

Intrinsic Linguistic Bias in Formal vs. Informal Bengali Pragmatics with Progressive Context Inflation

Md Tanzib Hosain^{1,3}, Md Kishor Morol^{2,3}

¹American International University-Bangladesh, ²Cornell University

³ELITE Research Lab

20-42737-1@student.aiub.edu, mmorol@cornell.edu

Abstract

The social biases inherent in language models necessitate a critical analysis of their social influence in many linguistic situations because of their extensive use. This study investigates gender bias in Bengali language models by highlighting the unique linguistic challenges posed by its complex morphology, dialectal variations, and distinctions between formal and informal language versions. While prior research on social bias in Bengali has provided foundational insights, it has not adequately addressed the nuances arising from these variations. This research extends to measuring intrinsic gender bias in both formal and informal Bengali, analyzing the impact of context lengths on bias detection, and proposing modifications to existing techniques to enhance their applicability to Bengali. Addressing these, the study aims to contribute to developing more inclusive and representative bias measurement methodologies for underrepresented languages. We open the source code and data at <https://github.com/kraritt/b-bias-ctxt>.

1 Introduction

Contextualized and context-free language models have both shown a growing number of human-like biases (Yu et al., 2024; Salewski et al., 2024; Caliskan et al., 2017; Bolukbasi et al., 2016). More complex techniques for bias identification have become required in tandem with the advent of innovative ideas, such as more complex language models (Wan et al., 2023; Esiobu et al., 2023; Lin and Ng, 2023; Koo et al., 2023; Guo and Caliskan, 2021; Kurita et al., 2019; May et al., 2019). Sentence-level bias detection techniques have consequently been developed. However, there has been little study on bias detection techniques in other languages, primarily focusing on English (Nangia et al., 2020) (Bolukbasi et al., 2016). Recent initiatives focus on detecting bias in German (Kurpicz-Briki, 2020), French (Kurpicz-Briki, 2020), Chi-

nese (Liang et al., 2020), Arabic (Lauscher et al., 2020) and Dutch (Mulsa and Spanakis, 2020). A thorough test of bias associated with binary gender (Pujari et al., 2019) and societal (Malik et al., 2021) in the Hindi language was carried out.

Although previously a kind of bias test has been done by (Sadhu et al., 2024) for the Bengali language, it is not sufficient to precisely point out bias in Bengali without analyzing the context in both formal and informal Bengali because of its complex linguistic structures, morphology, and dialectal variations in formal and informal versions. The distinctions are prominently marked by changes in gender usage through levels of pragmatics. Formal Bengali, used in written communications, official speeches, and media broadcasts, often adheres to a more standardized vocabulary and avoids colloquialisms. Informal Bengali, prevalent in everyday conversation, may include a variety of dialects, which are more expressive and personal. Although the gender neutrality of pronouns remains constant across these styles, the context in which gendered nouns are used can vary, reflecting the speaker's social and cultural nuances. For instance, it might be more common in formal settings to use titles and honorifics that specify gender, which adds a layer of respect or formality. In contrast, informal interactions might skip such formalities altogether.

This work addresses linguistic bias in Bengali by focusing on pragmatic variations between formal and informal registers—a dimension largely unexplored in bias studies for languages without grammatically gendered pronouns (e.g., Finnish, Turkish). While gender-neutral pronoun languages typically investigate semantic associations, Bengali's complex morphology, dialectal diversity, and register-specific gendered noun usage (e.g., honorifics in formal contexts vs. colloquialisms in informal) necessitate tailored methodologies. The study uniquely adapts bias tests (WEAT/SEAT/CEAT) to account for context span

variance across registers, revealing that informal Bengali exhibits slightly higher bias and that optimal context lengths (≈ 25 words) stabilize measurements—challenges absent in grammatically genderless languages where bias manifests primarily through lexical or contextual associations without register-based divergence. The study’s core innovations—probing context-length effects, creating parallel formal/informal datasets, and modifying bias metrics for morphological richness—are explicitly proposed as generalizable to other low-resource languages with similar linguistic complexities, advancing inclusive bias evaluation beyond English-centric approaches.

This study aims to address this constraint by focusing on gender bias and introducing formal and informal Bengali into the field of bias analysis. Throughout this work, we provide (i) an investigation for measuring intrinsic gender bias in both formal and informal versions, specifically with the development of a dataset, (ii) an analysis of how different context lengths affect bias measurement techniques, and (iii) discussions on the modifications required to apply current bias measurement techniques.

2 Experiment Methods

As a language, Bengali exhibits unique linguistic traits regarding gender representation, particularly in its use of pronouns and nouns. Unlike English, which employs gender-specific pronouns such as "he" and "she" to distinguish between male and female subjects, Bengali utilizes a gender-neutral pronoun "se" for both male and female entities. This characteristic facilitates a less gendered discourse in informal and formal communication, potentially reducing gender bias in language usage. However, Bengali does not entirely eschew gender distinctions; it mirrors English in its use of gendered common nouns for human referents, such as "chele" for boy and "meye" for girl, or "purusa" for man and "nari" for woman. For that, we employ ordinary nouns rather than pronouns when it’s crucial to mask the gendered term in a phrase. We have tested two distinct methods, embedding extraction and mask prediction, for measuring inherent bias for comparison in contextual contexts in this study.

For pragmatic bias test, our multi-stage methodological process initiates with separating the data based on different language variants. Subsequently, the separated language data undergoes

morphological processing. After this stage, the data enters a structured data preparation phase, which branches into two distinct paths: dataset creation and morphology handling. In the dataset creation pathway, sentence extraction and context length variation steps are undertaken to structure and diversify the data. Concurrently, the morphology handling path prepares data specifically for morphological evaluation. The processed datasets then flow into the analysis stage involving various statistical and semantic testing methodologies, including the embedding extraction methods: Word Embedding Association Test (WEAT), Sentence Embedding Association Test (SEAT), Contextual Embedding Association Test (CEAT), and the mask prediction methods: Probabilistic Logarithmic Bias Test.

