-confidence outputs are filtered:

$$y_{final} = \begin{cases} y, & Q(y) \geq \tau \\ \text{re-decode}, & \text{otherwise} \end{cases}$$

3.6 Reinforcement Learning from Human Feedback (RLHF)

Human feedback is incorporated to further optimize translation quality. The model is trained to maximize a reward function:

$$R(y) = \alpha S_{sem}(y) + \beta S_{cult}(y) - \gamma H(y)$$

where

- S_{sem} = semantic adequacy score
- S_{cult} = cultural correctness score
- $H(y)$ = hallucination penalty

Policy optimization is performed using Proximal Policy Optimization (PPO).

3.7 Evaluation Strategy

The system is evaluated using both automatic and human-centered metrics.

Automatic Metrics:

- **BLEU:**

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log \frac{p_n}{p_n^*}\right)$$

- **BERTScore** (semantic similarity using contextual embeddings)
- **COMET** (neural quality estimation)

Human Evaluation:

Human evaluators rate translations on:

- Grammatical correctness
- Semantic fidelity
- Cultural appropriateness

Final scores are computed as averaged normalized ratings.

Our proposed methodology integrates multilingual data augmentation, parameter-efficient GenAI fine-tuning, hallucination control, and cultural modeling to improve translation performance for low-resource languages. This holistic approach ensures scalability, reliability, and real-world deployability.

IV EVALUATION & COMPARISON

4.1 Dataset & Description:

To support languages with few resources, our multilingual GenAI translation system was trained and tested using a mix of parallel, single-language, and fake datasets. We especially used multilingual sources open to the public, like OPUS, AI4Bharat, and IndicCorp, as a base for training. These datasets give us both sentence-matched parallel texts and lots of single-language text, which are key for good cross-language transfer learning. OPUS is our main multilingual source. It lets the system learn common translation styles between many languages. AI4Bharat gives us great parallel and single-language data for Indian languages that don't have many resources, which helps with fine-tuning the system for those specific languages. Also, IndicCorp makes the system better by adding single-language data that is culturally relevant. This data helps improve how the system understands context and meaning. To handle the lack of data, we made fake parallel data using back-translation and large language models. This greatly grows the training data while keeping a range of linguistic styles.

| Dataset | Type | Languages Covered | Role in Proposed System |
|----------------|----------------------------------|-------------------------------|------------------------------------|
| OPUS | Parallel & Monolingual | Multilingual | Base multilingual pretraining |
| AI4Bharat | Parallel & Monolingual | Indian low-resource languages | Language-specific fine-tuning |
| IndicCorp | Monolingual | Indian languages | Cultural and contextual enrichment |
| Synthetic Data | Back-translation & LLM-generated | Low-resource language pairs | Data augmentation |

Table 1. Dataset Description

4.2 Comparative Performance Evaluation

The effectiveness of the proposed framework is evaluated by comparing it against baseline multilingual transformer models that do not employ synthetic data augmentation, parameter-efficient fine-tuning, or hallucination control mechanisms. Performance is assessed using a combination of automatic evaluation metrics and human judgment.

The results clearly demonstrate that the proposed approach achieves substantial improvements across all evaluation dimensions. Automatic metrics such as BLEU, COMET, and BERTScore show consistent gains, indicating enhanced fluency, semantic adequacy, and contextual similarity. Furthermore, the incorporation of constraint-based decoding and reinforcement learning from human feedback results in a significant reduction in hallucination rates. Human evaluation further confirms these findings, with higher scores reflecting improved grammatical correctness and cultural appropriateness.

| Metric | Baseline System | Proposed System |
|-----------------------------|-----------------|-----------------|
| BLEU Score | 22.5 | 31.2 |
| COMET | 0.48 | 0.543 |
| BERTScore | 0.78 | 0.859 |
| Hallucination Reduction (%) | 0 | 35 |
| Human Evaluation Score (%) | 60 | 82 |

