

# BANGLASOCIALBENCH: A Benchmark for Evaluating Sociopragmatic and Cultural Alignment of LLMs in Bangladeshi Social Interaction

Tanvir Ahmed Sijan<sup>1,†</sup> S. M Golam Rifat<sup>2</sup> Pankaj Chowdhury Partha<sup>3</sup>  
Md. Tanjeed Islam<sup>1</sup> Md. Musfique Anwar<sup>1</sup>

<sup>1</sup>Jahangirnagar University, Dhaka, Bangladesh,

<sup>2</sup>Rajshahi University of Engineering & Technology, Rajshahi, Bangladesh

<sup>3</sup>Bangladesh University of Engineering and Technology, Dhaka, Bangladesh,

{sijantanv, golamrifat, parthach25, tanjeedislamr}@gmail.com

manwar@juniv.edu <sup>†</sup>Corresponding author

## Abstract

Large Language Models have demonstrated strong multilingual fluency, yet fluency alone does not guarantee socially appropriate language use. In high-context languages, communicative competence requires sensitivity to social hierarchy, relational roles, and interactional norms that are encoded directly in everyday language. Bangla exemplifies this challenge through its three-tiered pronominal system, kinship-based addressing, and culturally embedded social customs. We introduce BANGLASOCIALBENCH, the first benchmark designed to evaluate sociopragmatic competence in Bangla through context-dependent language use rather than factual recall. The benchmark spans three domains: Bangla Address Terms, Kinship Reasoning, and Social Customs, comprising 1,719 culturally grounded instances written and verified by native Bangla speakers. We evaluate twelve contemporary LLMs in a zero-shot setting and observe systematic patterns of cultural misalignment. Models frequently default to overly formal address forms, fail to recognize multiple socially acceptable address pronouns, and conflate kinship terminology across religious contexts. Our findings show that sociopragmatic failures are often structured and non-random; for example, inappropriate addressing choices concentrate heavily in downward-hierarchy (Elder→Younger) and informal contexts. This reveals persistent limitations in how current LLMs infer and apply culturally appropriate language use in realistic Bangladeshi social interactions.

## 1 Introduction

Large Language Models (LLMs) have made rapid progress in multilingual generation, enabling fluent text production across a wide range of languages. However, fluency alone does not guarantee appropriate language use in social interaction, particularly in culturally grounded communi-

**Bn:** প্রেক্ষাপট: একজন বৃদ্ধ যাত্রী মেট্রো স্টেশনে দাঁড়ানো কিশোর ছেলেকে বিনয়ের সাথে লাইনে দাঁড়াতে বলছেন।  
এই প্রেক্ষাপটে, বাংলাদেশী সংস্কৃতিতে বৃদ্ধ লোকটি কিশোরকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবেন?  
এক কথায় উত্তর দিন — আপনি/তুমি/তুই?

**En:** Situation: An elderly passenger at a metro station is politely asking a teenage boy to stand in line.  
In this context, how would the elderly person conventionally address the teenage boy in Bangladeshi Culture?  
Answer in one word — apni/tumi/tui?

Setting Participants Ends Key Instruments Norms

Figure 1: Prompt design grounded in Hymes' SPEAKING model (Hymes, 1962). Each prompt operationalizes sociolinguistic context through explicit cues for setting, participants, gender, interactional goal, and social norms, allowing controlled evaluation of culturally appropriate Bangla Address Terms.

cation. In many languages, meaning is shaped not only by what is said, but by how it is said, to whom, and in what context. When models fail to account for these interactional constraints, they may produce outputs that are linguistically correct yet socially inappropriate. This distinction aligns with the notion of communicative competence, which emphasizes contextual appropriateness alongside grammatical correctness (Hymes, 1964).

This challenge is especially pronounced in high-context languages, where social relationships are encoded directly into everyday language use. Bangla, an Indo-Aryan language spoken by over 300 million people worldwide and used as the primary language by the vast majority of the population in Bangladesh (Sultana et al., 2023), exemplifies this phenomenon.<sup>1</sup> Rather than relying on

<sup>1</sup>We use **Bangla** for the language, **Bengali** for the ethnolinguistic group, and **Bangladeshi** for the civic nationality.

optional politeness markers, Bangla grammaticalizes social hierarchy through a three-tiered second-person pronominal system (apni, tumi, tui) and extensive kinship-based nominal addressing (Das, 1968a). Correct usage depends on factors such as power and formality, age difference, intimacy, emotion and inappropriate choices can sound rude or socially offensive despite correct propositional content (Uddin, 2019; Snigdha, 2022).

Despite the centrality of these phenomena to daily communication, existing evaluations only partially test whether models can reason over sociopragmatic distinctions in realistic interactional settings. Many existing multilingual evaluations rely on translated prompts, declarative questions, or factual recall tasks, which rarely require models to infer social roles or adapt language choices dynamically (Naous et al., 2024). As a result, models may appear culturally competent while still failing in direct social interaction.

To address this gap, we introduce **BANGLASOCIALBENCH**, a benchmark for evaluating sociopragmatic competence in Bangladeshi social contexts. The benchmark operationalizes language use as situated social behavior by embedding implicit sociopragmatic cues within realistic Bangla scenarios across three domains: Bangla Address Terms, Kinship Reasoning, and Social Customs. Using this framework, we evaluate twelve contemporary LLMs and observe consistent patterns of culturally inappropriate behavior. Across models, errors frequently arise from overuse of formal address forms, limited recognition of socially acceptable variation, and confusion in culturally grounded kinship usage, indicating systematic difficulties in resolving sociopragmatic ambiguity from context alone. In summary, our contributions are as follows:

- **BANGLASOCIALBENCH:** We introduce a benchmark that evaluates Bangla sociopragmatic competence through *situated social behavior* rather than static knowledge. It tests a model’s ability to infer hierarchy and intent across 1,719 human-written interactional scenarios. We make the benchmark publicly available to the research community.<sup>2</sup>
- **We demonstrate that sociopragmatic failure in LLMs is often systematic rather**

<sup>2</sup><https://huggingface.co/datasets/sijantanvir/BanglaSocialBench>

**than random** and conditioned on variables like age and setting. We find that culturally inappropriate address pronoun choices concentrate in downward-hierarchy (Elder→Younger) and informal contexts, exposing structural weaknesses in reasoning about social power.

- **We uncover characteristic patterns of cultural misalignment in twelve contemporary LLMs**, including pervasive over-politeness, probabilistic rigidity in address selection, and cross-religious kinship term conflation. These results suggest that models often rely on safe, default heuristics rather than socially grounded reasoning.

## 2 Related Work

### 2.1 Theoretical Frameworks and Factual Benchmarks

Recent work on cultural evaluation of LLMs distinguishes between culturally thin and culturally thick (Adilazuarda et al., 2024) conceptions of culture. Thin evaluations typically test recall of culture-related facts or alignment with aggregate statistics, offering broad but shallow signals of cultural awareness and often relying on outsider perspectives rather than lived experience. In contrast, culturally thick approaches view culture as an enacted and dynamic process, realized through language use in socially grounded interactional contexts (Zhou et al., 2025).

A large class of benchmarks operationalizes culture as region-specific factual knowledge, testing models on history, geography, literature, and linguistic facts using curated multiple-choice questions. Examples include SaudiCulture (Ayash et al., 2025), SANSKRITI (Maji et al., 2025), CLiCK (Kim et al., 2024), and IndoCulture (Koto et al., 2024). In the Bangla context, BnMMLU (Joy, 2025) evaluates massive multitask understanding using academic materials, while BLUCK (Kabir et al., 2025) focuses on linguistic competence and cultural knowledge. While effective at measuring recall, they largely neglect the dynamic, interactional competence required for social acceptability, leaving a gap in evaluating how models navigate social deixis and pragmatic ambiguity.

## 2.2 Value Alignment and Normative Reasoning

Another line of research evaluates cultural alignment through moral values and norms, comparing model responses against human population distributions derived from surveys or cultural indices (Zhao et al., 2024; AlKhamissi et al., 2024; Masoud et al., 2025). BengaliMoralBench (Ridoy et al., 2025) extends this approach to Bangla by auditing moral reasoning across culturally relevant subtopics. However, these approaches focus on abstract value alignment and largely overlook other semantic proxies of culture (Adilazuarda et al., 2024), particularly the realization of norms through interactional language use.

## 2.3 Sociopragmatic and Linguistic Evaluation

Recent work has moved toward evaluating sociopragmatic competence in language models, focusing on context-sensitive phenomena such as politeness, honorifics, and social etiquette. NormAd (Rao et al., 2025) evaluates whether models adapt outputs to sociocultural cues across scenarios, while PUB (Sravanthi et al., 2024) assesses pragmatic understanding across diverse contexts. The UNGGAH-UNGGUH benchmark (Farhan-syah et al., 2025) evaluates register-sensitive honorific usage in Javanese, primarily through sentence-level transformations. In South Asian languages, prior work on Bangla and Hindi honorifics restricts analysis to third-person forms, explicitly excluding second-person pronouns due to their sociolinguistic complexity (Mukherjee et al., 2023).

Several recent benchmarks examine sociopragmatic phenomena in high-context languages, where linguistic forms encode social hierarchy and interpersonal relations. For instance, *KinshipQA* (Sun and Kazakov, 2026) evaluates multi-hop reasoning over culturally grounded kinship structures derived from anthropological systems. *TAAROFBENCH* (Sadr et al., 2025) evaluates Persian *taarof*, a system of ritual politeness involving indirect offers and refusals.

Related work also explores pragmatic phenomena in East Asian languages. *The Mask of Civility* (Zhang et al., 2026) evaluates Chinese mock politeness and highlights challenges in detecting mismatches between polite expressions and impolite contexts. Likewise, the *KoBALT* dataset (Shin

et al., 2025) evaluates Korean across syntax, semantics, and morphology, revealing persistent difficulties with language-specific traits such as irregular verb conjugations and implicit interpersonal relations.

## 2.4 Addressing the Gap

Current multilingual evaluations often rely on English-centric or translated setups, which can obscure sociopragmatic distinctions encoded directly in native language use (Naous et al., 2024). As discussed, while recent efforts have introduced native Bangla benchmarks, they primarily target orthogonal capabilities. Datasets such as BnMMLU (Joy, 2025) and BLUCK (Kabir et al., 2025) focus heavily on massive multitask factual recall and static linguistic knowledge. Similarly, BengaliMoralBench (Ridoy et al., 2025) audits abstract moral reasoning, and prior work on Bangla honorifics restricts analysis strictly to third-person forms (Mukherjee et al., 2023).

While effective in their domains, existing benchmarks neglect the dynamic interactional competence required for everyday communication. In contrast, communicative competence in Bangla relies heavily on sensitivity to social hierarchy, relational roles, and interactional norms, requiring models to jointly reason over hierarchical address forms, complex kinship structures, and culturally embedded social customs. To address this, *BANGLASOCIALBENCH* introduces a culturally thick framework evaluating sociopragmatic competence in realistic Bangladeshi settings.

## 3 Cultural Context and Overview of *BANGLASOCIALBENCH*

*BANGLASOCIALBENCH* evaluates LLM sociopragmatic competence across three domains: Address Terms, Kinship Reasoning, and Social Customs.

### 3.1 Bangla Address Terms

Address forms are a primary mechanism through which Bangladeshi speakers manage social relations in interaction. Bangla employs a three-way second-person pronominal system that encodes different levels of social distance and respect. The form *apni* expresses the highest level of formality and is typically used for elders, authority figures, or strangers. The form *tumi* conveys neutral familiarity among peers or polite downward rela-

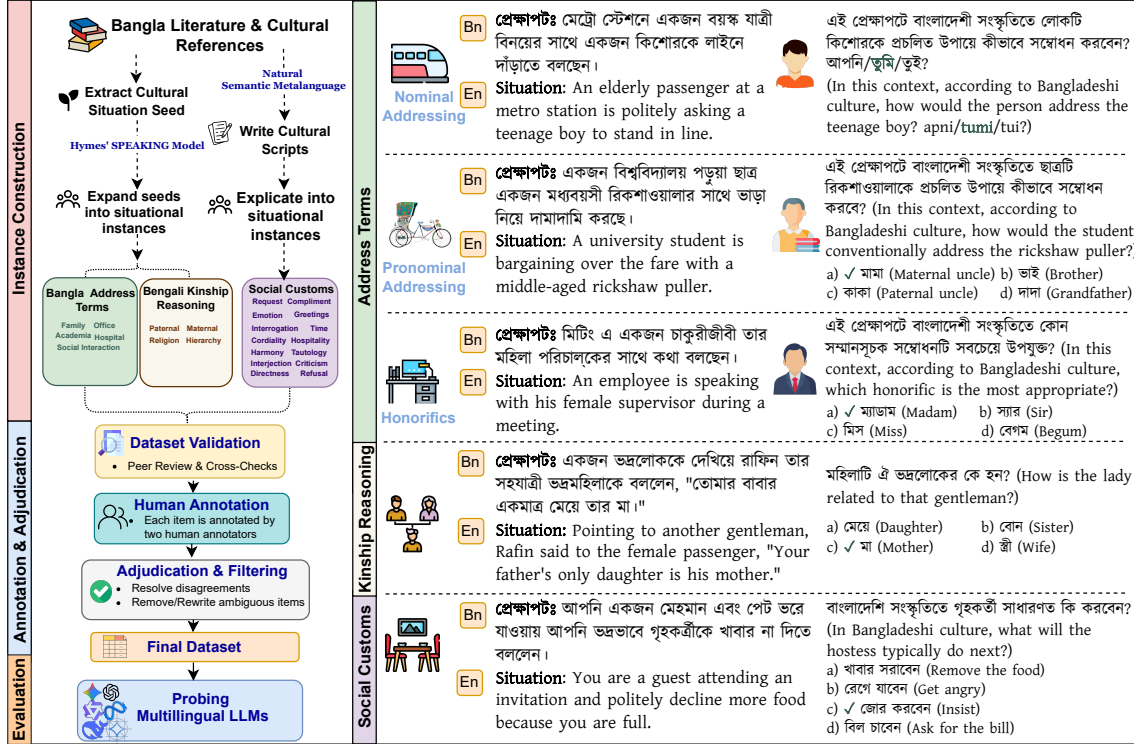


Figure 2: **Dataset creation pipeline for BANGLASOCIALBENCH.** The English prompts displayed in the diagram are translated for illustrative purposes; all model evaluations were conducted exclusively using Bangla prompts

tions. The form *tui* signals intimacy and is commonly used among close friends, siblings, or when addressing children. Because these forms directly encode hierarchy and relational distance, incorrect selection can sound rude or socially inappropriate despite grammatically correct content.

Beyond pronouns, Bangla widely uses kinship-based nominal addressing, where kinship terms are used both for relatives and for socially proximate non-relatives to signal age, gender, and hierarchy (Sultana et al., 2023). Honorific titles further refine these distinctions across professional, religious, and everyday settings (Das, 1968b). Together, these resources require speakers to infer appropriate address forms from contextual cues such as age difference, social role, and interactional goal.

To operationalize these distinctions, prompts are constructed following Hymes’ Ethnography of Communication and the SPEAKING framework (Hymes, 1964). Each instance specifies participants, setting, gender and interactional intent so that the correct address form must be inferred from sociopragmatic context rather than explicit lexical hints (Figure 1). Additional cultural context is pro-

vided in Appendix A, and example instances appear in Tables 4-9.

### 3.2 Bengali Kinship Reasoning

Bangla kinship terminology encodes detailed distinctions that require structured relational reasoning. Unlike English, which often collapses multiple relations into broad categories such as “uncle”, Bangla differentiates kinship roles based on lineage (paternal vs. maternal) and generational order (Kolenda and Aziz, 1981; Das, 1968b). For example, the language distinguishes *chacha* (father’s brother) from *mama* (mother’s brother), among many other lineage-specific terms.