2.1 Embedding Extraction

2.1.1 WEAT and SEAT

We use two well-known techniques based on the extraction of the embedding methodology for evaluating bias as our initial baselines: Word Embedding Association Test (WEAT) (Caliskan et al., 2017) and Sentence Encoder Embedding Association Test (SEAT) (May et al., 2019). WEAT is intended to indicate how strongly two word vectors associate statistically. By modifying the original dataset to suit the Bengali context, we create a dataset¹ especially for Bengali to carry out this experiment. As seen in Table 1, we employ separate sets of Target vs Attribute word categories. Word2vec (Mikolov, 2013) and GloVe (Pennington et al., 2014), two static word embedding models, are trained using Bangla2B+ (Bhattacharjee et al., 2021) to extract the associated embedding vectors. After that, we determine statistical significance by computing effect sizes, Cohen’s d , and associated p -values, with a significance threshold of $p < 0.07$. Cohen’s effect size metric, d indicates that medium effect sizes are denoted by $d > |0.5|$ and high effect sizes by $d > |0.8|$ (Rice and Harris, 2005).

The following equation presents the effect size.

$$d = \frac{\mu_{x \in X} s(x, A, B) - \mu_{y \in Y} s(y, A, B)}{\sigma_{w \in X \cup Y} s(w, A, B)}$$

¹We use the original WEAT categories, translating some words both in formal and informal versions verbatim and adapting others culturally

$$\mu_{x \in X} s(x, A, B) = \frac{1}{|X|} \sum_{x \in X} s(x, A, B)$$

$$\sigma_{w \in X \cup Y} s(w, A, B) = \sqrt{\frac{1}{|X \cup Y|} \sum_{w \in X \cup Y} (s(w, A, B) - \mu_{w \in X \cup Y} s(w, A, B))^2}$$

$$s(x, A, B) = \frac{1}{|A|} \sum_{a \in A} \cos(x, a) - \frac{1}{|B|} \sum_{b \in B} \cos(x, b)$$

The SEAT experiment allows for evaluating contemporary contextual embedding systems for bias by extending WEAT to apply to sentence embeddings. For the SEAT experiment, we utilize template phrases with Target vs Attribute terms from Table 1 for every category. The SEAT template sentences include every word of the target from the list of targets in WEAT. We employ the Bengali translation² of (May et al., 2019) semantically bleached templates. We employ BanglaBERT’s (Bhattacharjee et al., 2021) last layer to extract embeddings for every sentence. The effect size of the selected list of sentences based on the categories above is then determined using these embeddings.

2.1.2 CEAT

We test the Contextualized Embedding Association Test (CEAT) (Guo and Caliskan, 2021), an extension of WEAT, which measures intrinsic biases in Contextual Word Embeddings (CWE) by generating a representation of random effects (Hedges, 1983) in the effect size distribution on its input context variation. The random-effects model calculates the weighted mean of effect sizes and their statistical significances as a bias measure. A statistical model with random variables as its model parameters is called a random effects model. The concept is predicated on the idea that the data under analysis are taken from a hierarchy of distinct populations, the differences of which are related to that hierarchy. The contextualized variation effect size variations between two sets of target words based on their relative similarity to two sets of attribute words are assumed to be explained by a random variable uncorrelated with the independent variables in our CEAT computation.

The following formula represents the effect size, d_i , for the i^{th} sample.

$$d_i = \frac{\mu_{x \in X} s(x, A, B) - \mu_{y \in Y} s(y, A, B)}{\sigma_{w \in X \cup Y} s(w, A, B)}$$

²We use the Google Translate API and Bangla Academy Adhunik Bangla Abhidhan to do translations <https://translate.google.com/> <https://banglaacademy.gov.bd/>

The square of the $\sigma_{w \in X \cup Y} s(w, A, B)$ is the in-sample variance estimation, represented by V_i . The ANOVA technique estimates the between-sample variance, σ_b^2 . The following formula represents it.

$$\sigma_b^2 = \begin{cases} 0 & \text{if } Q < N - 1 \\ \frac{Q - (N - 1)}{c} & \text{if } Q \geq N - 1 \end{cases}$$

$$W_i = \frac{1}{v_i},$$

$$c = \sum W_i - \frac{\sum W_i^2}{\sum W_i},$$

$$Q = \sum W_i d_i^2 - \frac{\sum (W_i d_i)^2}{\sum W_i}.$$

Each effect size’s allocated weight for calculating the combined effect size, θ , is known as the weight v_i .

$$v_i = \frac{1}{V_i + \sigma_b^2}$$

$$\theta(X, Y, A, B) = \frac{\sum_{i=1}^N v_i d_i}{\sum_{i=1}^N v_i}$$

The hypothesis test is derived by computing the standard error, $\sigma_{\bar{x}}$ of θ . The following formula presents the standard error.

$$\sigma_{\bar{x}}(\theta) = \sqrt{\frac{1}{\sum_{i=1}^N v_i}}$$

The standard normal distribution is the limiting version of the $\frac{\theta}{\sigma_{\bar{x}}(\theta)}$ distribution according to the central limit theorem (Montgomery and Runger, 2010). We utilize a two-tailed p – value, which may assess the bias significance in two directions, because we discovered that some of the θ values are negative. The following formula produces the two-tailed p – value for the significance test, which supports the hypothesis that there is no difference between all contextualized versions of the two sets of target words in terms of their relative similarity to the two sets of attribute words.

$$P_c(X, Y, A, B) = 2 \times \{1 - \phi(|\frac{\theta}{\sigma_{\bar{x}}(\theta)}|)\}$$

From n_s extracted phrases for each stimulus s , we produce n_s CWE for a specific segment length l . We do this for certain phrase lengths, which we call segments. We randomly select each stimulus N times for every segment length l . We sample using replacement to maintain the distribution if the stimulus occurs in fewer than N sentences.

Table 1: WEAT stimuli list (F points Formal and I points Informal words, whereas (N) points names and (W) points words).

