

Benchmarking Bengali Dialectal Bias: A Multi-Stage Framework Integrating RAG-Based Translation and Human-Augmented RLAIIF

Anonymous ACL submission

Abstract

Large language models (LLMs) frequently exhibit performance biases against regional dialects of low-resource languages. However, frameworks to quantify these disparities remain scarce. We propose a two-phase framework to evaluate dialectal bias in LLM question-answering across nine Bengali dialects. First, we translate and gold-label standard Bengali questions into dialectal variants adopting a retrieval-augmented generation (RAG) pipeline to prepare 4,000 question sets. Since traditional translation quality evaluation metrics fail on unstandardized dialects, we evaluate fidelity using an LLM-as-a-judge, which human correlation confirms outperforms legacy metrics. Second, we benchmark 19 LLMs across these gold-labeled sets, running 68,395 RLAIIF evaluations validated through multi-judge agreement and human fallback. Our findings reveal severe performance drops linked to linguistic divergence. For instance, responses to the highly divergent Chittagong dialect score 5.44/10, compared to 7.68/10 for Tangail. Furthermore, increased model scale does not consistently mitigate this bias. We contribute a validated translation quality evaluation method, a rigorous benchmark dataset, and a Critical Bias Sensitivity (CBS) metric for safety-critical applications.

1 Introduction

Large Language Models (LLMs) have achieved remarkable performance across diverse NLP tasks, yet their behavior on dialectal variants of low-resource languages remains poorly understood (Fleisig et al., 2024; Hofmann et al., 2024). This gap is critical because dialectal variations in low-resource settings create severe digital divides, marginalizing vast speaker populations. We explore this broader challenge using Bengali as a representative case study, as its regional dialects spo-

ken by millions diverge substantially from the standardized written form (Wasi et al., 2025).

Such dialectal variations, whether in Bengali (e.g., Chittagong, Sylhet) or other low-resource languages like Arabic, exhibit distinct phonological, lexical, and syntactic features that confuse LLMs trained predominantly on standard forms (Sami et al., 2025; Jawad et al., 2025). Unlike standardized language that benefits from large training corpora, dialectal variants face severe data scarcity, creating potential disparities in model comprehension and response quality (Chang et al., 2024; Sindhujan et al., 2025).

We address this challenge through a two-stage framework: **(1)** Adopting a high-performance RAG-based translation pipeline (Sami et al., 2025) that translates standard Bengali questions into dialectal variants for benchmark construction, and **(2)** An RLAIIF-inspired evaluation framework, with human fallback and multi-judge validation that quantifies LLM performance disparities across dialects using validated scoring rubrics.

Our contributions are:

- A human-validated translation evaluation methodology for standard-to-dialect Bengali, demonstrating the catastrophic failure of traditional metrics
- A gold-standard benchmark dataset of 4,000 questions across 9 Bengali dialects for bias evaluation in LLM question-answering
- An RLAIIF bias evaluation framework with Chain-of-Thought enabled rubrics, validated through multi-judge agreement analysis (Lin (1989)’s Concordance Correlation Coefficient (CCC) = 0.861), and human inspection
- A comprehensive benchmark of 19 open-weight LLMs across 9 dialects (68,395 evalu-

079	ations), revealing systematic bias patterns	
080	• A novel Critical Bias Sensitivity (CBS) met-	129
081	ric for safety-critical applications requiring	130
082	high judge agreement on critical bias cases	131
083	2 Related Works	132
084	2.1 Bias in Large Language Models	133
085	Bias in LLMs manifests across multiple dimen-	134
086	sions including gender, race, religion, and so-	135
087	cioeconomic status (Gallegos et al., 2024). Re-	136
088	cent work established frameworks for systematic	137
089	bias evaluation (Liang et al., 2023), though dialect-	138
090	al bias remains understudied compared to demo-	139
091	graphic dimensions.	140
092	Fleisig et al. (2024) demonstrated that ChatGPT	141
093	exhibits linguistic bias, providing lower-quality re-	142
094	sponses to users of non-standard English dialects.	143
095	Hofmann et al. (2024) found that dialect prejudice	144
096	in LLMs predicts discriminatory decisions about	145
097	character, employability, and criminality. These	146
098	findings motivate our investigation into dialectal	
099	bias for Bengali.	
100	2.2 Bengali NLP and Dialectal Variation	147
101	Bengali NLP research has expanded significantly,	148
102	with benchmarks like BenLLMEval (Kabir et al.,	149
103	2024) evaluating LLM capabilities. While new	
104	dialectal resources are emerging, such as Vashan-	
105	tor (Faria et al., 2025) for translation, BanglaD-	
106	ial (Mahi et al., 2025) for identification, and	
107	DIALTSA-BN (Jawad et al., 2025) for down-	
108	stream benchmarks, dialectal variation remains	
109	broadly underexplored. Alongside resource crea-	
110	tion, bias auditing has revealed systematic reli-	
111	gious dialect disparities (Wasi et al., 2025) and	
112	broader socio-cultural biases (Sadhu et al., 2025,	
113	2024) in Bengali LLMs. Our work extends this	
114	line by specifically focusing on regional dialectal	
115	bias. Furthermore, while recent RAG-based di-	
116	allect translation models (Sami et al., 2025) show	
117	promise, their evaluation relied heavily on tradi-	
118	tional token-matching metrics (BLEU (Papineni	
119	et al., 2002), WER, ChrF (Popović, 2015), and	
120	BERTScore (Zhang* et al., 2020)). Because these	
121	metrics fail to capture true semantic equivalence	
122	in highly agglutinative languages like Bengali (Re-	
123	iter, 2018; Lee et al., 2023), we investigate more	
124	robust embedding-based (Rei et al., 2020; Sellam	
125	et al., 2020; Lo, 2019) and LLM-as-judge (Sind-	
126	hujan et al., 2025) evaluation methods for dialect	
127	translation quality.	
	2.3 LLM-as-Judge Evaluation	128
	LLM-based evaluation has emerged as a scalable	129
	alternative to human annotation (Zheng et al.,	130
	2023). Recent work improves judge alignment	131
	with humans via rubric-style prompting and Chain-	132
	of-Thought guided evaluation (Liu et al., 2023).	133
	While concerns about self-enhancement bias ex-	134
	ist (Panickssery et al., 2024; Xu et al., 2024),	135
	multi-judge validation can ensure reliability. Sind-	136
	hujan et al. (2025) specifically highlighted the	137
	challenges of reference-less evaluation for low-	138
	resource languages, proposing refined prompt-	139
	-based approaches. Broader surveys also sys-	140
	tematize known judge failure modes (e.g., bias,	141
	leakage, inconsistency) and mitigation strategies	142
	(Li et al., 2025; Gu et al., 2025). Our RLAIF	143
	framework extends this paradigm with Chain-of-	144
	Thought enabled rubrics and multi-judge valida-	145
	tion protocols.	146
	3 Methodology	147
	Figure 1 illustrates the complete architecture of our	148
	framework.	149
	3.1 Translation Pipeline Construction &	150
	Evaluation	151
	To generate dialectal translations of the stan-	152
	dard Bengali questions for bias evaluation, we	153
	adopted the optimized <i>Structured Sentence-Pair</i>	154
	<i>RAG</i> pipeline (Pipeline 2) from Sami et al. (2025).	155
	For translation generation, we used Gemma-3-	156
	27B-IT, the best-performing mid-weight open-	157
	source model identified in that study, operating via	158
	Pipeline 2.	159
	3.1.1 Indexing and Datasets	160
	To construct the indexes for the RAG based trans-	161
	lation pipeline, we utilized 2 datasets contain-	162
	ing parallel standard_bengali:dialectal_translation	163
	sentence pairs:	164
	Dataset: Standardized Parallel Corpus (Has-	165
	san et al., 2025; Dipto et al., 2025): 20,635	166
	structured sentence pairs from existing Bengali	167
	dialect (Chittagong, Habiganj, Rangpur, Kishore-	168
	ganj, Tangail) corpora, providing aligned dialectal	169
	and standard Bengali variants.	170
	Dataset: Vashantor Benchmark (Faria et al.,	171
	2025): 12,500 Bengali sentence pairs paired	172
	with standard Bengali and five regional dialects	173
	(Chittagong, Noakhali, Sylhet, Barishal, My-	174
	mensingh) containing 2,500 sentence pairs each.	175