To evaluate whether models can infer such relations, we construct a dataset of 345 kinship reasoning puzzles based on blood relations. Each instance presents a multi-hop lineage description and requires the model to identify the correct kinship term. The dataset spans maternal, paternal, and affinal relations, with distractor options designed to test sensitivity to lineage direction and generational hierarchy.

### 3.3 Bengali Social Customs

The Social Customs domain evaluates culturally grounded expectations that govern appropriate responses and behavior in everyday Bangladeshi interaction. Each instance presents a realistic scenario followed by four candidate responses that are grammatically plausible but differ in pragmatic appropriateness. Models must select the response that best aligns with commonly shared social norms.

Scenarios are derived from ethnographic descriptions of Bangladeshi social interaction and iteratively validated by native Bangla annotators. To structure these situations, we use cultural scripts (Goddard and Wierzbicka, 2004), formalized using Natural Semantic Metalanguage (Wierzbicka, 1996). This process identifies fourteen recurring scripts that capture interactional expectations (shown in Table 1).

### 3.4 BANGLASOCIALBENCH Development

#### 3.4.1 Instance Writing

Instance creation began with a cultural situation brainstorming phase conducted by six native Bangla annotators who are lifelong residents of Bangladesh and have extensive exposure to everyday Bangladeshi social interaction. Annotators first generated short descriptions of socially plausible interactions drawn from daily life.

To ensure cultural grounding, annotators were allowed to consult Bangladeshi literary works (short stories, drama, or novels) and ethnographic or sociological sources for inspiration. They were explicitly instructed not to copy text, characters, or narrative structures. All instances were written independently from scratch to reduce the possibility of memorization or data leakage from LLM training corpora.

These initial situation seeds were then expanded in a domain-specific manner. For **Address Terms**, scenarios were systematically varied along socio-pragmatic variables such as participant age, gender, key, end, and setting (Hymes, 1964). For **Kinship Reasoning**, annotators constructed relational puzzles requiring inference over Bangladeshi kinship structures, using external materials only as conceptual reference. For **Social Customs**, recurring cultural practices were abstracted into shared cultural scripts and instantiated as multiple-choice scenarios reflecting conventional Bangladeshi expectations. Detailed annotator recruitment and in-

Category	Subcategory	#
Bangla Address Term: Pronominal	Family	140
	Office	100
	Hospital	50
	Academia	100
	Social Interaction	200
Bangla Address Term: Nominal	Kinship	152
	Social Addressing	155
	Honorifics	85
Bengali Kinship Reasoning	Mixed Relations	345
Bengali Social Customs	Request	20
	Interrogation	14
	Compliment	40
	Refusal	20
	Tautology	20
	Interjection	20
	Indirect Directness	20
	Hospitality	40
	Greetings	20
	Emotion	20
	Harmony	20
	Cordiality	78
	Criticism	40
	Time	20
<b>Total</b>		<b>1719</b>

Table 1: Distribution of instances in BANGLASOCIALBENCH across address terms, kinship reasoning, and social customs domains.

stance creation procedures are described in Appendix C.

#### 3.4.2 Dataset Quality Control

**Annotation Procedure** Each instance was independently annotated by two annotators. Annotation guidelines emphasized selecting responses that reflect widely shared Bangladeshi social expectations rather than individual preferences. For pronominal address-term instances, annotators were allowed to specify both a primary and, when pragmatically acceptable, a secondary form.

#### Inter-Annotator Agreement Measurement

Agreement was measured using Cohen’s  $\kappa$  (Cohen, 1960). For single-answer subsets, the selected label was used directly. For pronominal subsets allowing two acceptable forms, agreement was computed using the *primary* label.

Across all applicable subsets, the global standard  $\kappa$  was 0.82, indicating strong agreement. For tasks allowing two valid pronouns, we additionally computed a binary  $\kappa$  (Artstein and Poesio, 2008) by treating each pronoun as an independent binary

decision and macro-averaging the results. This produced a score of 0.77. Detailed annotation protocols and subcategory-wise statistics are reported in Table 10 of Appendix D.

**Adjudication and Label Consolidation** All annotation disagreements were resolved through a two-stage process consisting of peer discussion followed by verification by senior annotators. Instances judged to be ambiguous, underspecified, or culturally contentious were revised or removed. For the Kinship Reasoning subset, which involves deterministic reasoning, standard agreement metrics were not applicable. Because each item corresponds to a single ground truth, any disagreement was treated as an annotation error and resolved through third-party adjudication.

**Dataset Filtering** Approximately 13% of initially drafted instances were excluded during quality control due to unresolved ambiguity or insufficient cultural grounding.

## 4 Experiments

**Models** We evaluate a total of 12 LLMs, comprising six closed-source and six open-weight models. The closed-source models include GPT-4o-mini, GPT-4o, Gemini Flash 2.0, Gemini Flash 2.5, Claude Haiku 3.5, Claude Sonnet 4. The open-weight models include LLaMA 3 8B Instruct, LLaMA 3.3 70B Instruct, Gemma 3 12B, Gemma 3 27B, Qwen 2.5 72B Instruct and DeepSeek V3.1. Together, these models span diverse scales, paradigms, and levels of Bangla support.

**Evaluation Protocol** All models are evaluated in a zero-shot setting using prompts written entirely in Bangla, without task-specific fine-tuning or in-context examples. A standardized prompt template is used across all models for comparability. Each prompt presents a contextual scenario and asks the model to select the culturally appropriate response in Bangladeshi social context. Decoding is performed with temperature set to 0 to ensure deterministic outputs. All task instructions are provided within the user prompt. Example prompts are shown in Appendix D.5.1.

**Dual Acceptable Forms** For Pronominal Addressing within (1) **Bangla Address Terms** domain, each instance is annotated with a primary pronoun and, when applicable, a secondary acceptable pronoun. This reflects the sociopragmatic re-

ality that multiple address forms may be culturally appropriate for the same interactional context. Accordingly, a model response is counted as correct if it matches either the primary or the secondary acceptable pronoun, ensuring that culturally appropriate variation in address forms is not penalized.

Because some instances specify two acceptable pronouns, we report an instance weighted random baseline of 42.2% for this subset of the benchmark. (Details can be found in Appendix D.5.2)

**Single Best Answer** For (2) **Kinship Reasoning**, Nominal Addressing within (1) **Bangla Address Forms**, and (3) **Social Customs**, evaluation is conducted using exact match accuracy over four-option Multiple-Choice Questions. A model response is considered correct if it selects the culturally appropriate option specified in the annotated reference. Under this setting, a uniform random baseline achieves 25% accuracy.

### 4.1 Evaluation of Language Models on BANGLASOCIALBENCH

**Performance varies across domains and model size** As shown in Table 2 and Figure 3, model performance varies substantially across domains, with Kinship Reasoning posing the greatest challenge. Accuracy ranges from 31.88% (LLaMA 3 8B) to 92.46% (Gemini 2.5 Flash) on Kinship Reasoning, compared to 46.23-81.16% for Address Terms and 46.67-89.49% for Social Customs. Results also reveal a consistent model-scale effect within model families. For example, within the LLaMA family, accuracy increases from 46.23% (8B) to 63.24% (70B) on Address Terms and from 31.88% to 51.01% on Kinship Reasoning. A similar trend appears for Gemma models (54.99% to 73.42% on Address Terms). To ensure that subcategory imbalance (e.g., 200 instances for Social Interaction vs. 50 for Hospital) does not skew our findings, we verified that micro-averaged and macro-averaged accuracies across the Address Terms domain deviate by < 2%, confirming our overall metrics accurately reflect true model competence.

Overall, larger proprietary models perform strongest, with Gemini 2.5 Flash achieving the highest accuracy on Address Terms and Kinship Reasoning, while Claude Sonnet 4 performs best on Social Customs. These results suggest that sociopragmatic competence improves with scale but remains uneven across domains.

Model	Address Term	Kinship Reasoning	Social Customs
GPT-4o mini	65.99	49.86	77.43
GPT-4o	78.51	67.83	86.15
Gemini 2.0 Flash	76.17	68.12	85.13
Gemini 2.5 Flash	<b>81.16</b>	<b>92.46</b>	86.41
Claude Haiku 3.5	61.91	44.06	71.02
Claude Sonnet 4	75.46	88.98	<b>86.67</b>
LLaMA 3 8B Instruct	46.23	31.88	42.82
LLaMA 3.3 70B Instruct	63.24	51.01	81.54
Gemma 3 12B	54.99	42.61	76.15
Gemma 3 27B	73.42	48.12	81.02
Qwen 2.5 72B Instruct	64.77	48.70	77.18
DeepSeek V3.1	64.77	64.06	83.33

Table 2: Overall zero-shot accuracy (%) across domains in BANGLASOCIALBENCH.

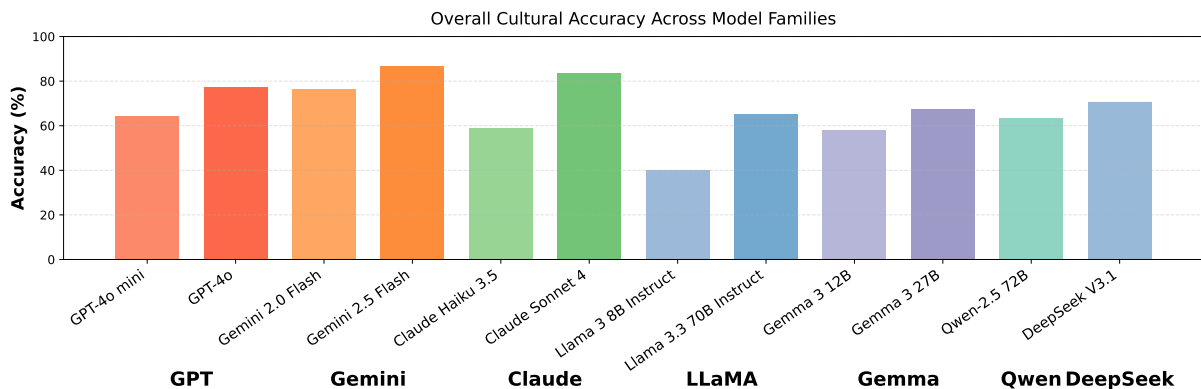


Figure 3: Overall benchmark accuracy of evaluated LLMs on BANGLASOCIALBENCH, computed as the macro-average of performance across Address Terms, Kinship Reasoning, and Social Customs.

### Most models exhibit systematic over-politeness

As shown in Figure 4, most models display a clear asymmetry in culturally inappropriate pronominal choices. A binomial test comparing over- and under-politeness errors shows a statistically significant imbalance for most models ( $p < 0.05$ ), with over-politeness occurring substantially more frequently. Several models, including DeepSeek V3.1 and Gemma 3 12B, frequently default to the highly formal pronoun *apni*, even in contexts where *tumi* or *tui* are culturally preferred. In contrast, under-politeness errors are rare across all model families. This pattern suggests that models adopt a conservative strategy, favoring safer or more polite forms rather than adapting flexibly to contextual cues such as age difference, intimacy, or interactional setting. This may reflect training data skew, since high-resource Bangla sources predominantly use formal registers.

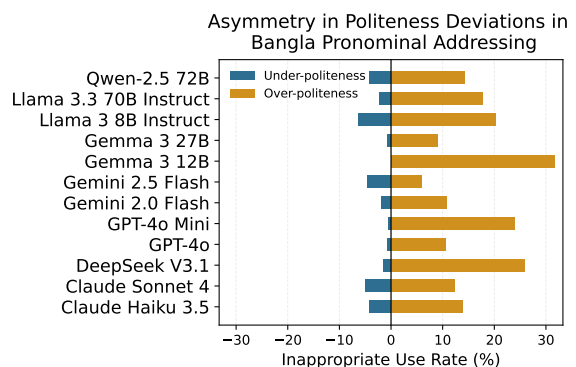


Figure 4: Asymmetry in inappropriate politeness use across LLMs in Bangla pronominal addressing, with over-politeness occurring substantially more frequently than under-politeness.

## 4.2 Analyzing Contextual Effects in Culturally Inappropriate Predictions

**Sociopragmatic factors linked to inappropriate Address Pronoun choice** To identify sys-

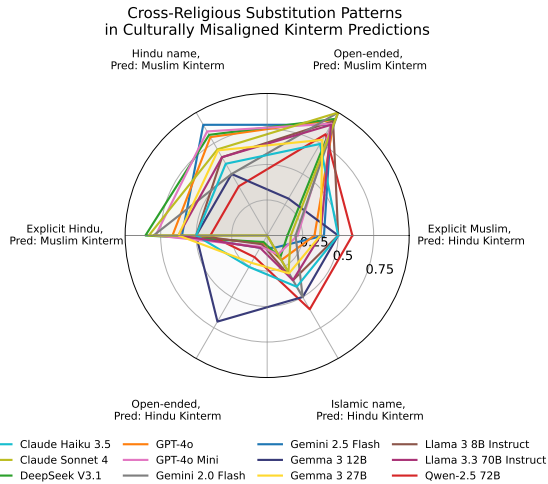


Figure 5: **Directional cross-religious kinship term misalignment.** Proportions of culturally inappropriate kinship term substitutions across explicit identity cues, implicit cues, and open-ended prompting. Misalignment is more pronounced toward substituting Muslim-associated kinterms in Hindu-marked contexts.

tematic sources of cultural misalignment, we conduct Pearson’s  $\chi^2$  tests of independence between culturally inappropriate pronominal address term choices and five sociopragmatic factors: *Participant Age*, *Setting*, *Key*, *End*, and *Participant Gender*, evaluated over  $N = 590$  instances per model. Standardized residual analysis further indicates that culturally inappropriate address choices concentrate in downward-hierarchy age relations (Elder→Younger) and informal or unstructured settings, whereas upward-hierarchy interactions and formal contexts exhibit relative robustness. Full statistics and dataset balance details are reported in Appendix D.5.3 and Appendix B (Figures 8, 9).

**Language Models conflate kinship terms across religious identities** Culturally inappropriate predictions display clear directional asymmetries across models. In all experiments, models were instructed to answer according to Bangladeshi cultural practices across multiple prompt strategies, including explicit religious identity markers (e.g., “a Muslim person”), religiously indicative personal names (e.g., “Arnab Roy,” a culturally Hindu name), and open-ended prompts without overt religious cues. Several models, including GPT-4o, DeepSeek V3.1, and Gemini 2.5 Flash, show a stronger tendency to substitute Muslim-associated kinship terms in explicitly or implicitly Hindu contexts, while fewer models exhibit the

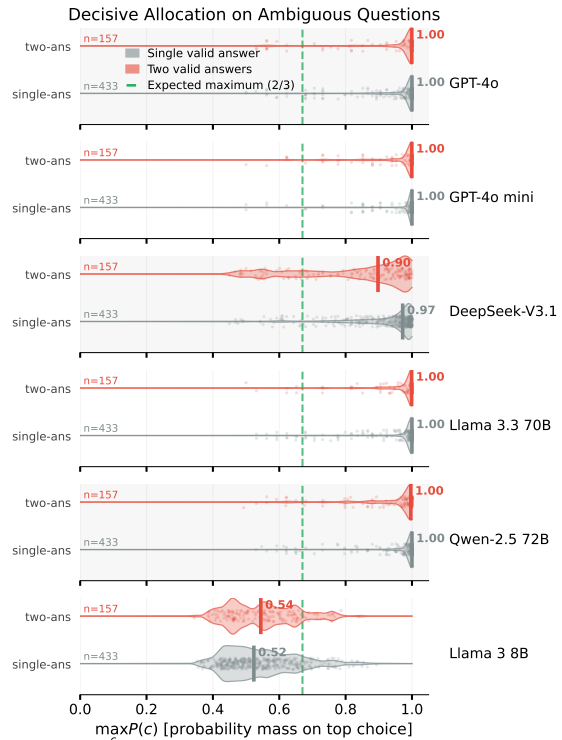


Figure 6: **Probability mass allocation between ambiguous and unambiguous questions.** Distribution of  $\max_c P(c)$ , the probability mass assigned to the most probable choice, over single-answer and two-answer questions, recovered from token-level log-probabilities via choice-restricted softmax. A well-calibrated model should show lower probability mass on two-answer questions than on single-answer ones.

reverse pattern. To quantify this asymmetry, we perform a  $\chi^2$  goodness-of-fit test for each model over all culturally inappropriate predictions involving religious mismatch. The number of such cases varies by model ( $N = 24$ – $65$ ). Eleven of twelve models exhibit statistically significant asymmetry ( $p < 0.05$ ), with a pronounced skew toward substituting Muslim-associated kinship terms in Hindu-marked contexts. Detailed results appear in Appendix D.10.1.