Category	Target	Attribute	Word	List (F)	List (I)
W1	Flowers/Insects	Pleasant/Unpleasant	Flowers Insects Pleasant Unpleasant	[golap, japa, shapla, genduk, yuthika, kamini, rajnigandha] [mashak, makshika, pipilica, murket, madhumakshika, tailpayika, utkun, patanga] [aadar, sbadhinata, susbasthya, bhalobasa, shanti, ananda, sukh, sundar, khushi] [apobyabohar, durghatana, asusthota, mri-tyu, dukkh, durgandh, lanchona, ghrina]	[golap, jaba, shapla, ganda, juin, kamini, rajnigandha] [masha, machi, pimpada, makadasa, mounmachi, telapoka, ukun, fading] [aadar, swadhinota, sushastho, bhalobasa, shanti, ananda, sukh, sundar, khushi] [apobyabohar, durghatana, asusthota, mri-tyu, dukkh, durgandh, lanchona, ghrina]
W2	Instruments/Weapons	Pleasant/Unpleasant	Instruments Weapons Pleasant Unpleasant	[gitara, haramoniyama, bina, behala, bansi, setara, ekatara, tabala] [tira, dhanusa, banduka, misaila, taloyara, raiphela, boma, ksuri] [aadar, sbadhinata, susbasthya, bhalobasa, shanti, ananda, sukh, sundar, khushi] [apobyabohar, durghatana, asusthota, mri-tyu, dukkh, durgandh, lanchona, ghrina]	[gitara, haramoniyama, bina, behala, bamsi, setara, ekatara, tabala] [tira, dhanuka, banduka, misaila, taloyara, raiphela, boma, churi] [aadar, swadhinota, sushastho, bhalobasa, shanti, ananda, sukh, sundar, khushi] [apobyabohar, durghatana, asusthota, mri-tyu, dukkh, durgandh, lanchona, ghrina]
W3	Male/Female (N)	Pleasant/Unpleasant	Male (N) Female (N) Pleasant Unpleasant	[mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba] [aadar, sbadhinata, susbasthya, bhalobasa, shanti, ananda, sukh, sundar, khushi] [apobyabohar, durghatana, asusthota, mri-tyu, dukkh, durgandh, lanchona, ghrina]	[mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba] [aadar, swadhinota, sushastho, bhalobasa, shanti, ananda, sukh, sundar, khushi] [apobyabohar, durghatana, asusthota, mri-tyu, dukkh, durgandh, lanchona, ghrina]
W4	Male/Female (N)	Career/Family	Male (N) Female (N) Career Family	[mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba] [byabasa, cakari, betana, aphis, karmas-thala, pesa] [bati, abhibhabaka, santana, paribara, bibaha, atmiya, sbajana]	[mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba] [byabasa, cakari, betana, aphis, karmas-thala, pesa] [bati, abhibhabaka, santana, paribara, biye, atmiya, sbajana]
W5	Male/Female (W)	Career/Family	Male (W) Female (W) Career Family	[chele, loka, purusa, bhratr, tata, mama, pu-tra, sbami] [matrka, mahila, nari, bhagini, kanya, mata, badhu, stri] [byabasa, cakari, betana, aphis, karmas-thala, pesa] [bati, abhibhabaka, santana, paribara, bibaha, atmiya, sbajana]	[chele, loka, purusa, bhai, caca, mama, pu-tra, sbami] [meye, mahila, nari, bona, kanya, ma, bau, stri] [byabasa, cakari, betana, aphis, karmas-thala, pesa] [bati, abhibhabaka, santana, paribara, biye, atmiya, sbajana]
W6	Math/Arts	Male/Female (W)	Math Arts Male (W) Female (W)	[ganita, jyamiti, ganana, sankhya, anka, samikarana, kona] [kabita, silpa, sahitya, upanyasa, nrtya, gana, calaccitra, abhinaya] [chele, loka, purusa, bhratr, tata, mama, pu-tra, sbami] [matrka, mahila, nari, bhagini, kanya, mata, badhu, stri]	[ganita, jyamiti, ganana, sankhya, anka, samikarana, kona] [kabita, silpa, sahitya, upanyasa, naca, gana, calaccitra, abhinaya] [chele, loka, purusa, bhai, caca, mama, pu-tra, sbami] [meye, mahila, nari, bona, kanya, ma, bau, stri]
W7	Math/Arts	Male/Female (N)	Math Arts Male (N) Female (N)	[ganita, jyamiti, ganana, sankhya, anka, samikarana, kona] [kabita, silpa, sahitya, upanyasa, nrtya, gana, calaccitra, abhinaya] [mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba]	[ganita, jyamiti, ganana, sankhya, anka, samikarana, kona] [kabita, silpa, sahitya, upanyasa, naca, gana, calaccitra, abhinaya] [mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba]
W8	Science/Arts	Male/Female (W)	Science Arts Male (W) Female (W)	[bijnana, prayukti, robata, padarthabijnana, rasayana, jibabijnana] [kabita, silpa, sahitya, upanyasa, nrtya, gana, calaccitra, abhinaya] [chele, loka, purusa, bhratr, tata, mama, pu-tra, sbami] [matrka, mahila, nari, bhagini, kanya, mata, badhu, stri]	[bijnana, prayukti, robata, padarthabijnana, rasayana, jibabijnana] [kabita, silpa, sahitya, upanyasa, naca, gana, calaccitra, abhinaya] [chele, loka, purusa, bhai, caca, mama, pu-tra, sbami] [meye, mahila, nari, bona, kanya, ma, bau, stri]
W9	Science/Arts	Male/Female (N)	Science Arts Male (N) Female (N)	[bijnana, prayukti, robata, padarthabijnana, rasayana, jibabijnana] [kabita, silpa, sahitya, upanyasa, nrtya, gana, calaccitra, abhinaya] [mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba]	[bijnana, prayukti, robata, padarthabijnana, rasayana, jibabijnana] [kabita, silpa, sahitya, upanyasa, naca, gana, calaccitra, abhinaya] [mohammad, ahmed, abdul, rahim, karim, ali, sheikh] [sharmin, jannatul, fatema, sadia, farajana, adiba]

2.2 Mask Prediction

2.2.1 Probabilistic Logarithmic Bias Scores Test

We investigate the mask prediction-based method using the (Kurita et al., 2019) framework. This technique evaluates bias in contextual models trained with a Masked Language Modeling (MLM) goal. As the training goal of BERT is to predict [MASK] tokens, we create unique template sentences for

each category of Target vs Attribute combinations (Table 1). We provide each category’s effect size using the expected values of the accompanying mask tokens. Any contrasting Target vs Attribute word combination can be used with our generic template sentences (Appendix A).

We calculate p_t and p_p to determine the bias, where [TARGET] is replaced by [MASK], the prob-

Table 2: SEAT test template sentences (F points Formal and I points Informal words; words highlighted with red point target words).