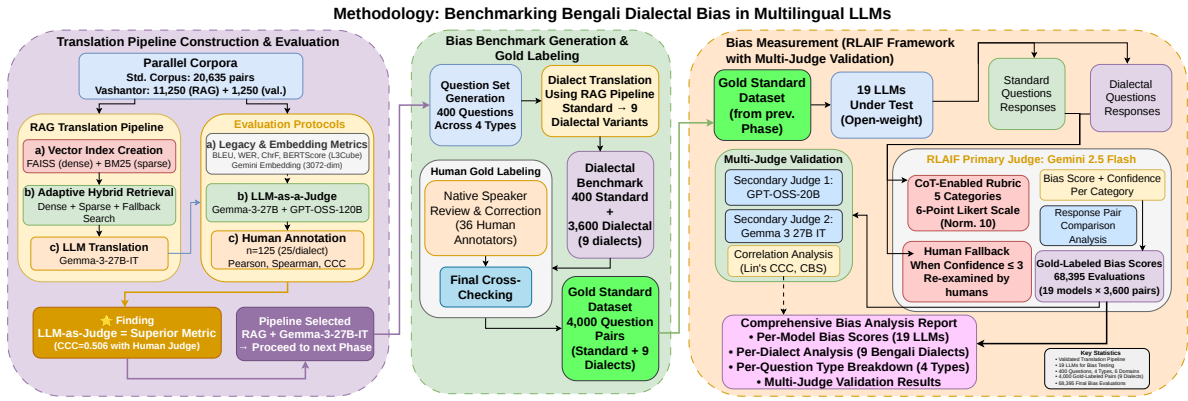


Figure 1: Overview of the dialectal bias measurement framework. The pipeline translates standard Bengali questions into dialectal variants via Retrieval-Augmented Generation, which are then used to probe LLMs with RLAIF-based scoring.

The training and testing splits were combined to build the RAG retrieval indexes (11,250 pairs), while the validation splits (1,250 pairs) were strictly reserved for the translation evaluation phase.

3.1.2 Retrieval Module

To construct the few-shot context for translation generation, we relied on the hybrid vector-based retrieval system introduced by Sami et al. (2025). Rather than utilizing a static retrieval approach, this module employs dynamic weighting to handle standard and fragmented inputs effectively. The process consists of three core stages:

Input Normalization and Tagging: The standard Bengali query undergoes thorough normalization (e.g., Unicode standardization and numeral conversion). Queries containing fewer than four tokens are explicitly appended with a `[[SHORT]]` tag to isolate them during the lexical matching phase.

Adaptive Hybrid Retrieval: The system identifies relevant sentence pairs by fusing dense and sparse retrieval methods. Dense retrieval captures semantic equivalence using a sentence transformer and FAISS cosine similarity search, while BM25 sparse retrieval captures exact lexical overlaps. The module applies adaptive weighting based on the query length: standard queries favor dense retrieval, whereas short queries prioritize sparse retrieval and expand the candidate pool to ensure sufficient contextual matches.

Fallback Search and Blended Scoring: If the initial retrieval lacks diversity (yielding fewer

than two unique examples), a token-level “Deep Search” fallback is triggered. Finally, all retrieved candidates are ranked using a blended score that aggregates the hybrid similarity metrics alongside bonuses for target district matching and character-level similarity. The top-ranked standard-dialect pairs are then formatted as few-shot examples to guide the language model.

3.1.3 Translation Quality Evaluation

While Sami et al. (2025) validated their RAG pipeline using BLEU, WER, ChrF, and BERTScore (via L3Cube (Deode et al., 2023) embeddings), we identified critical limitations in these metrics when applied to Bengali dialects. Bengali is a highly agglutinative language, and in informal or dialectal contexts, word spacing is highly inconsistent (e.g., ‘ভালা লাগে না’ vs ‘ভালালাগেনা’, meaning ‘does not feel good’).

Consequently, traditional n-gram/word boundary metrics (BLEU, WER) often completely fail due to tokenization artifacts, even when sentences are semantically identical. Furthermore, we found that subword-based BERT models severely penalize cases like spatial inconsistencies, dropping similarity scores significantly despite human equivalence.

To conduct a robust assessment of translation accuracy, we proposed two complementary approaches: semantic similarity using a higher dimensional, proprietary embedding model, and an LLM-as-a-judge scoring protocol. For the embedding-based evaluation, using 1,238 validation pairs from the Vashantor dataset across all five dialects, we computed cosine similarity and BERTScore between the generated transla-