**Log-probabilities reveal decisive probability allocation on ambiguous situations** For models with accessible token-level log probabilities (Figure 6), we recover log-probabilities via choice-restricted softmax over pronoun logits using the top- $k$  alternatives ( $k = 5$ ) returned by the API (Appendix D.11.1), and compute  $\max_c P(c)$  per question. Bangla tokenization inconsistencies were mitigated by mapping answer options to letter labels (A-C) and restricting model outputs to

these labels. Questions are split into unambiguous ( $n=433$ , single valid answer) and ambiguous ( $n=157$ , two valid answers). We use  $\frac{2}{3} \approx 0.67$  as a reference boundary, above which one form holds more than twice the probability mass of any alternative. Figure 6 shows that five of six models produce median  $\max_c P(c) = 1.00$  on ambiguous questions, identical to unambiguous ones, indicating they assign near-total probability mass to a single pronoun regardless of sociopragmatic ambiguity. Only Llama 3 8B deviates (median 0.54), though still near the 0.67 boundary. This suggests that current LLMs allocate probability mass as decisively on sociopragmatically ambiguous address forms as on unambiguous ones, without reflecting the variation acknowledged by human annotators.

## 5 Conclusion

Cultural competence in language models requires not only linguistic fluency but also sensitivity to sociopragmatic context. Using BANGLASOCIALBENCH, we show that even strong LLMs exhibit systematic limitations in Bangla address terms, kinship reasoning, and social customs, often defaulting to overly polite forms, lacking calibration in ambiguous contexts, and conflating kinship terminology across religious identities. Addressing the sociopragmatic weaknesses exposed by our evaluation requires both data-centric and algorithmic interventions. We hypothesize that deficits such as systematic over-politeness and cross-religious kinship conflation stem from a training data imbalance. Current multilingual LLMs are primarily exposed to formal or translated Bangla corpora, leaving informal contexts, downward-hierarchy interactions, and minority-community norms systematically under-represented. While training native Bengali LMs from scratch may offer deeper cultural embedding, the computational requirements for general reasoning parity remain prohibitive. A more pragmatic approach lies in culturally grounded alignment. Future research could utilize culturally specific adapters or Direct Preference Optimization (DPO) with contrastive pairs derived from the failure modes identified in BANGLASOCIALBENCH. Such methods would allow practitioners to leverage the reasoning capabilities of existing models while anchoring them more precisely to native Bengali relational norms.

## Limitations

While BANGLASOCIALBENCH provides a focused evaluation of sociopragmatic competence in Bangladeshi social interaction, several limitations remain.

**Standard Colloquial Bangla vs. dialect** All instances and annotations are grounded in Standard Colloquial Bangla, reflecting widely shared norms of address usage and social interaction in Bangladesh. Annotators were instructed to judge cultural appropriateness with respect to this interactional standard. As a result, the benchmark does not capture the regional and dialectal variation present across Bangla varieties. Address pronouns, kinship terms, and honorifics may differ in regional varieties such as Sylheti, Chittagonian, or Rangpuri, and culturally appropriate usage in these dialects may not always align with the norms encoded in our dataset.

We restrict the benchmark to Standard Colloquial Bangla for two reasons. First, it functions as a common interactional baseline used in interregional communication, education, media, and institutional settings. Second, regional dialects differ not only lexically but also in their pragmatic conventions, which would require dialect-specific annotation guidelines and expert annotators for each variety. Without sufficient dialectal coverage, such variation could be misinterpreted as model errors.

**Coverage of interactional contexts** Although the dataset spans multiple sociopragmatic settings, including family, office, hospital, academic, and everyday social interactions, it cannot exhaustively represent all interactional configurations present in Bangladeshi society. Factors such as rural-urban differences, sociolectal variation, and code-switching practices are not explicitly modeled.

Consequently, models that generate regionally appropriate but non-standard forms may be penalized under our evaluation framework. This should not be interpreted as a claim that Standard Colloquial Bangla is inherently more correct, but rather as a methodological constraint required to ensure annotation consistency and comparability across models. Future work may extend the benchmark to incorporate dialectal and regional sociopragmatic variation.

## Ethical Considerations

Annotators were compensated for their work at a rate consistent with local standards for similar annotation tasks. Prior to participation, annotators were informed about the purpose of the study and how their contributions would be used in the dataset.

To protect privacy, annotators were instructed not to include any personally identifiable information when writing scenarios. All collected instances were additionally reviewed by the research team to ensure that they did not contain harmful, unsafe, or sensitive content. The benchmark focuses on culturally grounded everyday interactions and does not include sexual, violent, or personally identifiable material.

## References

- Muhammad Farid Adilazuarda, Sagnik Mukherjee, Pradhyumna Lavania, Siddhant Singh, Alham Fikri Aji, Jacki O'Neill, Ashutosh Modi, and Monojit Choudhury. 2024. *Towards Measuring and Modeling "Culture" in LLMs: A Survey*. Preprint, arXiv:2403.15412.
- Badr AlKhamissi, Muhammad ElNokrashy, Mai Alkhamissi, and Mona Diab. 2024. *Investigating Cultural Alignment of Large Language Models*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12404–12422, Bangkok, Thailand. Association for Computational Linguistics.
- Ron Artstein and Massimo Poesio. 2008. *Survey Article: Inter-Coder Agreement for Computational Linguistics*. *Computational Linguistics*, 34(4):555–596.
- Lama Ayash, Hassan Alhuzali, Ashwag Alasmari, and Sultan Aloufi. 2025. *SaudiCulture: A Benchmark for Evaluating Large Language Models Cultural Competence within Saudi Arabia*. Preprint, arXiv:2503.17485.
- Anne Barron. 2003. *Acquisition in Interlanguage Pragmatics: Learning How to Do Things with Words in a Study Abroad Context*. Ph.D. thesis, J. Benjamins, Amsterdam Philadelphia (Pa.).
- Friederike Braun. 1988. *Terms of Address: Problems of Patterns and Usage in Various Languages and Cultures*. DE GRUYTER.
- Roger Brown and Albert Gilman. 1968. *THE PRO-NOUNS OF POWER AND SOLIDARITY*. In Joshua A. Fishman, editor, *Readings in the Sociology of Language*, pages 252–275. DE GRUYTER.
- Xinren Chen and Juanjuan Ren. 2020. *A memetic cultural practice: The use of generalized kinship terms in a research seminar attended by Chinese graduate students*. *Lingua*, 245:102942.
- Jacob Cohen. 1960. *A Coefficient of Agreement for Nominal Scales*. *Educational and Psychological Measurement*, 20(1):37–46.
- Sisir Kumar Das. 1968a. *Forms of address and terms of reference in Bengali*. *Anthropological Linguistics*, 10(4):19–31.
- Sisir Kumar Das. 1968b. *Forms of Address and Terms of Reference in Bengali*. *Anthropological Linguistics*, 10(4):19–31.
- Eleanor Dickey. 1997. *Forms of address and terms of reference*. *Journal of Linguistics*, 33(2):255–274.
- Mohammad Rifqi Farhansyah, Iwan Darmawan, Adryan Kusumawardhana, Genta Indra Winata, Alham Fikri Aji, and Derry Tanti Wijaya. 2025. *Do Language Models Understand Honorific Systems in Javanese?* In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 26732–26754, Vienna, Austria. Association for Computational Linguistics.
- Cliff Goddard and Anna Wierzbicka. 2004. *Cultural scripts: What are they and what are they good for?* *Intercultural Pragmatics*, 1(2).
- Dell Hymes. 1964. *Introduction: Toward Ethnographies of Communication*<sup>1</sup>. *American Anthropologist*, 66(6\_PART2):1–34.
- D.H. Hymes. 1962. *The Ethnography of Speaking*. The Ethnography of Speaking. Anthropological Society of Washington.
- Saman Sarker Joy. 2025. *BnMMLU: Measuring Massive Multitask Language Understanding in Bengali*. Preprint, arXiv:2505.18951.
- Daen Kabir, Minhajur Rahman Chowdhury Mahim, Sheikh Shafayat, Adnan Sadik, Arian Ahmed, Eunsu Kim, and Alice Oh. 2025. *BLUCK: A Benchmark Dataset for Bengali Linguistic Understanding and Cultural Knowledge*. Preprint, arXiv:2505.21092.
- Mohammad Hossein Keshavarz. 2001. *The role of social context, intimacy, and distance in the choice of forms of address*. *International Journal of the Sociology of Language*, 2001(148).
- Eunsu Kim, Juyoung Suk, Philhoon Oh, Haneul Yoo, James Thorne, and Alice Oh. 2024. *CLiCK: A Benchmark Dataset of Cultural and Linguistic Intelligence in Korean*. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3335–3346, Torino, Italia. ELRA and ICCL.

- Pauline Kolenda and K. M. Ashraful Aziz. 1981. [Kinship in Bangladesh](#). *Pacific Affairs*, 54(4):744.
- Fajri Koto, Rahmad Mahendra, Nurul Aisyah, and Timothy Baldwin. 2024. [IndoCulture: Exploring geographically-influenced cultural common-sense reasoning across eleven Indonesian provinces](#). *Transactions of the Association for Computational Linguistics*, 12:1515–1532.
- Arijit Maji, Raghvendra Kumar, Akash Ghosh, Anushka, and Sriparna Saha. 2025. [SANSKRITI: A comprehensive benchmark for evaluating language models’ knowledge of indian culture](#). *arXiv preprint arXiv:2506.15355*.
- Reem Masoud, Ziquan Liu, Martin Ferianc, Philip C. Treleaven, and Miguel Rodrigues Rodrigues. 2025. [Cultural Alignment in Large Language Models: An Explanatory Analysis Based on Hofstede’s Cultural Dimensions](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 8474–8503, Abu Dhabi, UAE. Association for Computational Linguistics.
- Anjishnu Mukherjee, Chahat Raj, Ziwei Zhu, and Antonios Anastasopoulos. 2023. [Global Voices, Local Biases: Socio-Cultural Prejudices across Languages](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 15828–15845, Singapore. Association for Computational Linguistics.
- Tarek Naous, Michael J. Ryan, Alan Ritter, and Wei Xu. 2024. [Having beer after prayer? measuring cultural bias in large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16366–16393. Association for Computational Linguistics.
- Fatma Hülya Özcan. 2016. [Choice of address terms in conversational setting](#). *International Journal of Human Sciences*, 13(1):982.
- Abhinav Sukumar Rao, Akhila Yerukola, Vishwa Shah, Katharina Reinecke, and Maarten Sap. 2025. [NormAd: A Framework for Measuring the Cultural Adaptability of Large Language Models](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2373–2403, Albuquerque, New Mexico. Association for Computational Linguistics.
- Shahriyar Zaman Ridoy, Azmine Toushik Wasi, and Koushik Ahamed Tonmoy. 2025. [BengaliMoralBench: A benchmark for auditing moral reasoning in large language models within bengali language and culture](#). *arXiv preprint arXiv:2511.03180*.
- Nikta Gohari Sadr, Sahar Heidariasl, Karine Megerdoo-mian, Laleh Seyyed-Kalantari, and Ali Emami. 2025. [We Politely Insist: Your LLM Must Learn the Persian Art of Taarof](#). *arXiv preprint*.
- Hyopil Shin, Sangah Lee, Dongjun Jang, Wooseok Song, Jaeyoon Kim, Chaeyoung Oh, Hyemi Jo, Youngchae Ahn, Sihyun Oh, Hyohyeong Chang, Sunkyoung Kim, and Jinsik Lee. 2025. [KoBALT: Korean Benchmark For Advanced Linguistic Tasks](#). *Preprint*, arXiv:2505.16125.
- Shayla Sharmin Snigdha. 2022. [Representation of social class and hierarchy in Bangla address terms](#). *Journal of Modern Languages*, 32(2):107–1529.
- Settaluri Sravanthi, Meet Doshi, Pavan Tankala, Rudra Murthy, Raj Dabre, and Pushpak Bhattacharyya. 2024. [PUB: A Pragmatics Understanding Benchmark for Assessing LLMs’ Pragmatics Capabilities](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 12075–12097, Bangkok, Thailand. Association for Computational Linguistics.
- Shaila Sultana, Mian Md. Naushaad Kabir, Md. Zulfeqar Haider, Mohammad Moninoor Roshid, and M. Obaidul Hamid. 2023. [Language in Society in Bangladesh and Beyond: Voices of the Unheard in the Global South](#), 1 edition. Routledge, New York.
- Tianda Sun and Dimitar Kazakov. 2026. [Kinship Data Benchmark for Multi-hop Reasoning](#). *Preprint*, arXiv:2601.07794.
- Md. Afaz Uddin. 2019. [Second Person Pronouns as Person Deixis in Bengali and English: Linguistic Forms and Pragmatic Functions](#). *International Journal of English Linguistics*, 10(1):345.
- Anna Wierzbicka. 1996. [Semantics: Primes and Universals](#). Oxford University PressOxford.
- Yitong Zhang, Yuhan Xiang, and Mingxuan Liu. 2026. [The Mask of Civility: Benchmarking Chinese Mock Politeness Comprehension in Large Language Models](#). *Preprint*, arXiv:2602.03107.
- Wenlong Zhao, Debanjan Mondal, Niket Tandon, Danica Dillion, Kurt Gray, and Yuling Gu. 2024. [WorldValuesBench: A Large-Scale Benchmark Dataset for Multi-Cultural Value Awareness of Language Models](#). *Preprint*, arXiv:2404.16308.
- Naitian Zhou, David Bamman, and Isaac L. Bleaman. 2025. [Culture is not trivia: Sociocultural theory for cultural NLP](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 25869–25886, Vienna, Austria. Association for Computational Linguistics.

## A Additional Cultural Context

### A.1 Bangla Address Forms

Address forms play a central role in shaping social interaction in a language. They organize how speakers relate to one another, influence the flow of conversation, and signal identity, relationship, and context. Address forms are culturally specific, and their use connects the semantic meaning of a term with the pragmatic demands of a situation. They embody the cultural concepts, values, and social principles of a community and they help reveal how members position themselves in interaction (Chen and Ren, 2020). Knowledge of variation in address forms is essential in intercultural settings because it reduces the risk of sociopragmatic failure (Barron, 2003). To avoid the communication failure, it is essential that LLMs not only know the address terms but also understand their pragmatics.