Category	Word (F & I)	Sentence
Male word	purusa	eti ekti purusa . oiti ekti purusa . ekthane ekti purusa ache. ai say purusa . sekhane ekjon purusa ache. purusa ekjon byakti.
Female word	nari	eti ekti nari . oiti ekti nari . ekthane ekti nari ache. ai say nari . sekhane ekjon nari ache. nari ekjon byakti.

ability is:

$$p_t = P([\text{MASK}] = [\text{TARGET}] \mid S)$$

And where both [TARGET] and [ATTRIBUTE] are replaced by [MASK], the prior probability is:

$$p_p = P([\text{MASK}] = [\text{TARGET}] \mid S)$$

The relationship between Target and Attribute is then calculated using our bias measure, $\log \frac{p_t}{p_p}$. We call p the Fill Bias Score, p_p the Prior Bias Score and $\frac{p_t}{p_p}$ the Probabilistic Logarithmic Bias Score or the Prior Corrected Score. To investigate how the variances affect the bias ratings, we also look at various phrase patterns with differing levels of context.

3 Preprocessing

3.1 Non Contextual Word

Words from each category in both formal and informal form, validated by several experts in the native Bengali, are shown in Table 1. These terms are used in each category in the WEAT trials, and models are used to extract their embeddings. Next, we compute bias detection. We employ this set of terms in various phrases and contexts for the remaining tests.

3.2 Non Contextual Sentence

We utilize template sentences to create sentences for the SEAT experiment, adding terms from Table 1 to each template. We employ the template sentences from the original SEAT experiment in their translated forms. Sentences that link to male and female words are shown in Table 2.

3.3 Contextual Word

We create the stimuli’s embeddings using widely-used language models that support Bengali (Table 3). We use the Bangla2B+ (Bhattacharjee et al., 2021) dataset to extract phrases with much context.

Then, using a pattern-matching technique, we extract sentences that include the Target vs Attribute terms from the unstructured raw data using the list of these words. In addition, we add extra sentences to the terms with a low sentence count to meet a minimal threshold to guarantee efficient data aggregation. It is inefficient and causes substantial data loss to only match root words due to the intricacy of Bengali word suffixes. Bengali word suffixes frequently include semantic values guaranteeing subject-verb agreement, resolving co-references, etc. Sometimes, even suffixes produce completely new words, changing the meaning of the text. To address this problem, we construct unique suffix groups corresponding to the most frequently related suffixes for every word in our chosen list. We create many versions of a root word and extract sentences that contain each variation by assigning each word to its corresponding set of suffixes. We input these phrases into a language model and use the target word embedding from the last layer to retrieve the matching embeddings. We employ almost three million phrases and more than two hundred fifty words from all categories to extract word embeddings and carry out the CEAT experiment. The goal of embedding extraction is to preserve semantic subtleties by attempting to maintain the complete word embedding, including its suffixes. After the model tokenizes the target word, we pool each fragment’s logit. We offer a more thorough examination of the dataset development process in Appendix B, with examples.

3.4 Contextual Sentence

We conduct tests for Bengali to compute probabilistic logarithmic bias using the context-based templating of (Kurita et al., 2019). To do this, we manually create five distinct kinds of context-aware phrase structures with placeholders for at-

Table 3: Language models list to extract embedding.

Model	Reference	Architecture	Objectives	# Layers	# Params	Details	Dimension
BanglaBERT Large (Generator)	(Bhattacharjee et al., 2021)	ELECTRA	MLM	24	52M	Final layer outputs	-
BanglaBERT Large (Discriminator)	(Bhattacharjee et al., 2021)	ELECTRA	RTD	24	339M	Final layer outputs	-
MuRIL Large (cased)	(Khanuja et al., 2021)	BERT	MLM and TLM	24	506M	First layer concealed units CWEs	1024
XLM-RoBERTa Large	(Conneau, 2019)	Transformer-based	Multilingual MLM	24	560M	First layer embeddings	1024

Table 4: An illustration of many sentence forms using positive and female words with varying degrees of context (S5 >> S1). Words highlighted with red color point target (subject) and blue color point attribute (object) words.

Category	Words	Sentences
S1	uchchakangkshi, narira	narira uchchakangkshi.
S2	uchchakangkshi, narira	narira khub uchchakangkshi prokritir hai.
S3	uchchakangkhar, narider	uchchakangkhar prati jhonk narider modhye beshi parilakshit hai.
S4	uchchakangksha, narider	uchchakangksha द्वारा चालित मनुश्रा, समजर सम्मिलितो उन्नयानेर पोरिबोटे प्रयाशै byaktigato safalya arjaner upper manoyog den, jekhane nijeder akankhake agradhikar dewa hai. narider ai uchchakangkshi prokritir karone samajer gotishilotayo ekti ullekhjoggyo poriborton ghotiche, jekhane byaktigato lakshya kakhano kakhano samajik kalyanke chadi yay.
S5	uchchakangkshi, narira	peshagato ebong rajnaitik kshetre, uchchakangkshi narira kakhano kakhano byaktigato un-nayanke agradhikar den ebong emon kichu kaushal prayog karen ya naitik maner sathe sangatipurna noy. sadharonoto, shaktishali uchchakangkhar द्वारा पोरिचालितो byaktira prayashai naari han, yara natun path tairi karchen kintu atmaswarth ebong samajik abodaner modhye varsamyo bajaay rakhar prashn tulechen.

tribute (Positive vs Negative)³ words and target (Male vs Female) words (Table 4), validated by several experts in the native Bengali. These vary from straightforward statements with no context, S1, to sentences extracted from the Bangla2B+ dataset with substantial context, S5. To incorporate changes in context length, we aimed to minimize the amount of structures used while introducing variety in subject and object placements inside phrases. We use 70 words with negative attributes and 110 words with positive attributes to create our test dataset. This approach produces a wide range of phrases representing different linguistic contexts. We also employ four distinct names for men and women. We produced 3600 sentences, which together reflect the various circumstances being tested.

4 Experiment Results

4.1 Contextual CEAT

Table 5 and Table 6 illustrate how we use CEAT to evaluate how contextual variation affects bias. As demonstrated by (Guo and Caliskan, 2021) findings, which indicate no discernible difference between samples of $N = 1000$ and $N = 10000$, the selection of sample size $N = 7000$ is supported. Our research aims to clarify how the magnitude of effect is impacted by contextual input duration. The effect size illustrates how the observed bias

varies according to segment length and stabilizes as contextual information grows. The dynamic variations in effect size between two models as context length fluctuates are depicted in Figure 1 and Figure 2. The ideal context length for reliable outcomes is moderate, at about 25 words. We select combinations for every CEAT experiment using fixed and random sets. Fixed sets enable cross-model comparisons, while random sets evaluate how context variation affects effect magnitude for a given segment length. There is no discernible difference in the effect magnitude between experiments with 7000 and 1000 samples; nevertheless, fewer instances produced statistically significant results. According to our findings, statistically significant bias varies depending on the model and occasionally shows bias in other ways, both in formal and informal Bengali. However, it shows that the bias in informal Bengali is slightly larger than in formal Bengali. Interestingly, the MuRIL Large (cased) model exhibits increased sensitivity to context for fixed samples in both formal and informal Bengali. According to Table 5 and Table 6, statistically insignificant findings typically occur in shorter segment lengths.