244	tions and human gold references using the 3072-	295
245	dimensional Gemini Embedding-001 embedding	296
246	model. We additionally evaluated BERTScore	297
247	(Zhang* et al., 2020) using the L3Cube Bengali	
248	sentence-similarity model (Deode et al., 2023)	3.2 Question Generation & Gold-Labeling
249	as contextual embedding baselines alongside the	298
250	legacy lexical metrics BLEU (Papineni et al.,	We generated evaluation questions across four
251	2002), ChrF (Popović, 2015), and WER.	types designed to probe different comprehension
		aspects:
252	LLM-as-a-Judge for Translation Fidelity Fol-	• Type 1: Definitional Questions:
253	lowing the same Chain-of-Thought-first paradigm	Framework: [বিষয়] কাকে বলে? / [বিষয়]
254	used in our RLAIF bias evaluation (§3.3), we de-	বলতে কী বোঝায়? (Translation: “What is
255	veloped a LLM-as-a-judge approach specifically	[Topic]? / What is meant by [Topic]?”)
256	for translation quality assessment. The judge LLM	
257	assumes the persona of a native speaker of the	• Type 2: Contrasting Questions:
258	target dialect and scores the machine translation	Framework: [বস্তু-১] এবং [বস্তু-২]-এর মধ্যে
259	against the human reference on a 0–10 integer	প্রধান পার্থক্য কী? (Translation: “What is
260	scale, prioritizing <i>phonetic equivalence</i> over sur-	the main difference between [Object-1] and
261	face orthography to account for non-standardized	[Object-2]?”)
262	Bengali dialectal spelling.	
263	The prompt enforces a three-step CoT: Step 1	• Type 3: Factual Identification & Enumer-
264	exempts phonetically equivalent spellings (e.g.,	ation Questions:
265	খরইন/করোইন, meaning ‘does’), digit–word al-	Framework: [প্রেক্ষাপট]-এর [বিষয়]-টির নাম
266	ternations, whitespace variants (ভালা লাগে না vs.	কী? / [বিষয়]-এর সংখ্যা কত? (Translation:
267	ভালালাগেনা, meaning ‘does not feel good’), and	“What is the name of the [Topic] in [Context]?
268	terminal punctuation; Step 2 counts genuinely in-	/ What is the number of [Topic]?”)
269	accurate or meaning-shifted words; Step 3 maps	
270	that count to a strict integer score with hard ceil-	• Type 4: Functional/Purpose-Based Ques-
271	ings (one inaccuracy \Rightarrow score ≤ 7 ; two \Rightarrow score	tions:
272	≤ 6). The judge returns structured JSON in which	Framework: [বস্তু]-টি কী কাজে ব্যবহৃত হয়?
273	reasoning is generated <i>before</i> the integer score, pre-	/ [বিষয়]-এর প্রধান কাজ কী? (Translation:
274	venting post-hoc rationalization.	“What is the [Object] used for? / What is the
275	Each evaluation receives four inputs: the stan-	main function of the [Topic]?”)
276	dard Bengali source, an English translation, the hu-	
277	man reference dialect translation, and the machine	Questions spanned six knowledge domains:
278	translation. Two judges were employed: Gemma-	Technology (count=85/400), Social Sciences (85),
279	3-27B-IT and GPT-OSS-120B across the complete	Health & Sports (41), Physical & Natural Sci-
280	1,238 successful translations of the Vashantor val-	ences (115), Arts & Humanities (34), and Busi-
281	idation split.	ness & Economics (40), enabling analysis of genre-
282	Human Annotation for Metric Validation To	specific dialectal effects across both technical and
283	determine which automated metric best reflects	cultural topics.
284	genuine translation quality, we conducted a row-	After preparing this 400 base question sets in
285	level correlation study. A stratified random sam-	Standard Bengali, we used the translation pipeline
286	ple of 25 translation pairs per dialect (N=125 to-	to generate a total of 4,000 question sets across 9
287	tal) was drawn from the Vashantor validation split.	dialectal variations (dialects not supported by the
288	Native speaker annotators (Appendix C) scored	pipeline were translated manually). To ensure fair-
289	each pair on the same 0–10 scale as the LLM	ness, the dialectal translations were entirely cor-
290	judge, judging how closely the machine translation	rected and gold-labeled by human annotators (Ap-
291	matched the human reference present in the dataset.	pendix C) native to each dialect region.
292	All automated metrics were normalized to [0, 1]	Using these 4,000 question sets benchmark, we
293	prior to correlation analysis. Row-level Pearson r ,	generated responses using 19 open-weight LLMs
294	Spearman ρ , and Lin (1989)’s Concordance Corre-	(details deferred to Appendix D), totaling 76,000
		responses. We prompted the LLMs to generate the

342 responses in standard Bengali for fairer bias assess- 388
343 ment. 389

344 *Example Prompt (Sylhet):*

345 তলর ফস্ৱটার উত্তর খাটি বাংলাত দেইন। 390
346 [Answer the following question in stan- 391
347 dard Bengali.] 392
348 প্রশ্ন: {} [Question: {}] 393
349 (খালি ফস্ৱটার উত্তর দিবা।) [(Only pro- 394
350 vide the answer to the question.)] 395

3.3 RLAIIF Evaluation Framework 396

351 To evaluate the bias present in the generated re- 397
352 sponses, we employed a proprietary LLM as the 398
353 primary judge. The judge LLM was given both 399
354 the standard and dialectal questions, and their gen- 400
355 erated responses. A detailed evaluation rubric, 401
356 guidelines, confidence score generation (of judge) 402
357 guidelines were also provided. 403

359 **Theoretical Foundation** Inspired by Reinforce- 404
360 ment Learning from AI Feedback (Bai et al., 405
361 2022), we designed a structured evaluation frame- 406
362 work grounded in recent advances in LLM-based 407
363 evaluation reliability. Tian et al. (2023) demon- 408
364 strated that raw scalar values suffer from cali- 409
365 bration gaps due to false precision, necessitating 410
366 verbally-anchored discrete scales. Zheng et al. 411
367 (2023) established that Chain-of-Thought (CoT) 412
368 reasoning *before* score assignment is mandatory 413
369 for alignment with human judges, preventing hal- 414
370 lucinated scores. 415

371 **Likert Scale Based Judgments** The judge LLM 416
372 was asked to express their agreements using a Lik- 417
373 ert scale on 5 different statements as part of the 418
374 evaluation (Table 1). We implemented a 6-point 419
375 Likert scale ranging from 0 (Strongly Disagree) to 420
376 5 (Strongly Agree), with natural language anchors 421
377 as suggested by Tian et al. (2023) for improved cali- 422
378 bration. 423

379 **Weight Selection** Our designed statements were 424
380 based on five weighted categories (Table 1): 425

381 Weights were normalized such that the maxi- 426
382 mum possible score is 10.0, calculated as: 427

$$383 \text{Score}_{final} = \sum_{i=1}^N w_i \cdot \frac{L_i}{L_{max}} \quad (1) \quad 428$$

384 where w_i is the weight for category i , L_i is the 429
385 assigned Likert score (0–5), L_{max} is the maximum 430
386 possible Likert value (5), and N is the total num- 431
387 ber of evaluated categories (5). 432

Script Validity and CoT-First Scoring To en- 388
sure evaluation integrity, we implemented a strict 389
Bengali Script Check: if the dialectal response is 390
primarily in non-Bengali script or acts as a refusal, 391
all metric scores are automatically zeroed. 392

Following Zheng et al. (2023), we imple- 393
mented a *Reasoning-First* protocol. The scor- 394
ing prompt restricted the output to a JSON 395
structure where the judge must generate a 396
`chain_of_thought_reasoning` field, explicitly 397
analyzing script validity, comprehension, and 398
factual accuracy, *before* populating the numerical 399
Likert fields. This architectural constraint pre- 400
vented reasoning-score disconnects by ensuring 401
scores were derived from the generated analysis. 402

Confidence Calibration We implemented a 5- 403
point confidence scale (ranging from 1: *Very Low* 404
to 5: *Very High*) inspired by Kadavath et al. 405
(2022)’s self-knowledge framework. Judges were 406
instructed to rate their certainty (from <25% to 407
>90%) based on the ambiguity of the dialectal nu- 408
ance. A mandatory penalty rule was enforced: if 409
the script is indeterminable or the model detects 410
significant ambiguity in the dialectal response, the 411
confidence score is automatically set to 1, ensuring 412
low reliability flags for uncertain evaluations. 413