The formality or informality of a situation and the level of politeness or deference a speaker intends to express are reflected directly in the chosen address form (Özcan, 2016). These expressions function as linguistic markers of social and cultural identity. For language models, understanding kinship, solidarity, and power distance is especially important in Bangladeshi communication, where these relationships are strongly mirrored in conversational patterns.

A nuanced understanding of the diversity and cultural significance of Bangla address terms is central to interpreting Bangladeshi culture. This makes it important to evaluate how well LLMs capture the cultural meanings and hierarchical cues encoded in these address forms, as address terms and honorifics together form an essential part of social deixis in the language (Keshavarz, 2001).

#### **The Pronoun of Power, Solidarity and Intimacy**

It is important to evaluate whether LLMs can recognize the social variables that shape pronominal address in Bangla, such as age, power, social distance, and intimacy. Languages differ widely in how they mark these distinctions. English relies mainly on nominal forms and lack a distinction in address pronouns (Dickey, 1997), while several Asian languages use personal pronouns sparingly. Many European languages follow the T-V distinction, which separates familiar and polite forms (Brown and Gilman, 1968).

Bangla offers a richer system with not two, but three second-person pronouns, *apni*, *tumi*, and *tui*, all meaning you but carrying different honorific and relational meanings. These forms index respect, equality, or intimacy and are closely tied to cultural expectations about how speakers negotiate hierarchy and solidarity (Snigdha, 2022). The three-way pronominal system therefore provides an informative setting for studying the sociolinguistic competence of LLMs, since correct pronoun selection requires an understanding of social hierarchy, solidarity, and interactional norms. Models that cannot infer these relational cues are likely to generate socially inappropriate or pragmatically incorrect responses.

#### **Nominal Addressing in Kinship and Social Contexts**

In Bangladeshi society, kinship terms reflect relational meaning between speakers and addressees, with distinctions shaped by paternal and maternal lineage as well as religious and cultural contexts (Sultana et al., 2023; Snigdha, 2022; Kolenda and Aziz, 1981). The pragmatics of nominal address forms show that the social meaning of a kinship term used in interaction often differs from its literal meaning (Braun, 1988). It is common in Bangladesh to address unfamiliar people with kinship-based terms that take into account age, gender, and perceived social position. Addressing elders, relatives, or even non-relatives directly by name or by surname alone is often considered disrespectful, which further reinforces the preference for kinship-based address. Because Bangladeshi society is deeply rooted in familial ties, the use of kinship terms for neighbors, acquaintances, and strangers has become a widespread and culturally expected practice.

Kinship terminology in Bangla also varies across religious communities. Muslim, Hindu, and other groups share many core kin terms but maintain distinct forms shaped by historical and ritual traditions, such as *thakur-dada* and *thakur-ma* in Hindu families, or *abba-jan* and *amma-jan* in traditional Muslim households (Snigdha, 2022). Capturing this variation is important for cultural knowledge, and our dataset therefore includes commonly used kinship terms across communities without overrepresenting any single register.

Kinship terms extend well beyond biological relations in everyday Bangladeshi interaction. They frequently function as social deictics for neigh-

bors, acquaintances, and strangers, with speakers choosing forms like *bhai* (brother), *apa* (elder sister), *chacha* (uncle), or *nana* based on perceived age, gender, and social position. This broad extension of kinship into general social life makes kinship reasoning a crucial component of culturally appropriate communication.

**Honorifics** Honorifics in Bangla signal respect, politeness, and social position, and they appear in both prefix and suffix forms (Das, 1968b). Some are influenced by English conventions, such as Mr., Ms., and Mrs., which are used in formal and educated settings. Bangla also has a rich set of local honorifics that vary across professional, religious, and social contexts. Common suffixes include *saheb*, *janab*, and *babu* for men, and *begum* or *debi* for women. In workplace or service settings, speakers often use forms like *sir*, *madam*, and *apa* to index respect and institutional hierarchy (Snigdha, 2022).

Honorific use also overlaps with kinship-based addressing, which is widespread in Bangladeshi society. Speakers frequently use terms such as *bhai*, *apa*, *bhabi*, *chacha*, or *dadu* for acquaintances or strangers based on perceived age, gender, and social position. Directly addressing elders, relatives, or even non-relatives by name alone, or by name plus surname, is often considered inappropriate, which strengthens the preference for kinship-inflected honorifics in everyday communication. These practices illustrate how Bangla honorifics encode politeness, hierarchy, and social affiliation, making them an informative component for evaluating the sociocultural competence of language models.

## A.2 Bengali Kinship Reasoning

Bangla kinship terminology is finely differentiated and socially encoded, reflecting lineage, marriage relations, generational hierarchy, and culturally grounded expectations of respect and obligation. Kolenda and Aziz (1981) shows that kin terms form a core interpretive framework through which interpersonal roles and social relationships are understood. Unlike English, which relies on broad categories such as uncle, aunt, and cousin, Bangla employs a highly granular system with dozens of specific terms such as *chacha* (father's brother), *mama* (mother's brother), *fupu* (father's sister), and *dada* (paternal grandfather). This precision encodes critical relational information in-

cluding lineage direction, generational rank, and relative age, making kinship not only descriptive but also a system that reinforces social hierarchy and regulates appropriate forms of respect within families and communities.

Evaluating LLMs on Bangla Kinship Reasoning is therefore valuable because correct interpretation and use of kin terms require cultural reasoning, social inference, and an understanding of relational structure that cannot be captured through grammar or vocabulary alone. A model that fails to grasp these distinctions risks producing socially inappropriate or culturally incoherent outputs.

## A.3 Bangladeshi Social Customs

The Social Customs domain targets culturally grounded expectations that govern appropriate behavior and response selection in everyday social interaction in Bangladesh. This subset focuses on interactional situations where multiple responses are possible, but cultural norms favor non-imposing, indirect, and harmony-preserving strategies. Instances were constructed using a script-guided instance writing approach. Cultural scripts (Goddard and Wierzbicka, 2004) were used as an abstraction layer to capture shared interactional expectations, such as softening requests, avoiding direct refusal, expressing hospitality, and managing social tension.

## B Definition and Example Instances Across Domains

### B.1 Annotation Definitions

#### B.1.1 Pronominal Addressing

For *Pronominal Addressing*, each instance is annotated along four sociopragmatic dimensions that jointly determine appropriate second-person pronoun choice (*apni, tumi, tui*) in Bangla.

**Setting** The social domain in which the interaction takes place, such as *Family, Office, Hospital, Academia, or Social Interaction*. The setting constrains expected norms of formality and power relations.

**Participant Age Relation** The relative age hierarchy between the speaker and the addressee:

- **Elder→Younger:** the speaker is older than the addressee
- **Younger→Elder:** the speaker is younger than the addressee
- **Peer→Peer:** the speaker and addressee are of comparable age

Age hierarchy is a primary determinant of pronominal choice in Bangla and often overrides lexical politeness markers.

**Participant Gender Relation** The gender configuration of the interaction:

- Female→Male
- Male→Female
- Female→Female
- Male→Male

Gender interacts with age and setting, particularly in family and institutional contexts.

**Key** The affective tone or manner of the interaction, following Hymes (1962)'s notion of *Key*:

- **Neutral:** emotionally unmarked interaction
- **Informal:** relaxed or intimate tone
- **Polite:** deferential or respectful tone
- **Angry:** emotionally charged or confrontational tone

**End** The communicative goal of the utterance:

- **Inform:** conveying information
- **Inquire:** asking a question
- **Request:** politely seeking action or permission
- **Command:** issuing a directive

Correct pronominal choice requires jointly inferring these variables rather than relying on surface lexical cues.

#### B.1.2 Nominal Addressing

For *Nominal Addressing*, instances are annotated to capture sociocultural variation in kinship and social titles, particularly where religious identity influences lexical choice.

**Identity** Whether the prompt specifies the interlocutor's religious identity:

- **Muslim**
- **Hindu**

**Type** How identity information is conveyed in the prompt:

- **Explicit:** religious identity is directly mentioned
- **Implicit:** identity is inferred through culturally indicative personal names or contextual cues

**Basis** The underlying reason a particular nominal form is appropriate:

- **Religion-based:** the address term differs across religious communities
- **Social-based:** the term reflects social hierarchy or familiarity independent of religion

#### B.1.3 Social Customs

The *Social Customs* domain evaluates culturally appropriate response selection in everyday Bangladeshi social interaction. Unlike Address Terms and Kinship Reasoning, which focus on lexical choice, this domain operationalizes culture as shared expectations about *how one should respond* in common social situations.

Following the Natural Semantic Metalanguage (NSM) framework (Goddard and Wierzbicka, 2004), we identify **14 recurrent cultural scripts** that reflect widely shared Bangladeshi customs governing indirectness, politeness, emotional restraint, hospitality, and interpersonal harmony. Each script is formulated as a simple, reductive mental model using universal semantic primitives, allowing cultural expectations to be expressed without theoretical jargon and providing a structured basis for data and instance creation.

<b>ID</b>	<b>Topic</b>	<b>Cultural Script (NSM)</b>
1.	Request	People think like this: when I want someone to do something, it is not good to say it strongly. If I say it strongly, this person can feel something bad. Because of this, it is good to say it softly, so this person can choose freely.
2.	Interrogation	People think like this: some things are personal. Asking about these things can feel intrusive. If I ask directly, people can think I have bad manners. Because of this, it is good not to ask directly.
3.1	Compliment	People think like this: when someone says something good about me, it is not good to agree openly. People can think I think too much of myself. Because of this, it is good to say little or deny it.
3.2	Compliment	People think like this: when someone praises me too much, this person may want something from me. Because of this, it is good to respond in a way that stops more praise.
4.	Refusal	People think like this: if I do not want to do something, I should not say “no” directly. If I do, this person can feel bad. It is better to say “I will try” or “I will see.”
5.	Tautology	People think like this: some things happen not because people want them. People cannot change these things. Saying words like this helps people feel calm.
6.	Interjection	People think like this: sometimes I feel something suddenly. It is good to say a small sound or word so other people know what I feel.
7.	Indirectness	People think like this: when I want someone to do something, it is not good to order them. This person can feel pushed. Because of this, it is good to say it indirectly.
8.1	Hospitality	People think like this: when someone comes to my place, it is good to offer food or drink. If I do not do this, people can think badly of me.
8.2	Hospitality	People think like this: when a guest says “I do not want to eat,” this person says this to be polite. Because of this, it is good to insist many times.
9.	Greetings	People think like this: when I see someone, it is good to say something first. Not saying anything can feel cold or bad.
10.	Emotion	People think like this: when I feel something strong, it is not always good to show it in public. Other people can feel uncomfortable. Because of this, it is good to control feelings.
11.	Harmony	People think like this: it is good when people feel good together. It is not good to create tension. Because of this, people speak and act gently.
12.	Cordiality	People think like this: people close to me are like family. I want to be close to them, help them, and spend time with them. Distance is not good.
13.	Criticism	People think like this: if someone does something wrong, it is not good to say this in front of others. This person can feel shame. It is better to say it softly or later.
14	Time	People think like this: when I say a time, it does not have to be exact. Being with people is more important than strict timekeeping.

Table 3: Cultural scripts in the Social Customs domain, expressed in Natural Semantic Metalanguage (NSM) and used as an abstraction layer to guide instance creation

Subcategory	Example Instances
<b>Address Terms: Pronominal</b>	
<b>Setting:</b> Family <b>Age:</b> Elder→Younger <b>Gender:</b> Female→Male <b>Key:</b> Inquire <b>End:</b> Angry	<p>এক মা তার কিশোর ছেলেকে রাগান্বিতভাবে তার পড়াশোনার অবস্থা সম্পর্কে জিজ্ঞাসা করছেন।  <i>(A mother angrily asks her teenage son about the state of his studies.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশের শহুরে সামাজিক ব্যবস্থায় মা ছেলেকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, within Bangladesh' urban social setting, how would a mother conventionally address her son?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Family <b>Age:</b> Elder→Younger <b>Gender:</b> Female→Female <b>End:</b> Inquire <b>Key:</b> Informal	<p>এক বৃদ্ধ খালা তার শিশু ভাগ্নিকে আদুরে স্বরে প্রশ্ন করছেন।  <i>(An elderly maternal aunt affectionately asks her young niece.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশের শহুরে সামাজিক ব্যবস্থায় খালা তার ভাগ্নিকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, within the urban social setting of Bangladesh, how would a maternal aunt conventionally address her niece?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Office <b>Age:</b> Younger→Elder <b>Gender:</b> Male→Male <b>End:</b> Request <b>Key:</b> Polite	<p>একজন মধ্যবয়স্ক সহকারী সচিব তার বয়স্ক উপসচিবকে সরকারি তহবিল ব্যবহারে বিনীত অনুরোধ করছেন।  <i>A middle-aged Assistant Secretary is politely requesting an older Deputy Secretary regarding the use of government funds.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী সংস্কৃতিতে সহকারী সচিব উপসচিবকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবেন?  <i>(In this context, within Bangladeshi culture, how would the Assistant Secretary conventionally address the Deputy Secretary?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Office <b>Age:</b> Elder→Younger <b>Gender:</b> Female→Female <b>End:</b> Inform <b>Key:</b> Angry	<p>একজন বয়স্ক নারী উপ-পরিচালক তার মধ্যবয়স্ক নারী সহকারি পরিচালককে পাসপোর্ট যাচাইকরণের নতুন নিয়মাবলী সম্পর্কে রাগান্বিত হয়ে অবহিত করছেন।  <i>(An elderly female Deputy Director is angrily informing her middle-aged female Assistant Director about the new regulations for passport verification.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী সংস্কৃতিতে উপ-পরিচালক সহকারি পরিচালককে প্রচলিত উপায় কিভাবে সম্বোধন করবেন?  <i>(In this context, within Bangladeshi culture, how would the Deputy Director conventionally address the Assistant Director?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Academia <b>Age:</b> Younger→Elder <b>Gender:</b> Male→Female <b>End:</b> Request <b>Key:</b> Polite	<p>স্কুলের একজন দ্বিতীয় শ্রেণীর ছাত্র বিনয়ের সাথে তার গণিত শিক্ষিকাকে অংকটি দেখাতে অনুরোধ করল।  <i>(A second-grade student politely asks his mathematics teacher to look at a problem.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী শহুরে সামাজিক ব্যবস্থায় দ্বিতীয় শ্রেণীর ছাত্রটি তার গণিত শিক্ষিকাকে প্রচলিত উপায় কিভাবে সম্বোধন করবে?  <i>(In this context, within Bangladesh' urban social setting, how would the second-grade student conventionally address his mathematics teacher?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Academia <b>Age:</b> Peer→Peer <b>Gender:</b> Male→Female <b>End:</b> Command <b>Key:</b> Angry	<p>পঞ্চম শ্রেণীর একজন ছাত্র তার একই ক্লাসের বান্ধবীকে রাগান্বিত স্বরে ডাকছে।  <i>(A fifth-grade male student is angrily addressing a female classmate from the same class.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী শহুরে সামাজিক ব্যবস্থায় তারা একে অন্যকে প্রচলিত উপায় কিভাবে সম্বোধন করবে?  <i>(In this context, within Bangladesh' urban social setting, how would the two conventionally address each other?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Hospital <b>Age:</b> Younger→Elder <b>Gender:</b> Female→Male <b>End:</b> Inquire <b>Key:</b> Neutral	<p>একজন তরুণ বয়সী মহিলা স্টাফ নার্স মধ্যবয়স্ক পুরুষ ইমার্জেন্সি মেডিকেল অফিসারকে ঠিকঠাক মাত্রা নিয়ে জিজ্ঞাসা করল।  <i>(A young female staff nurse asks a middle-aged male emergency medical officer about the medication dosage.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী সংস্কৃতিতে স্টাফ নার্স ইমার্জেন্সি মেডিকেল অফিসারকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবেন?  <i>(In this context, within Bangladeshi culture, how would the staff nurse conventionally address the emergency medical officer?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Social Interaction <b>Age:</b> Elder→Younger <b>Gender:</b> Female→Female <b>End:</b> Inform <b>Key:</b> Informal	<p>আফরিনের বাসায় প্রতিবেশীর ছোট একটি মেয়ে শিশু বেড়াতে আসলে আফরিন তাকে আদর করে স্কুলের মজার ঘটনার বিবরণ দিচ্ছে।  <i>(When a young girl from a neighboring family visits Afrin's home, Afrin affectionately tells her about a funny incident from school.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশের শহুরে সামাজিক ব্যবস্থায় আফরিন শিশুটিকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, within an urban Bangladeshi social setting, how would Afrin typically address the child?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>
<b>Setting:</b> Social Interaction <b>Age:</b> Younger→Elder <b>Gender:</b> Male→Male <b>End:</b> Command <b>Key:</b> Neutral	<p>একজন কিশোর বয়সী ছেলে বয়স্ক মুচিকে জুতা পরিষ্কার করতে বলছে।  <i>(A teenage boy is asking an elderly cobbler to clean his shoes.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশের শহুরে সামাজিক ব্যবস্থায় কিশোর মুচিকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, within an urban Bangladeshi social setting, how would the teenager typically address the cobbler?)</i></p> <p>A. আপনি (Apni)  B. তুমি (Tumi)  C. তুই (Tui)</p>

Table 4: Example instances from the Pronominal Addressing subset of *Address Terms* domain. Each instance is further annotated with sociopragmatic metadata, namely Setting, Participant Age, Participant Gender, Key, and End.