4.2 Contextual Probabilistic Logarithmic Bias Scores

(Kurita et al., 2019) proposed the template-based methodology, which shows improved consistency in human bias evaluation and provides a straightforward mechanism for querying models based on

³We use the MIT Ideonomy categories, translating words both in formal and informal version verbatim <https://ideonomy.mit.edu/essays/traits.html>

Table 5: Measures of formal Bengali d value of bias for various language models which includes θ pooling $N = 7000$ random-effects model samples (* points insignificant at $p < 0.007$).

Model	Length	W1		W2		W3		W4		W5		W6		W7		W8		W9	
		Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand
BanglaBERT Large (Generator)	9	1.223	1.223	-0.224	-0.224	0.180	0.180	0.637	0.638	0.261	0.261	0.256	0.256	-0.641	-0.641	0.364	0.364	-0.589	-0.589
	25	1.206	1.206	-0.237	-0.237	0.158	0.158	0.73	0.73	0.25	0.25	0.139	0.139	-0.637	-0.637	0.285	0.285	-0.641	-0.641
	75	1.204	1.204	-0.254	-0.254	0.163	0.163	0.728	0.728	0.252	0.252	0.126	0.126	-0.636	-0.636	0.267	0.267	-0.652	-0.652
	> 75	1.205	1.205	-0.256	-0.256	0.164	0.164	0.728	0.728	0.243	0.243	0.121	0.121	-0.635	-0.635	0.271	0.271	-0.649	-0.649
BanglaBERT Large (Discriminator)	9	0.549	0.549	-0.258	-0.269	0.002*	-0.001*	0.030	0.032	-0.034	-0.039	0.021*	0.011*	0.251	0.253	-0.038	-0.050	-0.113	-0.123
	25	0.464	0.464	-0.178	-0.177	0.013*	0.007*	0.021	0.016	-0.036	-0.031	0.059	0.064	0.267	0.264	-0.034	-0.042	-0.143	-0.140
	75	0.458	0.460	-0.156	-0.158	0.015	0.021	0.016	0.019	-0.028	-0.027	0.067	0.064	0.275	0.262	-0.048	-0.054	-0.151	-0.157
	> 75	0.459	0.451	-0.166	-0.156	0.016	0.016	0.014*	0.015	-0.021*	-0.026	0.061	0.061	0.268	0.286	-0.044	-0.047	-0.14	-0.151
MuRIL (cased)	9	1.191	1.192	0.475	0.476	0.477	0.482	0.014	0.022	0.224	0.228	0.419	0.412	-0.158	-0.149	-0.056	-0.058	-0.007*	-0.020
	25	1.213	1.212	0.37	0.378	0.633	0.628	-0.083	-0.09	0.254	0.252	0.424	0.43	-0.223	-0.228	0.005*	0.011*	-0.182	-0.181
	75	1.198	1.199	0.377	0.365	0.659	0.647	-0.087	-0.091	0.265	0.256	0.407	0.419	-0.270	-0.275	0.006*	0.004*	-0.211	-0.204
	> 75	1.206	1.199	0.375	0.37	0.649	0.659	-0.083	-0.081	0.262	0.255	0.406	0.415	-0.282	-0.268	0.008	0.011	-0.213	-0.202
XLM-RoBERTa Large	9	0.277	0.271	0.564	0.572	0.063	0.062	-0.208	-0.201	-0.138	-0.150	-0.101	-0.112	-1.258	-1.260	-0.074	-0.078	-0.294	-0.291
	25	0.476	0.470	0.736	0.747	0.061	0.048	-0.275	-0.261	-0.197	-0.212	-0.172	-0.172	-1.193	-1.200	-0.090	-0.082	-0.309	-0.305
	75	0.491	0.508	0.767	0.772	0.046	0.057	-0.290	-0.265	-0.208	-0.208	-0.141	-0.150	-1.189	-1.191	-0.102	-0.073	-0.311	-0.307
	> 75	0.493	0.48	0.765	0.764	0.054	0.038	-0.284	-0.324	-0.208	-0.214	-0.132	-0.131	-1.189	-1.192	-0.102	-0.104	-0.315	-0.31

Table 6: Measures of informal Bengali d value of bias for various language models which includes θ pooling $N = 7000$ random-effects model samples (* points insignificant at $p < 0.007$).