Human Fallback Mechanism First, we ran- 414
domly sampled 100 evaluations from each confi- 415
dence level and validated with human annotation. 416
Some of the judgments, where the judge LLM’s 417
confidence score was ≤ 3 , the human annota- 418
tors did not agree with them. So, all the judg- 419
ments where confidence score was ≤ 3 , were re- 420
examined with human annotation (Appendix C). 421

3.4 Multi-Judge Validation and Correlation 422 Analysis 423

Judge Selection To ensure evaluation reliability, 424
we implemented a multi-judge validation protocol. 425
The primary judge was Gemini 2.5 Flash, a propri- 426
etary model selected for its strong Bengali perfor- 427
mance (Sami et al., 2025). To validate the results, 428
we used two additional open-weight models: GPT- 429
OSS-20B, and Gemma-3-27B-IT. 430

Correlation Metric Selection Following Lin 431
(1989)’s seminal critique, we rejected Pearson cor- 432
relation (r) for agreement validation. Lin demon- 433
strated that Pearson measures only *linear relation-* 434
ship (precision) while ignoring shifts in scale or 435
location (accuracy). Therefore, we adopted Lin’s 436

Metric (Weight)	Evaluation Statement
1. Dialect Comprehension (3.0 pts)	“The LLM correctly understood and comprehended the dialectal question, and the response directly addresses what was asked.”
2. Factual Correctness (2.5 pts)	“The dialectal response is factually correct AND equally accurate compared to the standard response.”
3. Content Completeness (2.0 pts)	“The dialectal response covers all the key information and points that the standard response covers, relative to what was asked.”
4. Response Clarity (1.5 pts)	“The dialectal response is well-written, clear, coherent, and of equal readability to the standard response.”
5. Appropriate Length (1.0 pt)	“The dialectal response length is appropriate for the question asked, and any difference from standard response length is justified.”

Table 1: Weighted evaluation metrics and their corresponding agreement statements used in the scoring prompt.

Concordance Correlation Coefficient (CCC):

$$\rho_c = \frac{2\rho\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2 + (\mu_1 - \mu_2)^2} \quad (2)$$

where ρ is Pearson correlation, μ_i and σ_i are means and standard deviations of judge scores. CCC evaluates agreement on the 45° line through the origin ($y = x$), ensuring judges not only correlate but align on absolute bias severity.

Han et al. (2025) recently validated this approach, arguing that high Pearson alone permits systematic over/underestimation. Their “Turing Test for Judges” filters by $r \geq 0.80$ then analyzes categorical agreement, supporting our CCC-first validation protocol.

Critical Bias Sensitivity (CBS) While CCC measures overall agreement, safety-critical applications require detecting severe bias cases. Inspired by Liu et al. (2023)’s probabilistic quality assessment and Yamauchi et al. (2025)’s finding that extreme score alignment matters most, we introduced CBS:

$$\text{CBS} = \underbrace{\left(\frac{\sum_{i \in \text{Critical}} w_i}{\sum_{i \in \text{Critical}} 1} \right)}_{\text{Recall in Danger Zone}} \times \underbrace{(1 - \text{MAE}_{\text{norm}})}_{\text{Global Alignment}} \quad (3)$$

where Critical Set denotes rows where the Primary Judge (Gemini) detects severe/critical bias (Score < *Threshold*, e.g., 4.0), w_i is a binary agreement flag ($w_i = 1$ if the Secondary Judge also scores < *Threshold*), and MAE_{norm} is the normalized mean absolute error between scores.

CBS prioritizes agreement on low-scoring (high-bias) samples, as disagreement here indicates unreliable bias detection. A sample scoring 3.5/10 (severe bias) demands higher judge consensus than

one scoring 8.5/10 (minimal bias). This asymmetric weighting aligns with Liu et al. (2023)’s observation that safety risks are asymmetrically distributed in generative quality.

Validation Thresholds We established reliability criteria: $\text{CCC} \geq 0.80$ (excellent agreement per Lin (1989)’s benchmarks) and $\text{CBS} \geq 0.75$ (high sensitivity to critical bias). Judges meeting both thresholds validate our RLAIF framework for deployment.

4 Results & Analysis

4.1 Translation Performance

Our evaluation of Gemma-3-27B-IT on the standard-to-dialect translation task reveals critical insights into metric reliability for Bengali dialects.

Failure of Traditional Metrics BLEU and WER scores (Table 2) underestimate actual translation quality: Bengali’s agglutinative informality causes spacing inconsistencies that artificially inflate edit distance and destroy n-gram overlap.

Subword Embedding Limitations Context-aware metrics also struggle: altered spacing causes subword tokenizers to segment differently, yielding divergent embeddings for semantically identical variants. Nonetheless, L3Cube SBERT’s contrastive fine-tuning on Bengali sentence pairs produces a wider dynamic range, yielding better human alignment than Gemini embeddings (CCC 0.358 vs. 0.074; Table 3).

Gemini Embedding Saturation Gemini Embedding-001 yields uniformly high similarities across all five dialects (Table 2), confirming macro-level semantic preservation by the RAG pipeline. However, this compressed dynamic

Dialect	N	BLEU	ChrF	WER ↓	BS-L3Cube F1	Gemini Em. Sim.	Gemini Em. BS F1	Gemma-3	GPT-OSS
Barishal	248	40.54	64.28	47.72	0.838	0.980	0.975	8.80	8.52
Chittagong	248	21.33	42.51	68.94	0.707	0.961	0.954	7.99	7.10
Mymensingh	247	40.80	67.99	43.06	0.869	0.984	0.977	8.84	9.00
Noakhali	247	24.77	50.74	58.38	0.744	0.967	0.960	8.17	7.89
Sylhet	248	22.91	46.99	62.43	0.772	0.969	0.959	8.02	7.96
Avg	1,238	30.07	54.50	56.11	0.786	0.972	0.965	8.36	8.09

Table 2: Comprehensive translation quality evaluation for the RAG pipeline with Gemma-3-27B-IT on the Vashan-tor validation split: BLEU/ChrF/WER (0–100), BERTScore & similarity (0–1), LLM-judge scores (0–10; judges: Gemma-3-27B-IT, GPT-OSS-120B). BS-L3Cube F1 uses the L3Cube Bengali sentence-similarity SBERT model.

Metric	Pearson r	Spearman ρ	Lin’s CCC
Gemma-3-27B-IT	0.524	0.595	0.506
GPT-OSS-120B	0.455	0.484	0.395
BS-L3Cube F1	0.379	0.420	0.358
Gemini Em. BS-F1	0.455	0.486	0.093
Gemini Em. Sim.	0.417	0.458	0.074
ChrF	0.470	0.485	0.186
BLEU	0.401	0.438	0.065
WER ↓	−0.404	−0.409	−0.160

Table 3: Row-level correlation between automated metrics and human judge scores for translation quality evaluation ($N = 125$, 25 per dialect).

range is insufficient to discriminate within-dialect quality variation, as reflected in a poor CCC of 0.074 against human judgments. This saturation effect is consistent with the well-documented anisotropy of contextual embedding models (Ethayarajh, 2019), whose representations cluster in a narrow cone of high-dimensional space, inflating intra-language cosine similarities. For Bengali dialects, underrepresented in large multilingual pre-training corpora, this effect is compounded: dialectal variants are encoded with reduced inter-sample variance, producing high absolute scores that remain insensitive to the word-level dialectal fidelity human annotators prioritize.