Subcategory	Example Instances
	<b>Address Terms: Nominal</b>
<b>Subtopic:</b> Kinship <b>Identity:</b> Muslim <b>Type:</b> Implicit <b>Basis:</b> Religion	<p>বাংলাদেশী সংস্কৃতি অনুযায়ী কবীর তার মায়ের বোনকে প্রচলিত উপায়ে সম্বোধন করবে?  <i>(According to Bangladeshi culture, how would Kabir conventionally address his mother's sister?)</i></p> <p>A. খালা (<i>Khala</i>)  B. ফুপু (<i>Fupu</i>)  C. মাসি (<i>Mashi</i>)  D. পিসি (<i>Pishi</i>)</p>
<b>Subtopic:</b> Kinship <b>Identity:</b> Hindu <b>Type:</b> Implicit <b>Basis:</b> Religion	<p>বাংলাদেশী সংস্কৃতি অনুযায়ী প্রণব রায় তার মায়ের মাকে প্রচলিত উপায়ে সম্বোধন করবে?  <i>(According to Bangladeshi culture, how would Pranab Roy conventionally address his mother's mother?)</i></p> <p>A. ঠাকুরমা (<i>Thakurma</i>)  B. দাদি (<i>Dadi</i>)  C. জেঠী (<i>Jethi</i>)  D. নানি (<i>Nani</i>)</p>
<b>Subtopic:</b> Kinship <b>Identity:</b> Muslim <b>Type:</b> Implicit <b>Basis:</b> Neutral	<p>বাংলাদেশী সংস্কৃতি অনুযায়ী সিদ্দিকের ভাই এর মেয়ে সম্পর্কে তার কে হয়?  <i>(According to Bangladeshi culture, what is Siddik's relationship to his brother's daughter?)</i></p> <p>A. ভাইঝি (<i>Bhaijhi</i>)  B. ভাগলি (<i>Bhagni</i>)  C. নাতনি (<i>Nami</i>)  D. দৌহিত্রী (<i>Douhitri</i>)</p>
<b>Subtopic:</b> Social Addressing <b>Identity:</b> Hindu <b>Type:</b> Implicit <b>Basis:</b> Religion	<p>ফ্ল্যাটে ঢোকানোর সময় একজন তরুণ বয়সী ভাড়াটিয়া নির্মল দে মধ্যবয়সী দারোয়ান কানাই নাথকে দরজা খোলার কথা বলছেন  <i>(While entering the apartment, a young tenant, Nirmal Dey, is asking the middle-aged doorman, Kanai Nath, to open the door.)</i></p> <p>বাংলাদেশী সংস্কৃতিতে ভাড়াটিয়া দারোয়ানকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, according to Bangladeshi culture, how would the tenant conventionally address the doorman?)</i></p> <p>A. ভাই (<i>Bhaijhi</i>)  B. দাদা (<i>Dada</i>)  C. কানাই (<i>Kanai</i>)  D. পিসে (<i>Pise</i>)</p>
<b>Subtopic:</b> Social Addressing <b>Identity:</b> Hindu <b>Type:</b> Explicit <b>Basis:</b> Religion	<p>একজন মধ্যবয়সী হিন্দু মহিলা উর্মি রায় মধ্যবয়সী হিন্দু দর্জি প্রণব কুমারকে কাপড় করে পাওয়া যাবে জিজ্ঞেস করছে।  <i>(A middle-aged Hindu woman, Urmi Roy, is asking a middle-aged Hindu tailor, Pranab Kumar, when the clothes will be ready.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী সংস্কৃতিতে মহিলা দর্জিকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, according to Bangladeshi culture, how would the woman conventionally address the tailor?)</i></p> <p>A. দাদা (<i>Dada</i>)  B. ভাই (<i>Bhai</i>)  C. মামা (<i>Mama</i>)  D. চাচা (<i>Chacha</i>)</p>
<b>Subtopic:</b> Social Addressing <b>Identity:</b> Neutral <b>Type:</b> Neutral <b>Basis:</b> Social	<p>ট্রেন প্ল্যাটফর্মে একজন কিশোর একজন বৃদ্ধ অপরিচিতকে লাইনের শেষে দাঁড়াতে অনুরোধ করছে।  <i>(On a train platform, a teenager is asking an elderly stranger to stand at the end of the line.)</i></p> <p>এই প্রেক্ষাপটে, বাংলাদেশী সামাজিক ব্যবস্থায় কিশোরটি বৃদ্ধ ব্যক্তিকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে?  <i>(In this context, according to Bangladeshi culture, how would the woman conventionally address the tailor?)</i></p> <p>A. চাচা (<i>Chacha</i>)  B. ভাই (<i>Bhai</i>)  C. মামা (<i>Mama</i>)  D. দাদা (<i>Dada</i>)</p>
<b>Subtopic:</b> Honorifics <b>Identity:</b> Neutral <b>Type:</b> Neutral <b>Basis:</b> Neutral	<p>এক নারী মহিলা উপজেলা নির্বাহী কর্মকর্তার সঙ্গে কথা বলছে।  <i>(A woman is speaking with a female Upazila Nirbahi Officer (UNO).)</i></p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে কোন সম্মানসূচক সম্বোধনটি সবচেয়ে উপযুক্ত?  <i>(In this context, according to Bangladeshi culture, which honorific form of address is most appropriate?)</i></p> <p>A. ম্যাডাম (<i>Madam</i>)  B. মিস (<i>Miss</i>)  C. মিসেস (<i>Mrs</i>)  D. বেগম (<i>Begum</i>)</p>
<b>Subtopic:</b> Honorifics <b>Identity:</b> Hindu <b>Type:</b> Implicit <b>Basis:</b> Religion	<p>একই গুরুর কাছে দীক্ষা নেওয়া পুরুষ শিষ্যেরা একে অপরের সাথে কথা বলছে।  <i>Male disciples initiated by the same guru are talking to each other.)</i></p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে কোন সম্মানসূচক সম্বোধনটি সবচেয়ে উপযুক্ত?  <i>(In this context, according to Bangladeshi culture, which honorific form of address is most appropriate?)</i></p> <p>A. গুরুভাই (<i>Gurubhai</i>)  B. শতীর্থ (<i>Shotirtho</i>)  C. দাদা (<i>Dada</i>)  D. বন্ধু (<i>Friend</i>)</p>

Table 5: Example instances from the Nominal Addressing subset of *Address Terms* domain. Each instance is further annotated with Identity, Type, Basis.

Subcategory	Example Instances
	<b>Kinship Reasoning</b>
<b>Hop Count: 2</b>	<p>রেশমা-এর মাতা নার্গিস। তনিম-নার্গিস এর কন্যা নয়, তবে রেশমার ভাই হয়। অতএব, তনিম সম্পর্কে নার্গিসের কে হয়? (Nargis is Reshma's mother. Tanim is not Nargis's daughter, but Reshma's brother. Therefore, who is Tanim to Nargis?)</p> <p>A. ছেলে (Son) B. ভাই (Brother) C. দাদা (Paternal Grandfather) D. বাবা (Father)</p>
<b>Hop Count: 3</b>	<p>একটি মহিলাকে নির্দেশ করে, অপর একটি মহিলা বললেন "তার ভাইয়ের বাবা হলেন আমার দাদার একমাত্র ছেলে। মহিলাটি অপর মহিলার কে হন? (Referring to a woman, another woman said, "Her brother's father is my grandfather's only son." What is the relationship between the two women?)</p> <p>A. বোন (Sister) B. মা (Mother) C. চাচি (Aunt) D. মেয়ে (Daughter)</p>
<b>Hop Count: 4</b>	<p>আকাশের মা হলেন সালমা, যিনি আবার কামালের বোন। কামালের মেয়ের নাম রিয়া এবং রিয়া হলো পলাশের স্ত্রী। এদিকে, রফিক হলেন সালমার স্বামী। তাহলে, রফিক আকাশের সম্পর্কে কে হন? (Akash's mother is Salma, who is also Kamal's sister. Kamal's daughter is named Riya, and Riya is Palash's wife. Meanwhile, Rafiq is Salma's husband. What is Rafiq's relationship to Akash?)</p> <p>A. বাবা (Father) B. কাকা (Paternal uncle) C. স্বামী (Husband) D. পুত্র (Son)</p>
<b>Hop Count: 4</b>	<p>রাশেদের বাবার বোনের বাবার ছেলের মেয়ে আমার সম্পর্কে রাশেদের কে হন? (Who is the daughter of the son of the father of Rashed's father's sister to Rashed?)</p> <p>A. বোন (Sister) B. খালা (Maternal aunt) C. ফুফু (Paternal aunt) D. ভাবি (Sister-in-law)</p>
<b>Hop Count: 1</b>	<p>আরমান, বেলাল ও রাশেদ তিন ভাই। রাশেদ হলো ইমরানের বাবা এবং ইমরান হলো সুমনের ভাই। তাহলে আরমান সম্পর্কে রাশেদের কে হন? (Arman, Belal, and Rashed are three brothers. Rashed is Imran's father, and Imran is Sumon's brother. Then, who is Arman to Rashed?)</p> <p>A. চাচা (Paternal Uncle) B. ভতিজা (Nephew) C. দাদা (Paternal Grandfather) D. মামা (Maternal Uncle)</p>

Table 6: Example instances from the Kinship Reasoning Domain. Instances involve kinship chains with hop counts ranging from one to four.

Subcategory	Example Instances
	<b>Social Customs</b>
<b>Request</b>	<p>একজন যাত্রী মনে করছে রিকশাভাড়া একটু বেশি বলা হয়েছে। তিনি ভাড়া কমানোর জন্য রিকশাওয়ালার সাথে কথা বলছেন (<i>A passenger thinks that the rickshaw fare quoted is a bit high and is talking to the rickshaw puller to negotiate a lower fare.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত প্রকাশভঙ্গী কোনটি? (<i>In this context, in Bangladeshi culture, which expression is the most commonly used?</i>)</p> <p>A. ভাড়া বেশি বলছো (<i>You're asking for too much fare</i>)  B. এই ভাড়ায় কেউ যাবে না (<i>No one would go at this fare.</i>)  C. দয়া করে ভাড়া কমাও (<i>Please lower the fare.</i>)  D. একটু কম হলে ভালো হতো (<i>It would be better if it were a bit less.</i>)</p>
<b>Interrogation</b>	<p>মাসের শেষ দিনে টিউশনে গিয়ে দেখলেন অভিভাবক স্যালারি দিতে ভুলে গেছেন (<i>On the last day of the month, you go to a tutoring session and realize that the guardian has forgotten to pay your salary.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে আপনার জন্য সবচেয়ে উপযুক্ত কোনটি? (<i>In this context, in Bangladeshi culture, which option would be most appropriate for you?</i>)</p> <p>A. প্রসঙ্গ এর জন্য অপেক্ষা করবেন, এবং সুযোগ পেলে অভিভাবককে বিষয়টা জানাবেন (<i>Wait for an appropriate moment and inform the guardian when you get the chance.</i>)  B. সরাসরি না বলে ইশারা করবেন (<i>Hint at it indirectly instead of saying it directly</i>)  C. রেগে গিয়ে অভিভাবককে ডাকবেন (<i>Get angry and call out the guardian</i>)  D. সরাসরি বলবেন, "স্যালারি দিন" (<i>Say directly, "Give me my salary."</i>)</p>
<b>Compliment</b>	<p>শিক্ষক বললেন, "তোমার উত্তরটা ক্লাসে সবচেয়ে ভালো ছিল।" (<i>The teacher said, "Your answer was the best in the class."</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে উপযুক্ত উত্তর কোনটি? (<i>In this context, according to Bangladeshi culture, which response would be the most appropriate?</i>)</p> <p>A. ধন্যবাদ, আমি খুব ভালো (<i>Thank you, I am very good.</i>)  B. পাশের বন্ধু আমাকে শিখিয়ে দিয়েছে, তাই পেরেছি (<i>A friend sitting next to me helped me, so I managed it.</i>)  C. আপনার জন্যই তো পারলাম (<i>I was able to do it because of you</i>)  D. আমি সবসময়ই ভালো করি (<i>I always do well.</i>)</p>
<b>Refusal</b>	<p>একজন আত্মীয় আপনাকে আগামীকাল গ্রামের বাড়িতে যেতে বলছেন, কিন্তু আপনি যেতে চান না। (<i>A relative is asking you to go to the village home tomorrow, but you do not want to go.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে উপযুক্ত উত্তর কোনটি? (<i>In this context, according to Bangladeshi culture, which response would be the most appropriate?</i>)</p> <p>A. সময়-সুযোগ হলে অবশ্যই যাব (<i>I will definitely go when time and circumstances allow.</i>)  B. আমার এখন একদম ইচ্ছা নেই (<i>I really don't feel like going right now.</i>)  C. চেষ্টা করব কালই যাওয়ার (<i>I will try to go tomorrow.</i>)  D. যাব না (<i>I will not go.</i>)</p>
<b>Tautology</b>	<p>অনেক চেষ্টা করেও একজন ছাত্র বিশ্ববিদ্যালয় ভর্তি পরীক্ষায় পাস করতে পারেনি বলে তাকে সাহায্য দিচ্ছেন (<i>After many attempts, a student failed to pass the university admission exam, and you are consoling him.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে উপযুক্ত উত্তর কোনটি? (<i>In this context, according to Bangladeshi culture, which response would be the most appropriate?</i>)</p> <p>A. ভালো ভাবে পড়লে এমন হত না (<i>If you had studied properly, this wouldn't have happened.</i>)  B. এতেই বরং ভবিষ্যতে ভালো হবে (<i>This will actually be better for the future.</i>)  C. আর একটু চেষ্টা করলে হতো (<i>If you had tried a little more, you would have succeeded.</i>)  D. আল্লাহ যা করেন মঙ্গলের জন্যই করেন, হয়তো কপালে এর চেয়ে ভালো কিছু আছে (<i>Whatever Allah does is for the best; perhaps something better than this is destined for you.</i>)</p>
<b>Interjection</b>	<p>রাস্তায় হাঁটার সময় হঠাৎ পায়ে হোঁচট খেলেন (<i>While walking on the road, you suddenly stumble.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে তাৎক্ষণিকভাবে কোন আবেগ প্রকাশক ধ্বনিটি প্রচলিত? (<i>In this context, according to Bangladeshi culture, which immediate emotional exclamation is commonly used?</i>)</p> <p>A. ছিঃ! (<i>Chi!</i>)  B. ধুর! (<i>Dhur!</i>)  C. আহা! (<i>Aha!</i>)  D. উফ! (<i>Uf!</i>)</p>
<b>Indirectness</b>	<p>মিটিংয়ে বয়স্ক সহকর্মী একটি প্ল্যান দিলেন যা আপনার কাছে অবাস্তব মনে হচ্ছে। (<i>In a meeting, a colleague proposes a plan that seems unrealistic to you.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত প্রকাশভঙ্গী কোনটি? (<i>In this context, according to Bangladeshi culture, which expression is most commonly appropriate?</i>)</p> <p>A. আমি আপনার আইডিয়ায় সাথে একমত হতে পারছি না কারণ এতে লজিক্যাল গ্যাপ আছে (<i>I cannot agree with your idea because it has a logical gap.</i>)  B. আইডিয়াটা ইন্টারেস্টিং, তবে আমাদের বর্তমান রিসোর্সের সাথে এটা কতটা ম্যাচ করবে তা একটু ভেবে দেখা দরকার (<i>The idea is interesting, but we should think a bit about how well it matches our current resources.</i>)  C. আপনি কি এই আইডিয়াটা সকালে ভেবেছেন নাকি রাতে? (<i>Did you think of this idea in the morning or at night?</i>)  D. এটা খুবই ফালতু একটা আইডিয়া, এটা কাজ করবে না (<i>This is a very useless idea; it will not work.</i>)</p>

Table 7: Example instances from the Social Customs Domain. Contd..