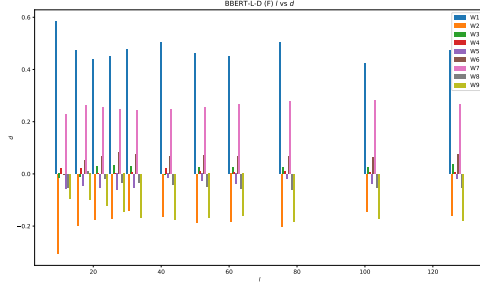
Model	Length	W1		W2		W3		W4		W5		W6		W7		W8		W9	
		Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand
BanglaBERT Large (Generator)	9	1.227	1.227	-0.228	-0.228	0.184	0.184	0.641	0.642	0.265	0.265	0.260	0.260	-0.645	-0.645	0.368	0.368	-0.593	-0.593
	25	1.21	1.21	-0.241	-0.241	0.162	0.162	0.734	0.734	0.254	0.254	0.143	0.143	-0.641	-0.641	0.289	0.289	-0.645	-0.645
	75	1.208	1.208	-0.258	-0.258	0.167	0.167	0.732	0.732	0.256	0.256	0.130	0.130	-0.640	-0.640	0.271	0.271	-0.656	-0.656
	> 75	1.209	1.209	-0.260	-0.260	0.168	0.168	0.732	0.732	0.247	0.247	0.125	0.125	-0.639	-0.639	0.275	0.275	-0.653	-0.653
BanglaBERT Large (Discriminator)	9	0.553	0.553	-0.262	-0.273	-0.002*	-0.005*	0.034	0.036	-0.038	-0.043	0.025*	0.015*	0.255	0.257	-0.042	-0.054	-0.117	-0.127
	25	0.468	0.468	-0.182	-0.181	0.017*	0.011*	0.025	0.2	-0.04	-0.035	0.063	0.068	0.271	0.268	-0.038	-0.046	-0.147	-0.144
	75	0.462	0.464	-0.160	-0.162	0.019	0.025	0.020	0.023	-0.032	-0.031	0.071	0.068	0.279	0.266	-0.052	-0.058	-0.155	-0.161
	> 75	0.463	0.455	-0.170	-0.160	0.020	0.020	0.018*	0.019	-0.025*	-0.030	0.065	0.065	0.272	0.290	-0.048	-0.051	-0.144	-0.155
MuRIL (cased)	9	1.195	1.196	0.479	0.480	0.481	0.486	0.018	0.026	0.228	0.232	0.423	0.416	-0.162	-0.153	-0.060	-0.062	-0.011*	-0.024
	25	1.217	1.216	0.374	0.382	0.637	0.632	-0.087	-0.094	0.258	0.256	0.428	0.434	-0.227	-0.232	0.009*	0.015*	-0.186	-0.185
	75	1.202	1.203	0.381	0.369	0.663	0.651	-0.091	-0.087	0.269	0.260	0.411	0.423	-0.274	-0.279	0.010*	0.008*	-0.215	-0.208
	> 75	1.210	1.203	0.379	0.374	0.653	0.663	-0.087	-0.085	0.266	0.259	0.410	0.419	-0.286	-0.272	0.012	0.015	-0.217	-0.206
XLM-RoBERTa Large	9	0.281	0.275	0.568	0.576	0.067	0.066	-0.212	-0.205	-0.142	-0.154	-0.105	-0.116	-1.262	-1.264	-0.078	-0.082	-0.298	-0.295
	25	0.48	0.474	0.74	0.751	0.065	0.052	-0.279	-0.265	-0.201	-0.216	-0.176	-0.176	-1.197	-1.204	-0.094	-0.086	-0.313	-0.309
	75	0.495	0.512	0.771	0.776	0.050	0.061	-0.294	-0.269	-0.212	-0.212	-0.145	-0.154	-1.193	-1.195	-0.106	-0.077	-0.315	-0.311
	> 75	0.497	0.484	0.769	0.768	0.058	0.042	-0.288	-0.328	-0.212	-0.218	-0.136	-0.135	-1.193	-1.196	-0.106	-0.108	-0.319	-0.314

Table 7: A few instances of positive and negative characteristics employed in the Probabilistic Logarithmic Bias Score Test (F points Formal and I points Informal words).

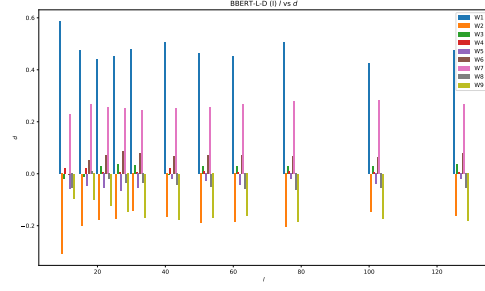
Category	Word (F)	Word (I)
Positive	[akritrim, atyadhunik, adwitiya, anugato, antadrishtipurna, abichal, abhiyojanyogya, abhedya, akarshaniya, atmabishwasi, atmasamalochak, adarshbadi, antorik, abegprobon, ashabadi, uchchakangkshi]	[akritrim, atyadhunik, adwitiya, anugato, antadrishtipurna, abichal, abhiyojanyogya, abhedya, akarshaniya, atmabishwasi, atmasamalochak, adarshbadi, antorik, abegprobon, ashabadi, uchchakangkshi]
Negative	[akritajna, agochalo, agya, aduradarshi, anubhutaheen, apachandaniya, apamanjanak, aparadhi, gowrn, abagyapurna, abibechak, abiswasto, ajouktik, alas, asangatipurna, asatark, asat, asantushta, asammanjanak]	[akritajna, agochalo, agya, aduradarshi, anubhutaheen, apachandaniya, apamanjanak, aparadhi, gowrn, abagyapurna, abibechak, abiswasto, ajouktik, alas, asangatipurna, asatark, asat, asantushta, asammanjanak]

modeling objectives. Two components make up the Fill Bias Score, which offers a direct look at model biases: the intrinsic language bias, which is measured by the prior bias score and the bias brought about by the presence of attributes, which is the actual bias measure known as the Prior Corrected Score or Probabilistic Logarithmic Bias Score. Models interact with genuinely occurring language in real-world situations. We pay attention to analyzing the negative and positive attributes in the BanglaBERT Generator situation in Figure 3 and Figure 4 respectively (further outcomes are presented in Figure 5 and Figure 6–Appendix C). If the corrected bias scores are distributed consistently across all phrase forms, then the difference in the prior bias distribution is caused by innate linguistic

bias. An enlarged range results from the preceding bias score for sentence structures S1 through S3, showing increasing inherent linguistic bias by adding new words in formal and informal Bengali. Values tend to cluster around a neutral point for S4 through S5, which is the opposite pattern. A change in the model’s behavior when the attribute adopts a more context-rich configuration is shown by the observed trend from S1 to S3, which highlights the model’s unique preferences in formal and informal Bengali. Additionally, several adjusted bias scores change from negative to positive when context increases (words lists are shown in Table 7). In formal and informal Bengali, a more organic language context is simulated by sentence forms S4 and S5. The model shifts attention and

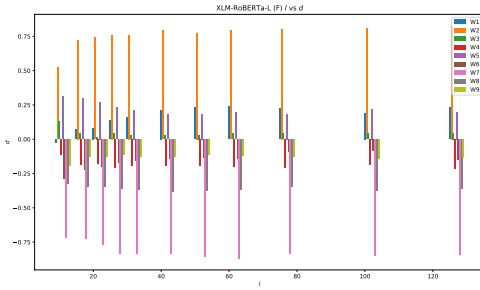


(a)

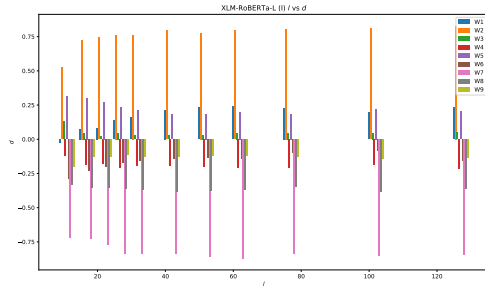


(b)

Figure 1: BanglaBERT Large (Discriminator) performance on l variation affects d value on formal and informal Bengali for category variations W1 through W9 (For a certain segment length using a sample size of $N=1000$, values that are statistically significant at $p < 0.007$ are presented).