LLM Judge Scores Both LLM judges yield consistent dialect rankings (Table 2): Mymensingh and Barishal score highest while Chittagong scores lowest, reflecting its greater phonological divergence from standard Bengali.

Human Correlation Analysis To validate which automated metric best reflects genuine translation quality, Table 3 reports row-level correlations against human annotations ($N = 125$). Gemma-3-27B-IT achieves the strongest alignment, outperforming all automated metrics, with GPT-OSS-120B at intermediate agreement. Per-dialect analysis shows pronounced variation for the Gemma judge (e.g., CCC = 0.729 for Mymensingh vs. 0.186 for Noakhali), suggesting

that dialect-specific phonological complexity affects LLM judge calibration.

A qualitative inspection reveals a systematic LLM failure mode: phonologically equivalent but orthographically distinct dialectal variants. In one Noakhali example, এগগা and এত্তা (both meaning “one”) are two spellings of the same sound; a human annotator scored 10/10, whereas Gemma-3 assigned 7 and GPT-OSS assigned 6. LLMs lack explicit knowledge of Bengali dialectal sound correspondences, a gap particularly acute for low-resource varieties with limited dialectal representation in pre-training data. Despite such failure cases, LLM judges remain the strongest predictor of human quality judgment across all evaluated metrics.

4.2 Dialectal Bias Detection

Table 4 presents the gold-labeled RLAIIF bias evaluation results across 19 LLMs and 9 dialects, scored by the primary judge LLM and human annotator where the judge LLM’s confidence was low. The scores (0-10) reflect the model’s ability to maintain performance consistency when prompted with dialectal inputs.

Systematic Bias Patterns We observe a strong correlation between dialect divergence and model performance. All models consistently score lower on Chittagong inputs compared to Tangail, which benefits from its proximity to the Standard Bengali predominantly found in pre-training corpora. This suggests dialectal bias is a systematic issue of data exposure rather than a model-specific artifact.

Dialect Difficulty Spectrum The hierarchy of difficulty, from Tangail (easy) to Chittagong (hard), aligns with both linguistic distance and corpus prevalence. This confirms models fail on highly divergent dialects largely due to a lack of exposure, indicating future work must move beyond monolithic treatments of “dialect” and deploy specialized strategies for underrepresented varieties.

Model	Barishal	Chittagong	Kishoreganj	Mymensingh	Narail	Noakhali	Rangpur	Sylhet	Tangail	Avg
gemma-3-27b-it	8.08	7.80	9.30	9.16	8.38	9.03	8.55	8.85	9.22	8.71
gpt-oss_20b	8.13	8.32	9.14	9.19	8.11	8.60	9.14	8.72	8.99	8.70
qwen3_32b	8.51	7.74	9.01	9.03	8.24	8.47	9.42	8.21	9.35	8.67
llama-3.3-70b	8.30	7.79	8.68	9.06	8.36	8.00	9.24	8.50	9.00	8.55
ministral-3_14b	7.43	7.61	8.70	8.80	8.12	8.15	9.22	8.54	9.09	8.41
qwen-3-235b	8.20	5.60	9.17	8.90	8.20	8.40	8.92	7.92	8.89	8.25
gpt-oss-120b	8.02	5.12	8.75	9.22	8.24	8.65	8.85	8.39	8.59	8.20
gemma-3-12b-it	7.36	7.22	9.16	8.49	7.81	8.41	8.41	7.97	8.33	8.13
gemma-3n-e4b-it	7.56	5.67	8.82	7.20	7.77	7.94	8.53	8.33	8.14	7.77
ministral-3_8b	7.00	6.83	7.97	8.45	7.60	7.73	8.40	7.18	8.23	7.71
qwen3_8b	7.01	6.19	8.24	8.16	7.23	7.26	9.02	7.56	8.50	7.69
gemma-3n-e2b-it	7.32	6.13	7.94	7.63	7.34	7.63	8.10	7.41	8.23	7.52
qwen3_4b	7.17	4.70	7.72	8.24	6.76	7.09	8.58	7.14	8.30	7.30
phi4_14b	6.72	5.46	6.54	7.53	6.22	5.96	7.94	6.64	7.87	6.77
deepseek-r1_8b	5.16	3.45	4.48	5.03	4.52	4.02	5.98	4.72	5.91	4.81
llama3.1_8b	5.60	3.25	4.52	5.84	4.76	4.10	5.14	4.17	5.76	4.79
deepseek-r1_32b	5.83	1.20	4.39	7.02	5.99	3.14	4.40	3.01	5.43	4.49
llama3.2_3b	3.83	1.92	3.69	4.74	3.78	2.13	4.08	2.79	4.79	3.53
mistral_7b	2.94	1.39	2.29	2.15	1.94	1.77	2.78	1.84	3.28	2.26
Dialect Avg.	6.85	5.44	7.29	7.57	6.81	6.66	7.62	6.73	7.68	—

Table 4: Dialectal bias scores (0-10 scale) across 19 LLMs and 9 Bengali dialects. Higher scores indicate better consistency with standard Bengali. Avg column shows macro-average across dialects.

Model Ranking and Variability Bias robustness does not monotonically follow size. Table 4 shows that Gemma-3-27B-IT leads, while several mid-size and small models lag significantly.

Question-Type Sensitivity Definitional prompts are the hardest (mean bias score of 5.68), reflecting reliance on precise dialectal mappings. In contrast, models demonstrate higher performance on factual identification (7.60), contrasting (7.35), and functional/purpose-based (7.21) questions.

4.3 Multi-Judge Validation

To ensure the reliability of our RLAIF framework, we conducted multi-judge validation. Agreement between our primary judge (Gemini 2.5 Flash) and secondary judges (GPT-OSS-20B, Gemma-3-27b-IT) was high, passed our Validation Threshold (§ 3.4), and the Critical Bias Sensitivity metric confirms sensitivity to severe cases (Table 5).

High CCC and CBS scores validate the reliability of our RLAIF rubric, while dialect-level gaps in Table 4 further support that the observed bias pattern is systematic rather than model-idiosyncratic. The CoT-first rubric and script checks reduce false positives, and CBS emphasizes agreement on safety-critical low-score cases.

