

“WHEN DATA IS SCARCE, PROMPT SMARTER”...APPROACHES TO GRAMMATICAL ERROR CORRECTION IN LOW-RESOURCE SETTINGS

BHASHA TASK-1 INDICGEC DEMONSTRATION

Somsubhra De¹, Harsh Kumar¹ and Arun Prakash A^{1,2}

¹ IIT Madras, ² AI4Bharat

- Grammatical correctness essential for clear and effective communication ...
- In the NLP field, developing systems that can automatically detect and fix sentence level grammatical errors is a significant area of research
- GEC has seen great success for languages like English, **applying these advancements to many other low-resource Indic languages has been challenging, why?**



Given a parallel corpus $D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N$, where $x^{(i)}$ is a noisy input sentence and $y^{(i)}$ is its grammatically corrected counterpart, the model is trained to minimize the negative log-likelihood (NLL) of the target sequence conditioned on the input:

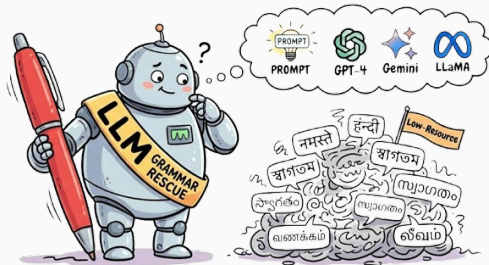
$$\mathcal{L}_{\text{GEC}}(\theta) = - \sum_{i=1}^N \sum_{t=1}^{|y^{(i)}|} \log P_{\theta} \left(y_t^{(i)} \mid y_{<t}^{(i)}, x^{(i)} \right) \quad (1)$$

What are we dealing with?



Figure 1: Examples from the GEC task dataset. Input sentence **X** (with errors in red) - ground truth **✓** (with corrections in blue) pairs. Error types have been mentioned, based on our understanding.

How Well Can LLMs Come to Indic Grammar Rescue?



We utilize three large language models – GPT4.1 mini, Gemini-2.5-Flash and Llama-4-Maverick 17B-128EInstruct in *zero and few-shot* prompting paradigms. These models are employed as *instruction-following LLMs with role-based prompts*.

Additionally for Hindi, we fine-tuned (we adopt LoRA for PEFT) the **Sarvam-M24B2** multilingual model using Hi-GEC dataset.

Prompt

""You are a <Language> Grammatical Error Correction assistant, in low resource settings. Your task is to accurately identify and correct grammatical errors in the given <Language> sentence. Correct all types of grammatical errors:

Verb usage: Correct conjugation, tense, aspect, and agreement with the subject.,

Pronouns: Usage of proper personal, possessive, and reflexive pronouns.,

Prepositions: Correct use of postpositions or prepositions in context.,

Fix spelling mistakes, diacritic marks (matras), and punctuation errors.,






Gender and number agreement: Ensure adjectives, nouns, and verbs match in gender (masculine/feminine) and number (singular/plural).,

The output should be ONLY the CORRECTED sentence, without any extra text or explanation. *If the input is already correct, return it unchanged.* Please ensure the corrections follow the rules and preserve the intended meaning.

Below are 10 random sentences for your reference.""

RESULTS

We ranked **1st** in Tamil (GLEU: 91.57) and Hindi (GLEU: 85.69)
2nd in Telugu (GLEU: 85.22), 4th in Bangla (GLEU: 92.86) and 5th in Malayalam (GLEU: 92.97).

Model	TAM			MAL			Hi		
	GLEU	$F_{0.5}$	BERT-score	GLEU	$F_{0.5}$	BERT-score	GLEU	$F_{0.5}$	BERT-score
 Gemini-2.5-Flash (<i>fs</i>)	91.57	87.82	97.83	92.97	88.48	97.89	84.61	88.01	95.69
 GPT-4.1 mini (<i>fs</i>)	86.00	78.97	96.52	91.78	84.72	97.08	85.69	87.86	95.76
 GPT-4.1 mini (<i>zs</i>)	85.51	78.45	96.38	92.34	84.62	97.24	85.37	87.80	95.53
 LLaMA-4 maverick (<i>zs</i>)	88.70	81.50	96.84	92.68	85.38	97.20	83.10	86.04	94.64
 LLaMA-4 maverick (<i>fs</i>)	85.62	77.75	95.98	90.75	83.22	96.65	85.37	87.35	95.56