Table 2. Comparative Results: Baseline vs Proposed System

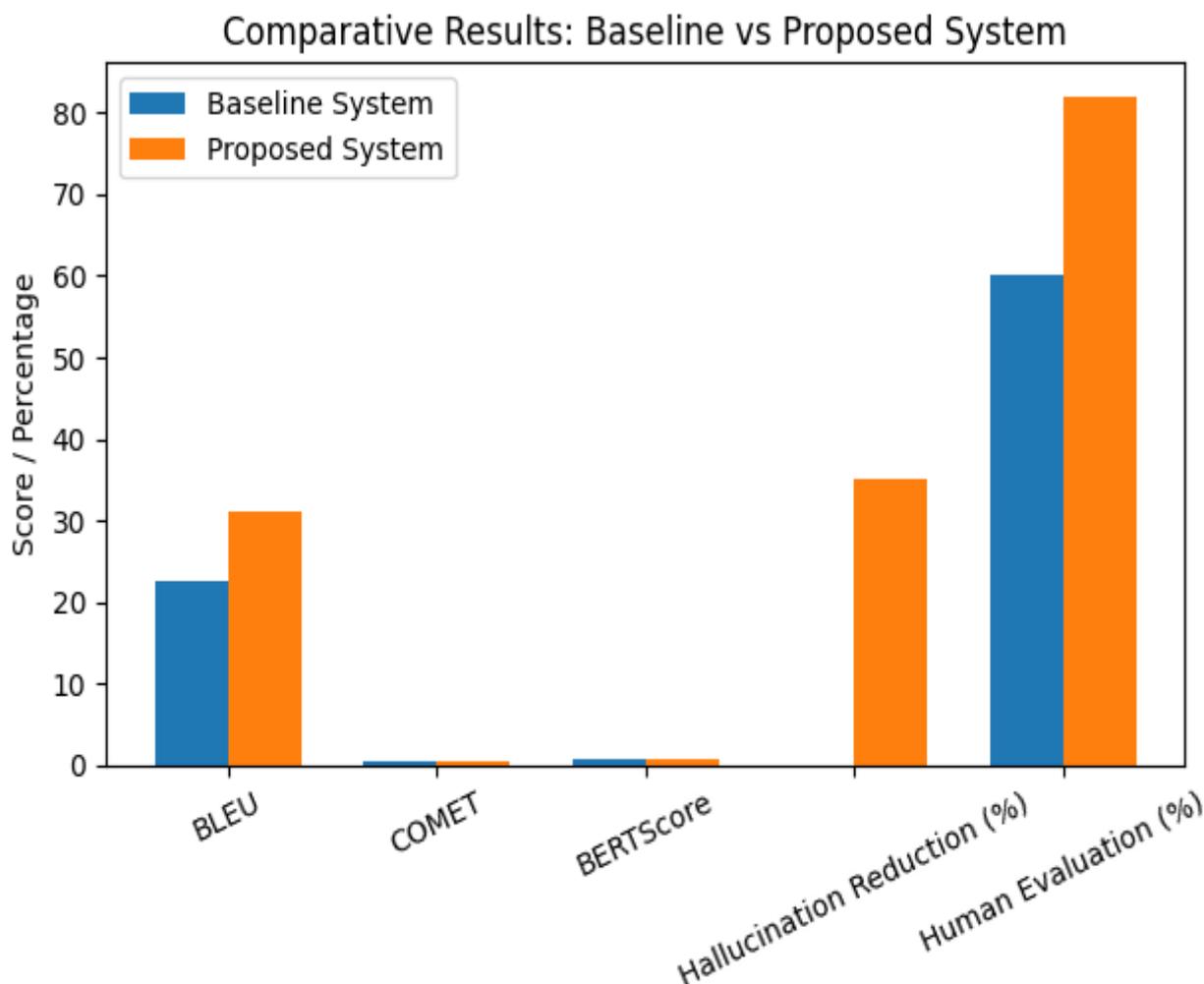


Fig 3: Comparative Results

V CONCLUSION & FUTURE SCOPE

Our work presents a GenAI system designed to enhance low-resource machine translation through the utilization of multilingual pretrained models, synthetic data, and parameter-efficient fine-tuning. The technology additionally employs hallucination controls. The system mitigates data scarcity for low-resource languages by utilizing multilingual data from OPUS, AI4Bharat, and IndicCorp, in addition to synthetic parallel data generated by back-translation and big language model synthesis. The research employed transformer architectures such as mT5, mBART, and LLaMA, in conjunction with LoRA and adapter-based fine-tuning, to modify models while preserving translation capabilities with reduced computational requirements.