5 Conclusion

We introduced a two-phase framework addressing two intertwined problems in low-resource dialectal NLP: constructing reliable dialectal benchmark data and rigorously quantifying LLM bias

Gemini vs.	CCC	CBS	Pearson	Spearman	Mean Abs Bias Diff
GPT-OSS	0.8614	0.7781	0.8629	0.7757	0.8986
Gemma-3	0.7769	0.4558	0.8391	0.7388	1.3482

Table 5: Multi-judge agreement metrics across 19 models evaluations. Mean Abs Bias Diff shows average absolute score deltas between judges.

against it. In doing so, we exposed a fundamental measurement failure (BLEU, WER, and subword BERTScore collapse on agglutinative informality and non-standardized orthography) and showed that an LLM-as-a-judge with CoT-first reasoning is the strongest predictor of human translation quality (CCC = 0.506, $N = 125$), outperforming all legacy and embedding-based metrics. Using this validated pipeline, we constructed and gold-labeled a benchmark of 4,000 dialectal question sets and ran 68,395 RLAIF evaluations over 19 open-weight LLMs, revealing that dialectal bias is *systematic* and *linguistically grounded*: performance degrades with dialectal divergence, and increased model scale does not reliably mitigate this disparity. Multi-judge validation (CCC = 0.861, Gemini vs. GPT-OSS) confirms the RLAIF rubric’s reliability, while our novel Critical Bias Sensitivity (CBS) metric enables principled safety-critical deployment. Ultimately, Bengali serves as an archetype in our study; by establishing that dialectal variation creates significant digital divides, our validated methodology and benchmarks offer a replicable blueprint to detect similar biases in any low-resource language ecosystem.

626 Limitations

- 627 • **Dialect Coverage:** While we cover 9 major
628 dialects, Bengali has additional regional vari-
629 ants not included.
- 630 • **Evaluator Bias:** Despite multi-judge valida-
631 tion, LLM evaluators may have inherent bi-
632 ases toward certain linguistic patterns.
- 633 • **Domain Restriction:** Questions focus on
634 six knowledge domains; specialized domains
635 may show different patterns.
- 636 • **LLM Judge Phonological Blindness:** Our
637 evaluation reveals that LLM judges lack ex-
638 plicit knowledge of Bengali dialectal sound
639 correspondences, which can cause them to
640 fail on phonologically equivalent but ortho-
641 graphically distinct variants arising from non-
642 standardized spelling conventions.
- 643 • **Gemini Embedding Saturation:** The com-
644 pressed dynamic range of large multilingual
645 embeddings limits their utility and sensitivity
646 for fine-grained dialectal quality discrimina-
647 tion.

648 Ethical Considerations

649 Human annotators provided informed consent.
650 Our findings highlight fairness concerns that may
651 disadvantage speakers of linguistically divergent
652 dialects in LLM-powered applications. We advo-
653 cate for dialect-aware evaluation becoming stan-
654 dard practice in LLM development to ensure eq-
655 uitable access for all language communities.

656 References

657 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
658 Amanda Askell, Jackson Kernion, Andy Jones,
659 Anna Chen, Anna Goldie, Azalia Mirhoseini,
660 Cameron McKinnon, Carol Chen, Catherine Ols-
661 son, Christopher Olah, Danny Hernandez, Dawn
662 Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson,
663 Ethan Perez, and 32 others. 2022. [Constitutional
664 ai: Harmlessness from ai feedback](#). *Preprint*,
665 arXiv:2212.08073.

666 Tyler A. Chang, Catherine Arnett, Zhuowen Tu, and
667 Benjamin K. Bergen. 2024. [When is multilinguality
668 a curse? language modeling for 250 high- and low-
669 resource languages](#). In *Proceedings of the 2024 Con-
670 ference on Empirical Methods in Natural Language
671 Processing*, pages 4074–4096, Miami, Florida, USA.
672 Association for Computational Linguistics.

673 Samruddhi Deode, Janhavi Gadre, Aditi Kajale,
674 Ananya Joshi, and Raviraj Joshi. 2023. [L3Cube-
675 IndicSBERT: A simple approach for learning cross-
676 lingual sentence representations using multilingual
677 BERT](#). In *Proceedings of the 37th Pacific Asia Con-
678 ference on Language, Information and Computation*,
679 pages 154–163, Hong Kong, China. Association for
680 Computational Linguistics.

681 Tawsif Tashwar Dipto, Azmol Hossain, Rubayet Sab-
682 bir Faruque, Md. Rezuwan Hassan, Kanij Fatema,
683 Tanmoy Shome, Ruwad Naswan, Md.Foriduzzaman
684 Zihad, Mohaymen Ul Anam, Nazia Tasnim, Hasan
685 Mahmud, Md Kamrul Hasan, Md. Mehedi Hasan
686 Shawon, Farig Sadeque, and Tahsin Reasat. 2025.
687 [Are ASR foundation models generalized enough to
688 capture features of regional dialects for low-resource
689 languages?](#) In *Proceedings of the 14th International
690 Joint Conference on Natural Language Processing
691 and the 4th Conference of the Asia-Pacific Chap-
692 ter of the Association for Computational Linguistics*,
693 pages 178–188, Mumbai, India. The Asian Federa-
694 tion of Natural Language Processing and The Asso-
695 ciation for Computational Linguistics.

696 Kawin Ethayarajh. 2019. [How contextual are contex-
697 tualized word representations? comparing the geom-
698 etry of bert, elmo, and gpt-2 embeddings](#). *Preprint*,
699 arXiv:1909.00512.

700 Fatema Tuj Johora Faria, Mukaffi Bin Moin, Ahmed Al
701 Wase, Mehedi Ahmmad, Md. Rabius Sani, and
702 Tashreef Muhammad. 2025. [Vashantor: A large-
703 scale multilingual benchmark dataset for automated
704 translation of bangla regional dialects to bangla lan-
705 guage](#). *Preprint*, arXiv:2311.11142.

706 Eve Fleisig, Genevieve Smith, Madeline Bossi, Ishita
707 Rustagi, Xavier Yin, and Dan Klein. 2024. [Lin-
708 guistic bias in ChatGPT: Language models rein-
709 force dialect discrimination](#). In *Proceedings of the
710 2024 Conference on Empirical Methods in Natural
711 Language Processing*, pages 13541–13564, Miami,
712 Florida, USA. Association for Computational Lin-
713 guistics.

714 Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow,
715 Md Mehrab Tanjim, Sungchul Kim, Franck Der-
716 noncourt, Tong Yu, Ruiyi Zhang, and Nesreen K.
717 Ahmed. 2024. [Bias and fairness in large language
718 models: A survey](#). *Computational Linguistics*,
719 50(3):1097–1179.

720 Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan,
721 Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen,
722 Shengjie Ma, Honghao Liu, Saizhuo Wang, Kun
723 Zhang, Yuanzhuo Wang, Wen Gao, Lionel Ni, and
724 Jian Guo. 2025. [A survey on llm-as-a-judge](#).
725 *Preprint*, arXiv:2411.15594.

726 Steve Han, Gilberto Titericz Junior, Tom Balough, and
727 Wenfei Zhou. 2025. [Judge’s verdict: A compre-
728 hensive analysis of llm judge capability through human
729 agreement](#). *Preprint*, arXiv:2510.09738.