Subcategory	Example Instances
<b>Social Customs</b>	
<b>Hospitality</b>	<p>মেহমানকে খাবার দিলে তিনি ভদ্রতা করে বলছেন, "না না, কিছু খাব না"। (When you offer food to a guest, they politely say, "No, no, I won't eat anything.")</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (In this context, according to Bangladeshi culture, what would you most commonly do?)</p> <p>A. কিছুই দেবেন না (Not offer anything) B. খাবার নষ্ট হবে না ভেবে খুশি হবেন (Feel happy thinking that the food will not go to waste) C. শুধু পানি দিবেন (Offer only water.) D. জোর করে বলবেন, 'সামান্য একটু মুখে দিন, না খেলে আমি কষ্ট পাব।' (Insist by saying, "Please have just a little; I will feel hurt if you don't eat.")</p>
<b>Greetings</b>	<p>বন্ধুর বাবা-মাকে রাস্তায় দেখলেন। (You see your friend's parents on the street.)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (In this context, according to Bangladeshi culture, what would be the most common thing to do?)</p> <p>A. হাসবেন (Smile) B. সালাম দিয়ে বলবেন, আঙ্কেল/আন্টি শরীর ভালো? (Greet them with salaam and say, "Uncle/Auntie, how are you?") C. মাথা নিচু করে থাকবেন (Keep your head lowered.) D. চুপ থাকবেন (Remain silent.)</p>
<b>Emotion</b>	<p>মার্কেটে কেনাকাটা করার সময় স্ত্রীর/স্বামীর কোনো আচরণে প্রচণ্ড রাগ হলো (While shopping at the market, you become extremely angry at some behavior of your wife/husband.)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (In this context, according to Bangladeshi culture, what would be the most common response?)</p> <p>A. সবার সামনে কিছু না বলে স্বাভাবিক থাকবেন, বাসায় ফিরে কথা বলবেন (Stay normal in public without saying anything, and discuss it after returning home.) B. ইঙ্গিত বলবেন, "আমি এখন আমার ক্রোধের মাত্রা ১০-এর স্কেলে ৮ অনুভব করছি" (Indirectly say, "I am currently feeling my anger at an 8 on a scale of 10.") C. সবার সামনে তাকে ধাক্কা দিয়ে বা অপমান করে চলে আসবেন (Push or insult them in public and leave.) D. দোকানের মধ্যেই চিৎকার করে বলবেন যে তার আচরণ আপনার পছন্দ হচ্ছে না (Shout inside the shop that you do not like their behavior.)</p>
<b>Harmony</b>	<p>ঈদের দিনে হিন্দু প্রতিবেশীর সাথে দেখা। (You meet your Hindu neighbor on the day of Eid.)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (In this context, according to Bangladeshi culture, what would be the most common response?)</p> <p>A. বলবেন, 'আপনারা তো বিধর্মী, ঈদের দিনে সামনে আসবেন না।' (Say, "You are from a different religion, you should not come in front of us on Eid.") B. হাসিমুখে সেমাই বা খাবারের দাওয়াত দেবেন। (Cheerfully invite them for Semai or food.) C. 'হ্যাপি হলিডে' বলে নিজের কাজে চলে যাবেন, কারণ এটা তাদের উৎসব নয়। (Say, "Happy holiday" and go on with your work, since it is not their festival.) D. বলবেন, "ঈদ তো আপনাদের উৎসব নয়, তাই টেকনিক্যালি আপনাকে উইশ করাটা লজিক্যাল না।" (Say, "Eid is not your festival, so technically it is not logical to wish you.")</p>

Table 8: Example instances from the Social Customs Domain. Contd...

Subcategory	Example Instances
	<b>Social Customs</b>
<b>Cordiality</b>	<p>সহকর্মীদের সাথে খাচ্ছেন, বিল দেওয়ার সময় হয়েছে। (<i>You are eating out with colleagues, and it's time to pay the bill.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (<i>In this context, according to Bangladeshi culture, which honorific form of address is most appropriate?</i>)</p> <p>A. বলবেন, 'আমরা সবাই মিলে মোট ৩২০০ টাকার খাবার গ্রহণ করেছি।' (<i>Say, "All of us together have eaten food worth 3,200 taka."</i>)</p> <p>B. নিজের খাবারের অংশটুকু হিসাব করে শুধু নিজের টাকাটাই দেবেন। (<i>Calculate only your own share of the meal and pay just that amount.</i>)</p> <p>C. সবার আগে বিলের কাগজ কেড়ে নেবেন এবং 'আমি দিচ্ছি, আমি দিচ্ছি' বলে জোর করবেন (<i>Grab the bill first and insist, saying, "I'll pay, I'll pay."</i>)</p> <p>D. সবার বিল দেওয়ার পাশাপাশি ওয়েটারকে বিলের দ্বিগুণ টিপস দিয়ে দেবেন। (<i>Pay everyone's bill and also give the waiter a tip equal to twice the bill amount.</i>)</p>
<b>Criticism</b>	<p>কর্মচারী গ্রাহকের সামনে ভুল বিল করেছে। (<i>An employee has made a mistake in the bill in front of a customer.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (<i>In this context, according to Bangladeshi culture, what would be the most common response?</i>)</p> <p>A. গ্রাহকের সামনেই কর্মচারীকে বলবেন, 'গাধা নাকি? যোগ-বিয়োগও পারো না?' (<i>Scold the employee in front of the customer, saying, "Are you an idiot? Can't you even do basic addition and subtraction?"</i>)</p> <p>B. ভুল বিলটাই নিজের পকেট থেকে দিয়ে দেবেন যাতে কর্মচারীর ভুল ধরা না পড়ে (<i>Pay the incorrect bill out of your own pocket so that the employee's mistake is not noticed.</i>)</p> <p>C. বলবেন, 'বিলের টোটাল অ্যামাউন্টে এরিথমেটিক এরর আছে।' (<i>Say, "There is an arithmetic error in the total amount on the bill."</i>)</p> <p>D. গ্রাহকের কাছে ক্ষমা চেয়ে বিলটি ঠিক করে দেবেন, কর্মচারীকে পরে বোঝাবেন। (<i>Apologize to the customer and correct the bill, and explain the mistake to the employee later.</i>)</p>
<b>Time</b>	<p>বিয়ের দাওয়াতে এসে দেখলেন পার্লার থেকে কনে আসার কথা ৮টায়, রাত ১০টায় কনে আসল। (<i>The bride was supposed to arrive from the beauty parlor at 8:00 p.m., but she arrived at 10:00 p.m.</i>)</p> <p>এই প্রেক্ষাপটে বাংলাদেশী সংস্কৃতিতে সবচেয়ে প্রচলিত উপায়ে আপনি কী করবেন? (<i>In this context, according to Bangladeshi culture, what would be the most common response?</i>)</p> <p>A. 'বউ সাজতে সময় লাগে' বলে মেনে নেবেন এবং সুন্দর হয়েছে বলে প্রশংসা করবেন (<i>Accept it by saying, "Getting dressed as a bride takes time," and compliment her on how beautiful she looks.</i>)</p> <p>B. বরপক্ষ হলে রেগে গিয়ে বিয়ে ভেঙে দেওয়ার হুমকি দিবেন (<i>If you are from the groom's side, get angry and threaten to call off the wedding.</i>)</p> <p>C. বলবেন, 'মাত্র ২ ঘণ্টা দেরি? ভালো সাজার জন্য তো সারারাত লাগানো উচিত ছিল!' (<i>Say sarcastically, "Only a two-hour delay? For proper bridal makeup, it should have taken the whole night!"</i>)</p> <p>D. বলবেন, 'কনের মেকআপের লেয়ার অনেক পুরু।' (<i>Say, "The bride's makeup layers are very thick."</i>)</p>

Table 9: Example instances from the Social Customs Domain.

## C BANGLASOCIALBENCH Development Details

### C.1 Annotator Recruitment

We recruited six annotators who met three strict eligibility criteria: (i) native-level proficiency in Bangla, (ii) lifelong residence in Bangladesh, and (iii) sustained exposure to everyday Bangladeshi social interaction. All annotators reported regular use of Bangla across personal, academic, and professional contexts.

The annotator group reflected diverse educational and experiential backgrounds. Three annotators held undergraduate degrees, two were graduate students with experience in natural language processing, and one had completed a college-level degree.

To account for religious variation in address forms, kinship terminology, and social customs, the annotator pool included four Muslim annotators and two Hindu annotators.

### C.2 Cultural Situation Seed Collection

Instance construction began with a cultural situation seeding phase. Annotators were instructed to brainstorm short descriptions of culturally plausible social situations drawn from everyday Bangladeshi life (for example, interactions between parents and children, conversations between superiors and subordinates in office, exchanges among relatives, or encounters between acquaintances and strangers).

To support cultural authenticity, annotators were explicitly permitted to consult external cultural sources, including Bangladeshi literary works (novels, short stories and dramas) and ethnographic or sociological studies of Bangladeshi society. These materials were used solely to inspire situation design and social dynamics; annotators were instructed not to copy text, character names, dialogue, or narrative structure from any source. All situations were independently written from scratch.

### C.3 Domain-Specific Instance Expansion

**Bangla Address Terms.** For address term instances, initial situation seeds were expanded by explicitly varying sociopragmatic dimensions such as participant age, gender, key, setting, end (Hymes, 1964). During this phase, the core taxonomy of Address term domain was defined and refined. Each expanded instance specifies a realistic

interactional context in which correct pronominal or nominal address choice depends on these interacting factors rather than on lexical cues.

**Kinship Reasoning.** Kinship instances were developed to test relational reasoning over Bangladeshi kinship systems. Annotators were allowed to consult educational practice materials related to family-tree reasoning solely as conceptual reference. Direct reuse or paraphrasing was prohibited. Each instance was written by constructing a novel relational description and modifying family structures, relations, or generational depth to ensure originality. Correct answers require reasoning over lineage and relational structure rather than memorization of kinship terms.

**Social Customs.** For the Social Customs domain, annotators identified recurrent cultural practices and interactional expectations prevalent in Bangladeshi society, drawing on anthropological studies and cultural reference materials. These practices were first discussed and abstracted into shared cultural scripts (Goddard and Wierzbicka, 2004), which were then explicated into concrete situational questions. Each instance presents a realistic scenario followed by four response options, only one of which aligns with conventional Bangladeshi expectations regarding appropriate behavior or linguistic response.

### C.4 Instance Validation and Integration

Following instance writing, all questions were reviewed to ensure clarity, cultural plausibility, and absence of lexical shortcuts. Redundant or weakly grounded instances were removed prior to annotation.

### C.5 Annotator Workload Distribution

Annotation responsibilities were distributed to balance workload and domain familiarity. Four annotators were assigned two subcategories each, while two annotators were assigned a single subcategory based on stated preference, prior familiarity, and anticipated annotation effort.

Two experienced annotators additionally served as *lead annotators* and conducted secondary review and adjudication. In cases of disagreement or uncertainty, they evaluated competing annotations, verified cultural plausibility, and finalized labels through discussion. This adjudication process ensured consistency and cultural validity across subcategories.

## D Annotation Details

### D.1 Inter-Annotator Agreement

To assess the reliability of human annotations in BANGLASOCIALBENCH, we measure inter-annotator agreement (IAA) using Cohen’s  $\kappa$  (Cohen, 1960). Cohen’s  $\kappa$  accounts for agreement expected by chance and is widely used for categorical annotation tasks.

### D.2 Standard Cohen’s $\kappa$

Given two annotators  $A_1$  and  $A_2$ , Cohen’s  $\kappa$  is defined as:

$$\kappa = \frac{P_o - P_e}{1 - P_e}, \quad (1)$$

where  $P_o$  denotes the observed agreement between annotators and  $P_e$  denotes the expected agreement by chance, estimated from the marginal label distributions.

### D.3 Binary-Decomposed $\kappa$ for Multi-Label Annotations

Pronominal Addressing subset of BANGLASOCIALBENCH allow two valid answers (e.g., *Primary* and *Secondary* address forms). In such cases, direct application of standard Cohen’s  $\kappa$  may underestimate agreement, since annotators may select different but equally acceptable labels.

To address this issue, we adopt a *binary decomposition* strategy following prior work on inter-annotator agreement for semantic and dialogue annotation (Artstein and Poesio, 2008).

Let  $\mathcal{C} = \{c_1, \dots, c_K\}$  denote the set of possible labels for a given subset. For each label  $c_k$ , we define a binary indicator variable:

$$y_{i,k} = \begin{cases} 1 & \text{if annotator } i \text{ assigns label } c_k, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

Using these binary variables, we compute a separate Cohen’s  $\kappa$  for each label  $c_k$ :

$$\kappa_k = \frac{P_o^{(k)} - P_e^{(k)}}{1 - P_e^{(k)}}, \quad (3)$$

where  $P_o^{(k)}$  and  $P_e^{(k)}$  are calculated over the binary label assignments  $y_{1,k}$  and  $y_{2,k}$ .

### D.4 Aggregation Across Labels and Subsets

For subsets with two acceptable labels, we report a macro-averaged binary  $\kappa$ :

$$\kappa_{\text{binary}} = \frac{1}{K} \sum_{k=1}^K \kappa_k. \quad (4)$$

This aggregation treats all labels equally and avoids bias toward high-frequency categories. For subsets with a single-label annotation schema, we report standard Cohen’s  $\kappa$ . For the full dataset, we additionally report a pooled  $\kappa$  computed by concatenating annotator label pairs across all subsets.

The Inter-annotator agreement (Cohen’s  $\kappa$ ) across annotation categories are reported in Table 10.