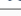
Model	BN			TEL		
	GLEU	$F_{0.5}$	BERT-score	GLEU	$F_{0.5}$	BERT-score
 Gemini-2.5-Flash (<i>fs</i>)	92.23	89.61	97.10	84.16	76.96	94.92
 GPT-4.1 mini (<i>fs</i>)	92.86	89.27	97.35	85.22	77.68	95.28
 GPT-4.1 mini (<i>zs</i>)	91.62	86.98	96.79	84.74	76.75	95.15
 LLaMA-4 maverick (<i>zs</i>)	90.39	86.48	96.30	83.01	74.20	94.28
 LLaMA-4 maverick (<i>fs</i>)	92.00	88.02	97.19	82.02	74.28	94.08

Figure 2: Performance of different approaches on the test set across languages. Highlighted cells indicate the best-performing model for each language, while underlined values denote the overall best score in the task.

The finetuned Sarvam-M significantly failed to capture the correct edits, achieving only **13.81** GLEU in Hindi.

- Experiments revealed that even simple prompting strategies could lead to impressive results, even with limited data.
- LLMs when guided by well-designed prompts, substantially outperform fine-tuned Indic-language models like Sarvam-22B thereby illustrating the exceptional multilingual generalization capabilities of contemporary LLMs for GEC.
- Shows that instead of needing massive amounts of language-specific data to train models from scratch, we can leverage the existing knowledge within large, general-purpose LLMs.
- Where did the LLMs fail? We did some qualitative analysis and found out recurring patterns.
- A high Fertility Score has two direct consequences for GEC: increased latency and Context Window Reduction. To quantify this impact, we looked into the tokenizer fertility.

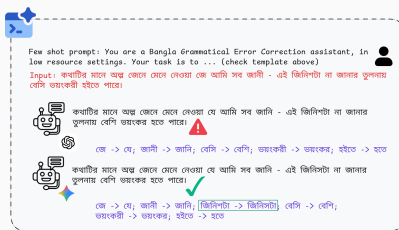


Figure 3: Comparison of model outputs on a **multi-correction** example (Transl. Knowing only a little about something and assuming that I know everything can be more dangerous than not knowing at all.) from test set. *Gemini’s output fully aligns with the gold standard, while GPT omits one necessary edit.*

Language	Script Family	GPT-4.1 Mini	Gemini 2.5 Flash	Llama 4 Maverick
HI	Devanagari	1.44	1.31	1.55
BN	Eastern Nagari	2.32	1.76	2.77
TAM	Dravidian	3.09	2.54	<u>5.88</u>
TEL	Dravidian	2.97	2.87	4.32
MAL	Dravidian	3.20	3.10	4.58

Figure 4: Cross-Model Tokenizer Fertility Comparison

INPUT SENTENCE	CHARS	WORDS	TOKENS
HI हम बस पीछे-पीछे चलते-चलते तमिऴों में घुस ही रहे थे कि हमसे एक लड़की को बहुत गुस्सा कपड़े सोचते हुए देखा।	104	22	35
MAL ഇന്ന് സംസ്ഥാനമാർ തമ്മിൽ തടവിലാക്കുന്ന ഈ കാലഘട്ടത്തിൽ നാം ലക്ഷക്കണക്കിന് നിന്നും ബുദ്ധിമുട്ടുകൾക്കുതീർക്കേണ്ടതുണ്ട്.	104	11	34
TAM திரும்ப வேறாட்டத்துக்கு சொள்ள டாக்ஸி கிடைக்காததால், அவர்கள் பெற்றோருடன் என் பெற்றோர்களை அனுப்பி வைத்தீர்கள்.	104	11	28
BN উহা বিলাস বাস বাস মিচ , এই কোন মিচি বাস বাস মিচি উহা বাস বাস মিচি বাস মিচি।	104	19	24
TEL మన పెన్ పనితీరు క్లుప్తం కూడా పెక్కుతే తీవ్రతతో ఉన్న క్లుప్తం పనితీరు పెక్కు తీవ్రతతో ఉంటుంది.	104	16	34

Figure 5: Tokenization density across the architectures

Overall, the study highlights the immense potential of large language models and prompt-based techniques for grammatical error correction in low-resource settings.



For the paper, data and codes, please scan the QR.