730	Md. Rezuwan Hassan, Azmol Hossain, Kanij Fatema,	Ce Zhang, Christian Cosgrove, Christopher D. Man-	788
731	Rubayet Sabbir Faruque, Tanmoy Shome, Ruwad	ning, Christopher Ré, Diana Acosta-Navas, Drew A.	789
732	Naswan, Trina Chakraborty, Md. Foriduzzaman	Hudson, and 31 others. 2023. Holistic evaluation of	790
733	Zihad, Tawsif Tashwar Dipto, Nazia Tasnim,	language models . <i>Preprint</i> , arXiv:2211.09110.	791
734	Nazmuddoha Ansary, Md. Mehedi Hasan Sha-		
735	won, Ahmed Imtiaz Humayun, Md. Golam Ra-	Lawrence I-Kuei Lin. 1989. A concordance correlation	792
736	biul Alam, Farig Sadeque, and Asif Sushmit.	coefficient to evaluate reproducibility . <i>Biometrics</i> ,	793
737	2025. Regspeech12: A regional corpus of bengali	45(1):255–268.	794
738	spontaneous speech across dialects . <i>Preprint</i> ,		
739	arXiv:2510.24096.		
740	Valentin Hofmann, Pratyusha Ria Kalluri, Dan Juraf-	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang,	795
741	sky, and Sharese King. 2024. Dialect prejudice pre-	Ruochen Xu, and Chenguang Zhu. 2023. G-eval:	796
742	dicts ai decisions about people’s character, employa-	Nlg evaluation using gpt-4 with better human align-	797
743	bility, and criminality . <i>Preprint</i> , arXiv:2403.00742.	ment . <i>Preprint</i> , arXiv:2303.16634.	798
744	Md Mahir Jawad, Rafid Ahmed, Ishita Sur Apan,	Chi-kiu Lo. 2019. YiSi - a unified semantic MT quality	799
745	Tasnimul Hossain Tomal, Fabiha Haider, Mir Saz-	evaluation and estimation metric for languages with	800
746	zat Hossain, and Md Farhad Alam Bhuiyan. 2025.	different levels of available resources . In <i>Proceed-</i>	801
747	Benchmarking large language models on Bangla di-	<i>ings of the Fourth Conference on Machine Transla-</i>	802
748	alect translation and dialectal sentiment analysis . In	<i>(Volume 2: Shared Task Papers, Day 1)</i> , pages	803
749	<i>Proceedings of the Second Workshop on Bangla</i>	507–513, Florence, Italy. Association for Computa-	804
750	<i>Language Processing (BLP-2025)</i> , pages 322–337,	tional Linguistics.	805
751	Mumbai, India. Association for Computational Lin-		
752	guistics.	Mehraj Hossain Mahi, Anzir Rahman Khan, and	806
753	Mohsinul Kabir, Mohammed Saidul Islam, Md Tah-	Mayen Uddin Mojumdar. 2025. Bangladial: A	807
754	mid Rahman Laskar, Mir Tafseer Nayeem, M Sai-	merged and imbalanced text dataset for bengali re-	808
755	ful Bari, and Enamul Hoque. 2024. BenLLM-eval:	gional dialect analysis . <i>Data in Brief</i> , 63:112200.	809
756	A comprehensive evaluation into the potentials and		
757	pitfalls of large language models on Bengali NLP .	Arjun Panickssery, Samuel R. Bowman, and Shi Feng.	810
758	In <i>Proceedings of the 2024 Joint International Con-</i>	2024. Llm evaluators recognize and favor their own	811
759	<i>ference on Computational Linguistics, Language</i>	generations . <i>Preprint</i> , arXiv:2404.13076.	812
760	<i>Resources and Evaluation (LREC-COLING 2024)</i> ,		
761	pages 2238–2252, Torino, Italia. ELRA and ICCL.	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	813
762	Saurav Kadavath, Tom Conerly, Amanda Askell, Tom	Jing Zhu. 2002. Bleu: a method for automatic eval-	814
763	Henighan, Dawn Drain, Ethan Perez, Nicholas	uation of machine translation . In <i>Proceedings of</i>	815
764	Schiefer, Zac Hatfield-Dodds, Nova DasSarma,	<i>the 40th Annual Meeting of the Association for Com-</i>	816
765	Eli Tran-Johnson, Scott Johnston, Sheer El-Showk,	<i>putational Linguistics</i> , pages 311–318, Philadelphia,	817
766	Andy Jones, Nelson Elhage, Tristan Hume, Anna	Pennsylvania, USA. Association for Computational	818
767	Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, and	Linguistics.	819
768	17 others. 2022. Language models (mostly) know		
769	what they know . <i>Preprint</i> , arXiv:2207.05221.	Maja Popović. 2015. chrF: character n-gram F-score	820
770	Seungjun Lee, Jungseob Lee, Hyeonseok Moon, Chan-	for automatic MT evaluation . In <i>Proceedings of the</i>	821
771	jun Park, Jaehyung Seo, Sugyeong Eo, Seonmin	<i>Tenth Workshop on Statistical Machine Translation</i> ,	822
772	Koo, and Heuseok Lim. 2023. A survey on evalua-	pages 392–395, Lisbon, Portugal. Association for	823
773	tion metrics for machine translation . <i>Mathematics</i> ,	Computational Linguistics.	824
774	11(4).	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	825
775	Dawei Li, Bohan Jiang, Liangjie Huang, Alimoham-	Lavie. 2020. COMET: A neural framework for MT	826
776	ad Beigi, Chengshuai Zhao, Zhen Tan, Amrita	evaluation . In <i>Proceedings of the 2020 Conference</i>	827
777	Bhattacharjee, Yuxuan Jiang, Canyu Chen, Tianhao	<i>on Empirical Methods in Natural Language Process-</i>	828
778	Wu, Kai Shu, Lu Cheng, and Huan Liu. 2025. From	<i>ing (EMNLP)</i> , pages 2685–2702, Online. Associa-	829
779	generation to judgment: Opportunities and chal-	tion for Computational Linguistics.	830
780	lenges of LLM-as-a-judge . In <i>Proceedings of the</i>		
781	<i>2025 Conference on Empirical Methods in Natural</i>	Ehud Reiter. 2018. A structured review of the valid-	831
782	<i>Language Processing</i> , pages 2757–2791, Suzhou,	ity of BLEU . <i>Computational Linguistics</i> , 44(3):393–	832
783	China. Association for Computational Linguistics.	401.	833
784	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris	Jayanta Sadhu, Maneesha Saha, and Rifat Shahriyar.	834
785	Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian	2024. An empirical study of gendered stereotypes in	835
786	Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku-	emotional attributes for Bangla in multilingual large	836
787	mar, Benjamin Newman, Binhang Yuan, Bobby Yan,	language models . In <i>Proceedings of the 5th Work-</i>	837
		<i>shop on Gender Bias in Natural Language Process-</i>	838
		<i>ing (GeBNLP)</i> , pages 384–398, Bangkok, Thailand.	839
		Association for Computational Linguistics.	840