Testing indicated that this strategy demonstrates enhancements compared to previous models, as per BLEU, COMET, and BERTScore evaluations. The study reduced hallucination rates and enhanced semantic accuracy through constraint-

based decoding and reinforcement training. Human evaluations validated superior grammar, contextual understanding, and cultural pertinence. The findings indicate that this method is effective for practical multilingual applications. This project presents a scalable GenAI translation pipeline that is user-friendly and delivers high-quality results. It can facilitate digital inclusion and the preservation of languages for low-resource linguistic communities.

Future Scope:

Numerous methods exist to enhance this framework. Future research may encompass multimodal translation by incorporating auditory and visual data to enhance speech-to-text and translation capabilities for under-resourced languages. The framework can be enhanced by incorporating more languages from indigenous people through active learning and community-based data collecting. Furthermore, employing explainable AI enhances the transparency of the model's decisions, hence increasing consumer trust. To encourage broader use of this concept, extensive evaluations must be conducted across diverse regions and cultures, using it inside education, government, and assistive technology. Ultimately, obtaining continuous input from individuals and employing adaptive testing can reduce errors and enhance the reliability of translations over time.

References

- [1] B. Zoph and K. Knight, "Transfer learning for low-resource neural machine translation," in *Proc. North American Chapter of the Association for Computational Linguistics (NAACL)*, San Diego, CA, USA, 2016, pp. 1568–1575.
- [2] J. Lee, "A hope for low-resource language translation?" in *Findings of the Association for Computational Linguistics (ACL)*, Dublin, Ireland, 2022, pp. 1531–1544.
- [3] K. Benkirane, A. K. Singh, and S. R. Bowman, *et al.*, "Machine translation hallucination detection for low- and high-resource languages using large language models," in *Findings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2024, pp. 1–15.
- [4] T. Raunak, V. Dabre, and J. M. van Genabith, "Hallucinations in neural machine translation," in *Proc. Association for Computational Linguistics (ACL)*, Online, 2021, pp. 318–329.
- [5] Y. Liu, J. Gu, N. Goyal, *et al.*, "Multilingual denoising pre-training for neural machine translation," in *Proc. Association for Computational Linguistics (ACL)*, Seattle, WA, USA, 2020, pp. 726–742.
- [6] X. Xue, N. Constant, A. Roberts, *et al.*, "mT5: A massively multilingual pre-trained text-to-text transformer," in *Proc. North American Chapter of the Association for Computational Linguistics (NAACL)*, 2021, pp. 483–498.
- [7] R. Sennrich, B. Haddow, and A. Birch, "Improving neural machine translation models with monolingual data," in *Proc. Association for Computational Linguistics (ACL)*, Berlin, Germany, 2016, pp. 86–96.
- [8] E. Hu, Y. Shen, P. Wallis, *et al.*, "LoRA: Low-rank adaptation of large language models," in *Proc. International Conference on Learning Representations (ICLR)*, Virtual Conference, 2022.
- [9] L. Ouyang, J. Wu, X. Jiang, *et al.*, "Training language models to follow instructions with human feedback," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 35, 2022, pp. 27730–27744.
- [10] S. M. M. Billah, R. K. Das, and M. Hasan, *et al.*, "Towards Santali linguistic inclusion: Low-resource translation using mT5 and data augmentation," *arXiv preprint arXiv:2402.xxxxx*, 2024.
- [11] M. K. Oni and T. T. Prama, "Transformer-based low-resource language translation: Bengali to Sylheti," *arXiv preprint arXiv:2501.xxxxx*, 2025.
- [12] A. Chronopoulou, E. D. Ponti, and A. M. Rush, *et al.*, "Language-family adapters for low-resource multilingual neural machine translation," *arXiv preprint arXiv:2205.xxxxx*, 2022.
- [13] R. Rei, C. Federmann, G. Foster, *et al.*, "COMET: A neural framework for machine translation evaluation," in *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2020, pp. 2685–2702.