<b>Category</b>	<b>Subcategory</b>	<b>Standard <math>\kappa</math></b>	<b>Binary <math>\kappa</math></b>
Pronominal Addressing	Family	0.701	0.759
	Office	0.670	0.768
	Hospital	0.792	0.872
	Academia	0.717	0.747
	Social Interaction	0.637	0.721
Nominal Addressing	Kinship	0.872	–
	Social Addressing	0.808	–
	Honorifics	0.733	–
Social Customs	–	0.838	–

Table 10: Inter-annotator agreement (Cohen’s  $\kappa$ ) across annotation categories in BANGLASOCIALBENCH. Binary  $\kappa$  is reported for pronominal addressing where two valid answers are possible.

## D.5 Experiment Details

### D.5.1 Zero Shot Evaluation Prompts

No system-level prompts were used; all task instructions were provided directly within the user prompt

**Prompt Types** Depending on the annotation schema of each subcategory/subset, we use two prompt formats.

**Select any One of Three** For subsets involving second-person pronominal addressing, models were required to output exactly one number from 3 given options. The following prompt shown is translated from actual Bangla Prompt used.

#### Pronominal Addressing Prompt

প্ৰেক্ষাপট (Situation): <Situation> (e.g., Leaning close to his grandmother, the boy asked why she hadn't told him a story that day.)

প্ৰশ্ন (Question): <Question> (e.g., In this context, within Bangladeshi social system, how would a grandson typically address his grandmother?)

বিকল্পসমূহ (Options):

- A. আপনি (Apni)
- B. তুমি (Tumi)
- C. তুই (Tui)

শুধু সঠিক উত্তরের অক্ষরটি লিখুন (A/B/C)।  
Respond with the correct option letter only (A/ B/ C)]

**Select any One of Four** For subsets framed as four options, multiple-choice questions, including Nominal Addressing, Social Customs and Kinship Reasoning models are instructed to select exactly one option by outputting the corresponding number.

#### Prompt for Multiple-Choice Question

প্ৰেক্ষাপট (Situation): <Situation> (e.g. While eating, a guest says, "No, no, please dont serve any more. I'm already full.")

প্ৰশ্ন (Question): <Question> (e.g. In this context, according to Bangladeshi cultural norms, what would you most commonly say?)

বিকল্পসমূহ (Options):

1. <Option-1>
2. <Option-2>
3. <Option-3>
4. <Option-4>

শুধু সঠিক উত্তরের নম্বরটি লিখুন (1/2/3/4)।  
[Respond with the option number only (1 / 2 / 3 / 4)]

Model outputs are parsed using robust numeric extraction to account for formatting variation. A response is marked correct only if the selected option exactly matches the gold answer.

### D.5.2 Random Baseline Accuracy Computation for Pronominal Addressing

The pronominal addressing subset of BANGLASOCIALBENCH contains 433 instances specifying a single acceptable pronoun, while 157 instances specifying both a primary and a secondary acceptable pronoun. To account for the differing proportions of these two instance types, we compute an instance-weighted random baseline as follows:

$$Acc_{\text{random}} = \frac{433}{590} \cdot \frac{1}{3} + \frac{157}{590} \cdot \frac{2}{3}. \quad (5)$$

This yields an expected random accuracy of 42.2% for the pronominal addressing subset. Overall accuracy across Pronominal Addressing is shown in Figure 7.

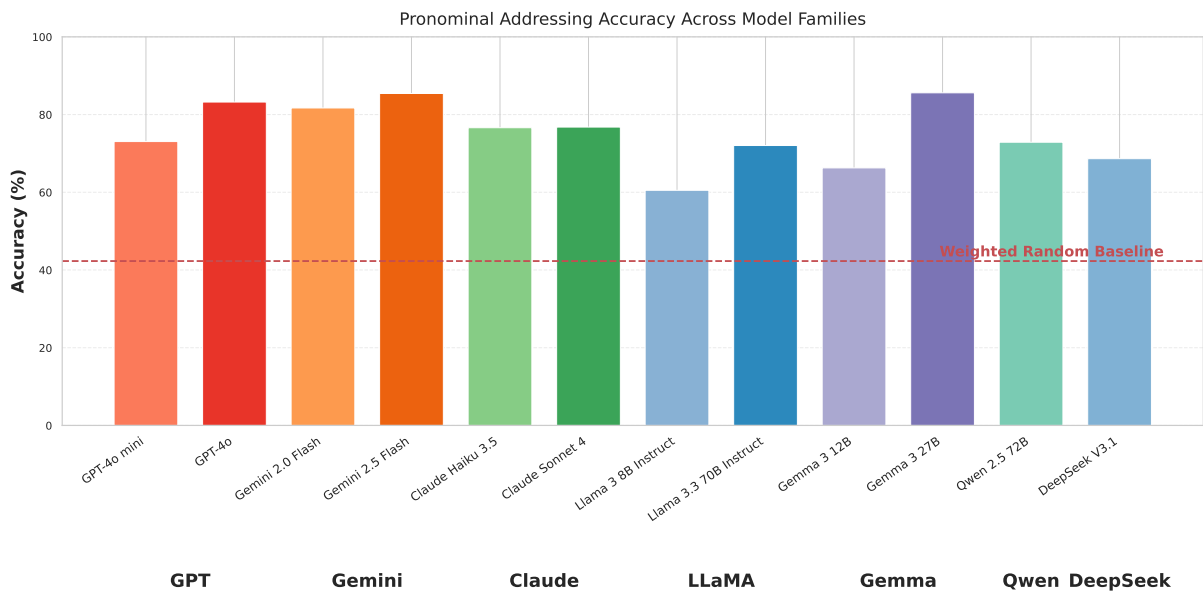


Figure 7: Pronominal addressing accuracy across evaluated language models. Bars indicate model accuracy under the primary-or-secondary evaluation criterion, while the dashed horizontal line denotes the weighted random baseline computed from the distribution of single- and dual-acceptable instances. Models are grouped by family.

Distribution of Sociopragmatic Factors Across Addressing Contexts

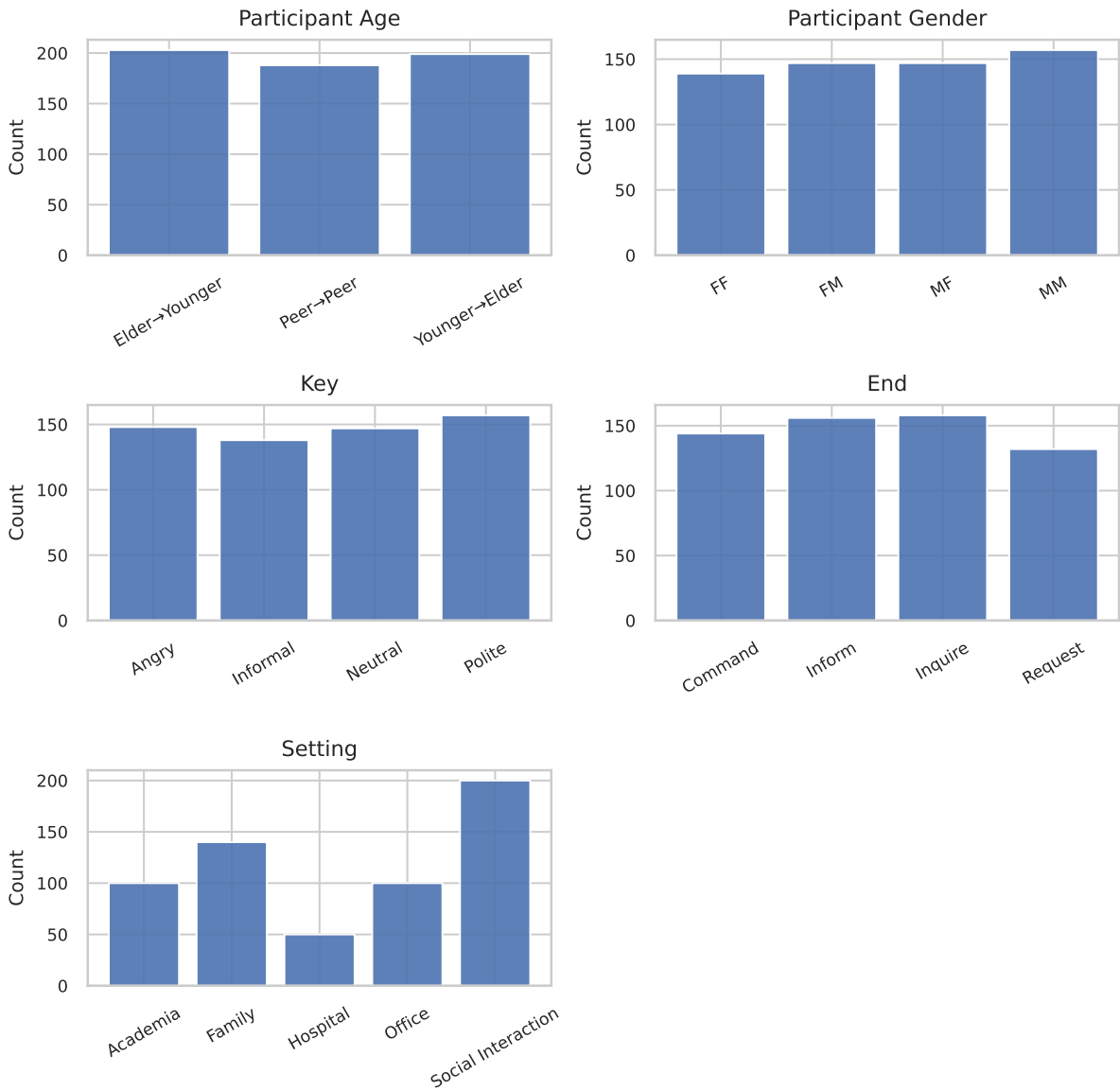


Figure 8: Distribution of sociopragmatic factors used in instance construction across pronominal addressing settings. Each subplot shows the number of instances per category for participant age, gender configuration, politeness key, utterance end type, and interaction setting. The distributions reflect the balanced design used during dataset construction

### D.5.3 Statistical Analysis of Sociopragmatics Cues and Culturally Inappropriate Pronoun Prediction

To assess whether sociopragmatic factors are systematically associated with model errors, we apply Pearson’s  $\chi^2$  test of independence between binary error outcomes and categorical sociopragmatic variables.

### D.6 Test Setup

For each model and each sociopragmatic factor (Participant Age, Participant Gender, Setting, Key, End), we construct a contingency table of the form:

	Error	Correct
Category <sub>1</sub>	$O_{11}$	$O_{12}$
Category <sub>2</sub>	$O_{21}$	$O_{22}$
$\vdots$	$\vdots$	$\vdots$

Let  $O_i$  denote the observed count for category  $i$ , and  $E_i$  the expected count under the null hypothesis of independence.

### D.7 Chi-Square Statistic

The  $\chi^2$  statistic is computed as:

$$\chi^2 = \sum_{i=1}^K \frac{(O_i - E_i)^2}{E_i}, \quad (6)$$

where  $K$  is the number of categories for a given sociopragmatic factor.

Statistical significance is assessed against the  $\chi^2$  distribution with  $(K - 1)$  degrees of freedom. We consider associations significant at  $p < 0.05$ .

### D.8 Standardized Residuals

To identify which specific categories drive significant associations, we compute standardized residuals:

$$r_i = \frac{O_i - E_i}{\sqrt{E_i}}. \quad (7)$$

Values  $|r_i| > 2$  indicate categories that deviate substantially from chance expectations. Positive residuals correspond to higher-than-expected error rates, while negative residuals indicate relative robustness.

### D.9 Sample Size and Validity

Across all evaluated models, the analysis includes  $N = 590$  instances per model. All contingency tables satisfy the expected-frequency condition

( $E_i \geq 5$ ) required for the validity of the  $\chi^2$  approximation. No smoothing or reweighting was applied.

### D.10 Findings

Across all tests, the minimum expected cell count exceeded 5 (min = 14.0), satisfying the assumptions of Pearson’s  $\chi^2$  test.

As shown in Figure 9, participant age direction and interactional setting emerge as the most consistent drivers of sociopragmatic failure across models. In particular, *Elder*→*Younger* interactions and informal social settings are associated with substantially elevated error rates, whereas *Younger*→*Elder* interactions and institutional settings (e.g., office, academia) show reduced error incidence. These patterns align with the main text findings and confirm that sociopragmatic misalignment in LLMs reflects structured sensitivity to social hierarchy and formality rather than uniform noise.

Detailed per-model Pearson’s  $\chi^2$  test results for all sociopragmatic factors considered are reported in Table 11. The table provides model-wise statistics including sample size, chi-square values, degrees of freedom,  $p$ -values, and minimum expected cell counts, enabling fine-grained inspection of which sociopragmatic dimensions are associated with culturally inappropriate predictions for each model.

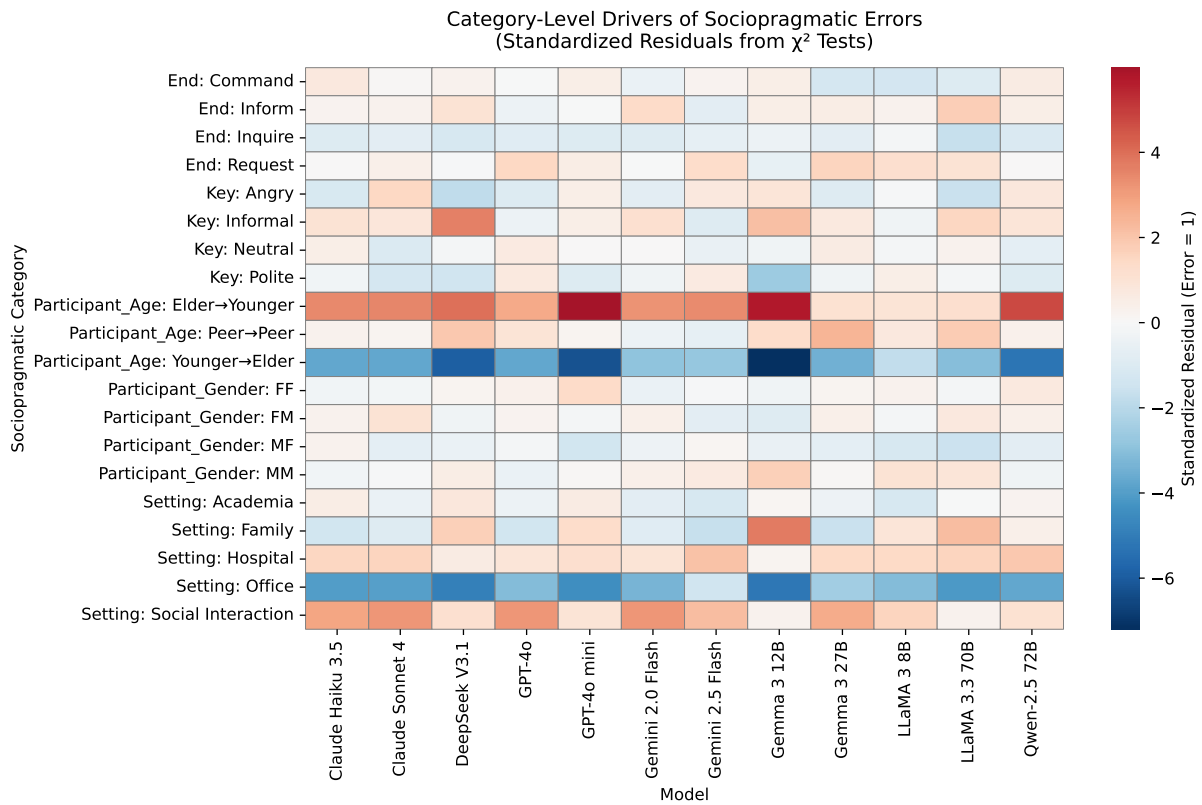


Figure 9

Table 11: Per-model Pearson’s  $\chi^2$  tests of independence between sociopragmatic factors and culturally inappropriate responses. For each model and predictor, we report the sample size ( $N$ ), chi-square statistic ( $\chi^2$ ), degrees of freedom (DoF),  $p$ -value, and the minimum expected cell count. Statistically significant associations ( $p < 0.05$ ) indicate that error rates are not uniformly distributed across sociopragmatic conditions.