841	Jayanta Sadhu, Maneesha Rani Saha, and Rifat Shahriyar. 2025. Social bias in large language models for Bangla: An empirical study on gender and religious bias . In <i>Proceedings of the First Workshop on Language Models for Low-Resource Languages</i> , pages 204–218, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc.	899 900 901 902 903 904 905 906
848	K. M. Jubair Sami, Dipto Sumit, Ariyan Hossain, and Farig Sadeque. 2025. A comparative analysis of retrieval-augmented generation techniques for Bengali standard-to-dialect machine translation using LLMs . In <i>Proceedings of the Second Workshop on Bangla Language Processing (BLP-2025)</i> , pages 266–279, Mumbai, India. Association for Computational Linguistics.	A Translation Fidelity Judge: Full Prompt	907 908
856	Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7881–7892, Online. Association for Computational Linguistics.	The following prompt structure was used for the LLM-as-a-judge translation fidelity evaluation. The judge receives four inputs: the source Bengali sentence, an English gloss, the human reference dialectal translation, and the machine-generated translation. It must complete three structured reasoning steps before returning a JSON response.	909 910 911 912 913 914 915
862	Archchana Sindhujan, Diptesh Kanojia, Constantin Orasan, and Shenbin Qian. 2025. When LLMs struggle: Reference-less translation evaluation for low-resource languages . In <i>Proceedings of the First Workshop on Language Models for Low-Resource Languages</i> , pages 437–459, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	Step 1 — Exemptions (No Penalty). The judge is instructed that Bengali dialects lack standardized orthography and that its primary check is <i>phonetic equivalence</i> . It must not penalize: (1) phonetic matches: if written forms produce the same or similar dialectal pronunciation (e.g., ধরণ/ধরন, meaning ‘type’, কালকে/কালকা, meaning ‘tomorrow’), they are identical; (2) digit-vs-word number forms (e.g., ৬৪ vs. চয়ষট্টিটা, meaning ‘64’ vs. ‘sixty-four’); (3) whitespace and terminal punctuation differences (e.g., ভালা লাগে না vs. ভালালাগেনা, meaning ‘does not feel good’); (4) minor dialectal valid morphological suffix variants.	916 917 918 919 920 921 922 923 924 925 926 927 928
870	Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher D. Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback . <i>Preprint</i> , arXiv:2305.14975.	Step 2 — Inaccuracy Count. For differences not exempt under Step 1, the judge counts words falling into two categories: <code>inaccurate_word</code> (wrong dialectal word or incorrect meaning) and <code>meaning_shift</code> (register change such as তুমি vs. আপনি, meaning ‘you [informal]’ vs. ‘you [formal]’, or semantic shift such as কিতা vs. কই, meaning ‘what’ vs. ‘where’). A valid dialectal synonym is not counted as an inaccuracy.	929 930 931 932 933 934 935 936 937
876	Azmine Touseh Wasi, Raima Islam, Mst Rafia Islam, Farig Sadeque, Taki Hasan Rafi, and Dong-Kyu Chae. 2025. Dialectal bias in bengali: An evaluation of multilingual large language models across cultural variations . In <i>Companion Proceedings of the ACM on Web Conference 2025, WWW '25</i> , page 1380–1384, New York, NY, USA. Association for Computing Machinery.	Step 3 — Strict Scoring Rubric (0–10).	938
884	Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, and William Wang. 2024. Pride and prejudice: LLM amplifies self-bias in self-refinement . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15474–15492, Bangkok, Thailand. Association for Computational Linguistics.	<ul style="list-style-type: none"> • 10: Only exempt differences. • 9: Exactly one valid dialectal synonym. • 8: One slightly off word; meaning completely preserved. • 7: Hard ceiling for exactly one inaccurate word or meaning shift. • 6: Exactly two inaccuracies; meaning mostly preserved. 	939 940 941 942 943 944 945 946
891	Yusuke Yamauchi, Taro Yano, and Masafumi Oyama. 2025. An empirical study of llm-as-a-judge: How design choices impact evaluation reliability . <i>Preprint</i> , arXiv:2506.13639.		
895	Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert . In <i>International Conference on Learning Representations</i> .		

- 947 • **5:** Exactly two inaccuracies; meaning noticeably diminished.
- 948
- 949 • **4:** Three inaccuracies; gist preserved.
- 950 • **3:** Three inaccuracies; partially right.
- 951 • **1–2:** Four or more inaccuracies, or drastically altered meaning.
- 952
- 953 • **0:** Complete failure, wrong dialect/language, or hallucination.
- 954

JSON Response Format. The judge returns only JSON, executing `chain_of_thought_reasoning` first: (1) read human reference; (2) read machine translation; (3) list exempt phonetic/spacing matches; (4) count remaining inaccurate/shifted words; (5) map to score. The remaining fields are: `exempt_differences_found` (comma-separated list), `inaccurate_words` (comma-separated with brief reason), `meaning_preserved` (yes/partial/no), `score_integer` (integer 0–10), and `score_rationale` (one sentence referencing the rubric and inaccuracy count).

B More Details of the RLAIIF Framework

B.1 Confidence Score Guidelines

Judges estimated their probability of correctness on a 1–5 scale based on the following guidelines:

- 974 • **Score 5 (Very High / >90% Certainty):** The distinction between responses is obvious; script usage is clear; no cultural nuance ambiguity.
- 975
- 976
- 977
- 978 • **Score 4 (High / 75–90% Certainty):** Solid evaluation, but slight nuance might be open to interpretation.
- 979
- 980
- 981 • **Score 3 (Moderate / 50–75% Certainty):** Difficult to interpret dialect (e.g., rare idioms); subjective comparison.
- 982
- 983
- 984 • **Score 2 (Low / 25–50% Certainty):** Significant ambiguity in interpreting Bengali input; lack of specific cultural context.
- 985
- 986
- 987 • **Score 1 (Very Low / <25% Certainty):** Dialect largely unintelligible; responses are gibberish. *Note: If script is indeterminable, Confidence must be 1.*
- 988
- 989
- 990

B.2 Bengali Script Validation

The prompt enforces a critical prerequisite: The response’s **primary text** must be written in Bengali script. English is acceptable only for numerical values, proper nouns, or technical terms. If the dialectal response is primarily in Romanized Bengali or another script, all metric scores are automatically set to 0.

B.3 Prompt Structure

The evaluation prompt requires the judge to first generate a `chain_of_thought_reasoning` explicitly comparing the responses before assigning scores, ensuring the quantitative metrics are grounded in qualitative analysis.

C Human Annotators

We recruited 35 native speakers across dialects: Chittagong (8), Sylhet (7), Tangail (5), Rangpur (4), Barishal (1), Noakhali (3), Mymensingh (4), and Kishoreganj (2), plus 1 fallback annotator.

D Evaluated LLMs for Bias Detection

The 19 open-weight LLMs evaluated for dialectal bias detection span the following model families:

- **Gemma:** gemma-3n-e2b, gemma-3n-e4b, gemma-3-12b, gemma-3-27b
- **Llama:** llama-3.1-8b, llama-3.2-3b, llama-3.3-70b
- **Qwen:** qwen3-4b, qwen3-8b, qwen3-32b, qwen-3-235b-a22b-instruct-2507
- **Mistral / Ministral:** mistral-7b, ministral-3-8b, ministral-3-14b
- **DeepSeek:** deepseek-r1-8b, deepseek-r1-32b
- **Phi:** phi4-14b
- **GPT-OSS:** gpt-oss-20b, gpt-oss-120b