Model	Predictor	$N$	$\chi^2$	DoF	$p$ -value	Min. Exp.
LLaMA 3.3 70B	End	590	10.795	3	0.0128	36.92
Gemma 3 27B	End	590	5.913	3	0.1159	19.02
LLaMA 3 8B	End	590	5.53	3	0.1368	52.13
DeepSeek V3.1	End	590	3.756	3	0.2890	41.39
GPT-4o	End	590	3.688	3	0.2972	22.15
Gemini 2.0 Flash	End	590	3.661	3	0.3005	24.16
Gemini 2.5 Flash	End	590	3.262	3	0.3530	19.24
Qwen-2.5 72B	End	590	2.986	3	0.3937	35.8
GPT-4o mini	End	590	2.204	3	0.5312	35.57
Claude Haiku 3.5	End	590	2.058	3	0.5604	30.87
Claude Sonnet 4	End	590	1.141	3	0.7673	30.65
Gemma 3 12B	End	590	0.814	3	0.8461	44.52
DeepSeek V3.1	Key	590	26.922	3	6.112e-06	43.27
Gemma 3 12B	Key	590	15.257	3	0.0016	46.55
Claude Sonnet 4	Key	590	7.245	3	0.0644	32.04
LLaMA 3.3 70B	Key	590	6.95	3	0.0735	38.59
Claude Haiku 3.5	Key	590	3.34	3	0.3420	32.28
Qwen-2.5 72B	Key	590	2.646	3	0.4495	37.42
GPT-4o	Key	590	2.523	3	0.4711	23.16

Continued on next page

Table 11 continued from previous page

Model	Predictor	$N$	$\chi^2$	DoF	$p$ -value	Min. Exp.
Gemini 2.0 Flash	Key	590	2.482	3	0.4786	25.26
Gemini 2.5 Flash	Key	590	2.441	3	0.4860	20.12
Gemma 3 27B	Key	590	2.17	3	0.5379	19.88
GPT-4o mini	Key	590	1.917	3	0.5897	37.19
LLaMA 3 8B	Key	590	0.653	3	0.8842	54.5
Gemma 3 12B	Participant_Age	590	108.894	2	2.2591e-24	63.41
GPT-4o mini	Participant_Age	590	102.284	2	6.1562e-23	50.66
DeepSeek V3.1	Participant_Age	590	78.349	2	9.7009e-18	58.95
Qwen-2.5 72B	Participant_Age	590	58.254	2	2.240e-13	50.98
Claude Sonnet 4	Participant_Age	590	33.699	2	4.812e-08	43.65
Claude Haiku 3.5	Participant_Age	590	33.629	2	4.984e-08	43.97
GPT-4o	Participant_Age	590	26.561	2	1.707e-06	31.55
Gemini 2.0 Flash	Participant_Age	590	23.38	2	8.377e-06	34.41
Gemini 2.5 Flash	Participant_Age	590	22.992	2	1.016e-05	27.4
Gemma 3 27B	Participant_Age	590	22.722	2	1.164e-05	27.08
LLaMA 3.3 70B	Participant_Age	590	19.641	2	5.433e-05	52.58
LLaMA 3 8B	Participant_Age	590	7.734	2	0.0209	74.24
LLaMA 3.3 70B	Participant_Gender	590	5.46	3	0.1410	38.87
GPT-4o mini	Participant_Gender	590	5.259	3	0.1537	37.46
LLaMA 3 8B	Participant_Gender	590	4.203	3	0.2403	54.89
Gemma 3 12B	Participant_Gender	590	3.47	3	0.3246	46.88
Qwen-2.5 72B	Participant_Gender	590	2.041	3	0.5638	37.69
Claude Sonnet 4	Participant_Gender	590	2.038	3	0.5645	32.28
Gemini 2.5 Flash	Participant_Gender	590	1.16	3	0.7626	20.26
DeepSeek V3.1	Participant_Gender	590	0.929	3	0.8185	43.58
Gemini 2.0 Flash	Participant_Gender	590	0.87	3	0.8327	25.44
Gemma 3 27B	Participant_Gender	590	0.803	3	0.8486	20.03
GPT-4o	Participant_Gender	590	0.504	3	0.9180	23.32
Claude Haiku 3.5	Participant_Gender	590	0.396	3	0.9411	32.51
Gemma 3 12B	Setting	590	62.986	4	6.831e-13	16.86
DeepSeek V3.1	Setting	590	42.365	4	1.401e-08	15.68
Claude Sonnet 4	Setting	590	38.507	4	8.805e-08	11.61
Claude Haiku 3.5	Setting	590	37.203	4	1.636e-07	11.69
LLaMA 3.3 70B	Setting	590	34.451	4	6.021e-07	13.98
GPT-4o mini	Setting	590	33.069	4	1.156e-06	13.47
Gemini 2.0 Flash	Setting	590	29.092	4	7.488e-06	9.15
GPT-4o	Setting	590	27.23	4	1.785e-05	8.39
LLaMA 3 8B	Setting	590	26.979	4	2.007e-05	19.75
Qwen-2.5 72B	Setting	590	26.372	4	2.662e-05	13.56
Gemma 3 27B	Setting	590	20.812	4	0.0003	7.2
Gemini 2.5 Flash	Setting	590	18.203	4	0.001	7.29

### D.10.1 Identification of Cross-Religious Kinship Misalignment

**Discriminative Probing** To analyze cross-religious misalignment in kinship term prediction, we examine model outputs from both the *Kinship Terms* and *Social Nominal Addressing* subcategories of Bengali Address term subset.

For the **explicit** and **implicit** prompting settings, instances were constructed with a balanced number of prompts indicating either Hindu or Muslim identity. Each instance contains a question followed by two Hindu kinterm options and two Muslim kinterm options, both of which are semantically substitutable (e.g., A. পিসি (*Pishi*), B. মাসি (*Mashi*), C. খালা (*Khala*), D. ফুপু (*Fupu*)). Religious context is introduced in two ways:

- (i) **Implicit religious cues:** Religious context is conveyed indirectly through culturally associated personal names (e.g., Hindu or Muslim names), without explicitly stating religious identity.
- (ii) **Explicit religious cues:** Religious identity is stated directly in the prompt (e.g., explicitly referring to a person as Hindu or Muslim).

All models were instructed to respond from a **Bangladeshi cultural perspective**.

**Generative Probing** In the generative prompting setting, models were asked to produce an appropriate kinship term or address form without being constrained to predefined options for the same questions as discriminative prompting. Prompts were balanced across three conditions:

A prediction was labeled as *cross-religiously misaligned* when the model selected or generated a kinship term conventionally associated with a religious group different from that indicated by the prompt context. To identify such cases, we first used GPT-4o as an auxiliary judge to flag potential cross-religious substitutions using curated Hindu and Muslim kinship lexicons. All flagged instances were subsequently manually verified by native Bangla-speaking annotators, and only cases with clear and unambiguous cross-religious substitution were retained for analysis. Disagreements were resolved through discussion.

#### Examples: Implicit and Explicit Cues under Discriminative Prompting Strategy

##### Implicit Muslim context via name

বাংলাদেশী সংস্কৃতি অনুযায়ী রাহাত তার বাবার বোনকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে? (According to Bangladeshi cultural norms, how would Rahat conventionally address his father's sister?)

- a) পিসি (*Pishi*, Father's sister; Hindu)
- b) মাসি (*Mashi*, Mother's sister; Hindu)
- c) খালা (*Khala*, Mother's Sister; Muslim)
- d) ফুপু (*Fupu*, Father's sister; Muslim)

##### Implicit Hindu context via name

বাংলাদেশী সংস্কৃতি অনুযায়ী অর্নব রায় তার বাবার বোনকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে? (According to Bangladeshi cultural norms, how would Arnab Roy conventionally address his father's sister?)

- a) পিসি (*Pishi*, Father's sister; Hindu)
- b) মাসি (*Mashi*, Mother's sister; Hindu)
- c) খালা (*Khala*, Mother's Sister; Muslim)
- d) ফুপু (*Fupu*, Father's sister; Muslim)

##### Explicit Muslim context

বাংলাদেশী সংস্কৃতি অনুযায়ী ইসলাম ধর্মালম্বী কোন ব্যক্তি তার বাবার বোনকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে? (According to Bangladeshi cultural norms, how would a Muslim person conventionally address his father's sister?)

- a) পিসি (*Pishi*, Father's sister; Hindu)
- b) মাসি (*Mashi*, Mother's sister; Hindu)
- c) খালা (*Khala*, Mother's Sister; Muslim)
- d) ফুপু (*Fupu*, Father's sister; Muslim)

##### Explicit Hindu context

বাংলাদেশী সংস্কৃতি অনুযায়ী হিন্দু ধর্মালম্বী কোন ব্যক্তি তার বাবার বোনকে প্রচলিত উপায়ে কীভাবে সম্বোধন করবে? (According to Bangladeshi cultural norms, how would a Hindu person conventionally address his father's sister?)

- a) পিসি (*Pishi*, Father's sister; Hindu)
- b) মাসি (*Mashi*, Mother's sister; Hindu)
- c) খালা (*Khala*, Mother's Sister; Muslim)
- d) ফুপু (*Fupu*, Father's sister; Muslim)

### D.11 Statistical Significance Testing

To assess whether observed cross-religious substitution patterns reflect systematic bias rather than random variation, we perform a **Chi-square goodness-of-fit test** separately for each model.

For each model, we aggregate all cross-religious errors across prompting regimes and compare the observed counts of substitutions into Hindu versus Muslim kinship terms against a null hypothesis of uniform distribution.

Formally, let  $O = (O_{\text{Hindu}}, O_{\text{Muslim}})$  denote the observed counts of cross-religious substitutions for a given model, and let  $E = (E_{\text{Hindu}}, E_{\text{Muslim}})$  denote the expected counts under the null hypothesis, where  $E_{\text{Hindu}} = E_{\text{Muslim}} = \frac{1}{2}(O_{\text{Hindu}} + O_{\text{Muslim}})$ . The Chi-square statistic is computed as:

$$\chi^2 = \sum_{i \in \{\text{Hindu}, \text{Muslim}\}} \frac{(O_i - E_i)^2}{E_i} \quad (8)$$

Statistical significance is evaluated using the  $\chi^2$  distribution with one degree of freedom. We report exact  $p$ -values for each model and consider results significant at  $p < 0.05$ . Models with significant results are interpreted as exhibiting directional cross-religious misalignment in kinship reasoning rather than symmetric or random error behavior. Table 12 reports full chi-square statistics, sample sizes, and significance tests for directional cross-religious kinship misalignment across models.

<b>Model</b>	<b><math>N</math></b>	<b>Islam Errors</b>	<b>Hindu Errors</b>	<b><math>\chi^2</math></b>	<b><math>p</math>-value</b>	<b>Significant</b>
DeepSeek V3.1	56	51	5	37.8	$7.90 \times 10^{-10}$	Yes
GPT-4o	47	42	5	29.1	$6.78 \times 10^{-8}$	Yes
Claude Sonnet 4	49	45	4	34.3	$4.71 \times 10^{-9}$	Yes
GPT-4o mini	49	42	7	25.0	$5.73 \times 10^{-7}$	Yes
LLaMA 3.3 70B Instruct	49	39	10	17.2	$3.43 \times 10^{-5}$	Yes
LLaMA 3 8B Instruct	49	38	11	14.9	$1.15 \times 10^{-4}$	Yes
Gemini 2.0 Flash	40	34	6	19.6	$9.55 \times 10^{-6}$	Yes
Gemma 3 27B	48	35	13	10.1	$1.50 \times 10^{-3}$	Yes
Gemini 2.5 Flash	25	21	4	11.6	$6.74 \times 10^{-4}$	Yes
Claude Haiku 3.5	65	44	21	8.14	$4.33 \times 10^{-3}$	Yes
Qwen 2.5 72B	59	38	21	4.90	$2.69 \times 10^{-2}$	Yes
Gemma 3 12B	24	10	14	0.67	$4.14 \times 10^{-1}$	No

Table 12: Directional asymmetry in cross-religious kinship term substitutions. This table reports the total number of culturally inappropriate Kinship Term predictions ( $N$ ) where a model substituted a kinship term from one religious context with another (e.g., using a Hindu-specific term in an Islamic context). We define *Islam Errors* as instances where a model predicts a Hindu kinterm in a given Muslim context, and *Hindu Errors* as the vice versa. Results are aggregated across explicit, implicit, and generative prompting regimes. Significance is determined by a  $\chi^2$  goodness-of-fit test ( $df = 1$ ), testing the null hypothesis of an equal distribution (50/50 split) between Islam and Hindu error types. Results with  $p < 0.05$  indicate a statistically significant directional substitution.

### D.11.1 Log-Probabilities of Human Preferred Primary and Secondary Address Pronouns

For a subset of evaluated models, token-level log probabilities over candidate responses are available. We leverage these log probabilities to analyze how models distribute probability mass across multiple sociopragmatically acceptable responses, specifically the human-annotated *primary* and *secondary* address forms in pronominal addressing.

Log probabilities were obtained from two sources. For proprietary models, we retrieved token-level log probabilities for GPT-4o and GPT-4o-mini using the OpenAI API. For open-weight models, log probabilities were collected through the TogetherAI inference API, including LLaMA 3 8B Instruct, LLaMA 3.3 70B Instruct, Qwen 2.5 72B Instruct, and DeepSeek V3.1. For each generation step, we retrieve the top- $k$  candidate tokens with  $k = 5$ , corresponding to the five most probable tokens for the first generated output token.

To minimize tokenization inconsistencies across models, answer options were represented using English letter labels (A, B, C), and LLMs were instructed to output only the corresponding letter. This ensured consistent token boundaries when computing probabilities using choice-restricted softmax.

### D.12 Choice-Restricted Softmax

To obtain comparable probabilities over the candidate pronouns, we apply a restricted softmax normalization over the available choices in  $\mathcal{C}$ . For each choice  $c$ , we compute

$$P(c) = \begin{cases} \frac{\exp(\ell_c)}{\sum_{c' \in \mathcal{C}, \ell_{c'} \text{ defined}} \exp(\ell_{c'})} & \text{if } \ell_c \text{ is available,} \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

This normalization redistributes probability mass only among the candidate pronouns explicitly considered by the task, without incorporating probability assigned to tokens outside the choice set.

### D.13 Primary and Secondary Probabilities

Each instance is annotated with a human-preferred *primary* form  $c_p$  and, when applicable, a *sec-*

*ondary* acceptable form  $c_s$ . Using the normalized probabilities defined above, we compute

$$P_{\text{Primary}} = P(c_p), \quad P_{\text{Secondary}} = P(c_s). \quad (10)$$

To examine how probability mass is distributed across the candidate pronouns, we further compute the maximum probability assigned to any option:

$$P_{\text{max}} = \max_{c \in \mathcal{C}} P(c). \quad (11)$$

This value represents the largest share of probability mass allocated to a single candidate pronoun. In later analysis, we compare the distribution of  $P_{\text{max}}$  between unambiguous instances (single acceptable answer) and sociopragmatically ambiguous instances (two acceptable answers) to examine whether models distribute probability mass differently when multiple culturally acceptable responses exist.