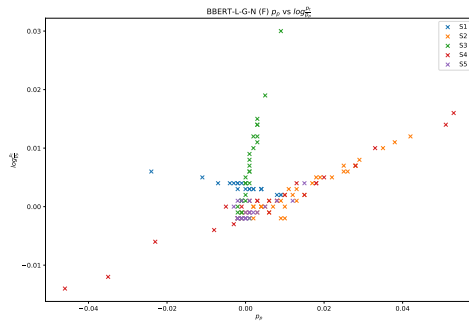


(a)

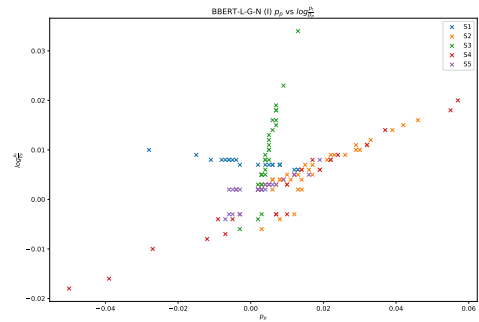


(b)

Figure 2: XLM-RoBERTa Large performance on l variation affects d value on formal and informal Bengali for category variations W1 through W9 (For a certain segment length using a sample size of $N=1000$, values that are statistically significant at $p < 0.007$ are presented).



(a)



(b)

Figure 3: BanglaBERT Large (Generator) relation between the p_p and $\log \frac{p_t}{p_p}$ value of bias for negative attributes for sentence constructions S1 through S5.

reduces the difference between the probabilities for male and female target words when there is too much information, which allows the model to assign greater probabilities to non-target phrases. Figure 3, Figure 4, Figure 5, and Figure 6 charts make this behavior clear. The charts show that

both the corrected and prior bias scores have values closely grouped around the neutral point. However, it shows that the bias in informal Bengali is slightly larger than in formal Bengali.

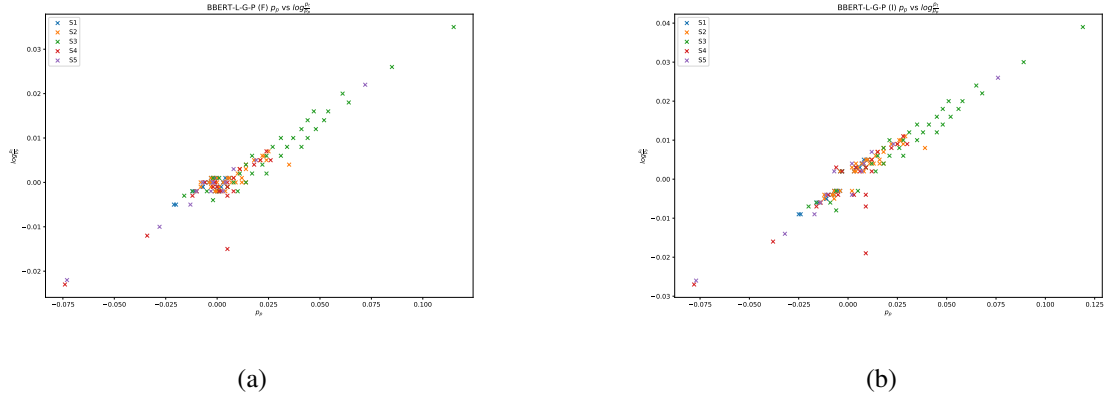


Figure 4: BanglaBERT Large (Generator) relation between the p_p and $\log \frac{p_t}{p_p}$ value of bias for positive attributes for sentence constructions S1 through S5.

5 Conclusion

This study aims to investigate bias in formal and informal Bengali in Bengali language models by building a curated dataset. We contend that the quantity of context included in templates affects the bias result for formal and informal Bengali. Additional research on other low-resource languages can be done. To reduce detrimental bias in Bangla embeddings, we intend to explore the impact of bias on downstream applications of Bangla language models in the future. We also hope to expand these efforts to generative models by creating language-specific debiasing techniques.

Limitations

Certain limits point to areas for further research. Since most of our datasets are derived from pre-existing datasets, they are synthetic to comply with accepted bias measurement techniques. Furthermore, the primary emphasis of our research is gender bias in both formal and informal Bengali. For this specific project, we are motivated to address gender bias only for three reasons. First, Bengali exhibits complex linguistic structures, rich morphology, and dialectical variations in formal and informal forms. Secondly, bias against gender is pervasive everywhere. Thirdly, gender bias shows far more subtle differences than the others, which makes it a fertile field for research. Some shortcomings in our methods for assessing inherent bias have previously been identified (Blodgett et al., 2020). Rather than concentrating on the failures of the approaches that have already been used, we aimed to provide the framework for further research on Bangla bias. Future developments can investigate

more testing with different biases, including social, religious, political, etc. Moreover, using bias analysis static templates without taking downstream applications into account is another drawback of our research. Additionally, not covering the bias properties of generative language models is another limitation of our study. Future research might investigate these areas with corresponding debiasing techniques.

Ethics Statement

People may find our study potentially upsetting since it focuses on gender bias and statistics associated with this social prejudice. Nonetheless, this study must be carried out to guarantee equity in the natural language model sector. We also recognize that although our study emphasizes gender as a binary entity, non-binary entities may warrant more research.

References

- Abhik Bhattacharjee, Tahmid Hasan, Wasi Uddin Ahmad, Kazi Samin, Md Saiful Islam, Anindya Iqbal, M Sohel Rahman, and Rifat Shahriyar. 2021. Banglabert: Language model pretraining and benchmarks for low-resource language understanding evaluation in bangla. *arXiv preprint arXiv:2101.00204*.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of "bias" in nlp. *arXiv preprint arXiv:2005.14050*.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. Man is to computer programmer as woman is to home-maker? debiasing word embeddings. *Advances in neural information processing systems*, 29.

- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.
- A Conneau. 2019. Unsupervised cross-lingual representation learning at scale. *arXiv preprint arXiv:1911.02116*.
- David Esiobu, Xiaoqing Tan, Saghar Hosseini, Megan Ung, Yuchen Zhang, Jude Fernandes, Jane Dwivedi-Yu, Eleonora Presani, Adina Williams, and Eric Michael Smith. 2023. Robbie: Robust bias evaluation of large generative language models. *arXiv preprint arXiv:2311.18140*.
- Wei Guo and Aylin Caliskan. 2021. Detecting emergent intersectional biases: Contextualized word embeddings contain a distribution of human-like biases. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages 122–133.
- Larry V Hedges. 1983. A random effects model for effect sizes. *Psychological bulletin*, 93(2):388.
- Simran Khanuja, Diksha Bansal, Sarvesh Mehtani, Savya Khosla, Atreyee Dey, Balaji Gopalan, Dilip Kumar Margam, Pooja Aggarwal, Rajiv Teja Nagipogu, Shachi Dave, and 1 others. 2021. Muril: Multilingual representations for indian languages. *arXiv preprint arXiv:2103.10730*.
- Ryan Koo, Minhwa Lee, Vipul Raheja, Jong Inn Park, Zae Myung Kim, and Dongyeop Kang. 2023. Benchmarking cognitive biases in large language models as evaluators. *arXiv preprint arXiv:2309.17012*.
- Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. *arXiv preprint arXiv:1906.07337*.
- Mascha Kurpicz-Briki. 2020. Cultural differences in bias? origin and gender bias in pre-trained german and french word embeddings.
- Anne Lauscher, Rafik Takieddin, Simone Paolo Ponzetto, and Goran Glavaš. 2020. Araweat: Multidimensional analysis of biases in arabic word embeddings. *arXiv preprint arXiv:2011.01575*.
- Sheng Liang, Philipp Dufter, and Hinrich Schütze. 2020. Monolingual and multilingual reduction of gender bias in contextualized representations.
- Ruixi Lin and Hwee Tou Ng. 2023. Mind the biases: Quantifying cognitive biases in language model prompting. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5269–5281.
- Vijit Malik, Sunipa Dev, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2021. Socially aware bias measurements for hindi language representations. *arXiv preprint arXiv:2110.07871*.
- Chandler May, Alex Wang, Shikha Bordia, Samuel R Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. *arXiv preprint arXiv:1903.10561*.
- Tomas Mikolov. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 3781.
- Douglas C Montgomery and George C Runger. 2010. *Applied statistics and probability for engineers*. John Wiley & sons.
- Rodrigo Alejandro Chávez Mulsa and Gerasimos Spanakis. 2020. Evaluating bias in dutch word embeddings. *arXiv preprint arXiv:2011.00244*.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*.
- Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Arun K Pujari, Ansh Mittal, Anshuman Padhi, Anshul Jain, Mukesh Jadon, and Vikas Kumar. 2019. Debiasing gender biased hindi words with word-embedding. In *Proceedings of the 2019 2nd International Conference on Algorithms, Computing and Artificial Intelligence*, pages 450–456.
- Marnie E Rice and Grant T Harris. 2005. Comparing effect sizes in follow-up studies: Roc area, cohen’s d, and r. *Law and human behavior*, 29:615–620.
- Jayanta Sadhu, Ayan Antik Khan, Abhik Bhattacharjee, and Rifat Shahriyar. 2024. An empirical study on the characteristics of bias upon context length variation for bangla. *arXiv preprint arXiv:2406.17375*.
- Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. 2024. In-context impersonation reveals large language models’ strengths and biases. *Advances in Neural Information Processing Systems*, 36.
- Yuxuan Wan, Wenxuan Wang, Pinjia He, Jiazhen Gu, Haonan Bai, and Michael R Lyu. 2023. Biasasker: Measuring the bias in conversational ai system. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 515–527.
- Yue Yu, Yuchen Zhuang, Jieyu Zhang, Yu Meng, Alexander J Ratner, Ranjay Krishna, Jiaming Shen, and Chao Zhang. 2024. Large language model as attributed training data generator: A tale of diversity and bias. *Advances in Neural Information Processing Systems*, 36.

A Probabilistic Logarithmic Bias Scores Test Sentences

Example sentences for the log probability bias experiment are shown in Table 8. The target and attribute words are highlighted in each phrase; the probability bias score and impact magnitude are determined by methodically masking these terms. We underline the attribute words in blue and the target terms in red to improve readability. Following the templating procedure, we compute the logarithmic differences between these probabilities and the p_t and $\log \frac{p_t}{p_p}$.

B CEAT Extraction

For the CEAT test, the Bangla2B+ dataset is used to extract sentences. Table 9 shows instances of sentences that occur spontaneously in the sample. The presence and uniqueness of suffixes significantly impact how sentences based on root words are extracted in Bengali. Bengali words generally contain suffixes, which need careful attention in sentence extraction approaches, in contrast to English, where words frequently occur in their basic form. For example, Table 10 shows that root words like "ma" have a variety of suffixes, including "_yeder," "_yera," "_yer," and "_ke," which allow for precise identification in phrases like "mayeder abadan" and "mayer ranna." But using the wrong suffixes ("_tir," "_ti," and "_ra") results in inaccurate matches, as in "matir basne" and "matite shikad." A thorough system associating every word with a different group of suffixes was created to overcome this issue. This system consists of twenty-one groups, each with two to fifteen suffixes.

C Extended Experiment Results

A lower sample size, $N = 1000$, is the subject of this section. Table 11 and Table 12 displays our findings. These data points, however, typically only exhibit slight alterations concerning Table 5 and Table 6 in Section 4.1. Table 11 and Table 12 still clearly display the main features of the model, which were emphasized for the $N = 7000$ sample in Table 5 and Table 6 in Section 4.1. This shows that comparable outcomes may be obtained even with a smaller sample size, particularly when resources are limited. However, it is important to understand that although the general patterns in the model's behavior do not change, there are subtle changes in the statistical importance of specific data points.

We suggest that increasing the sample size might be a viable line of inquiry in light of these findings. Finding the ideal sample size that guarantees accurate results without compromising statistical significance presents an intriguing topic for further study and detailed investigation. A wholesome level of our intrinsic bias measurements have been indexed in Table 13 and Table 14.

Table 8: Probabilistic Logarithmic Bias Score Test template sentences (F points Formal and I points Informal words; words highlighted with red color point target and blue color point attribute words).

Category	Target vs Attribute (F & I)	Sentence
Instruments/Weapons vs Pleasant/Unpleasant	Instruments vs Pleasant	behala er awaz shunley shanti er katha mane pade.
	Instruments vs Unpleasant	behala er awaz shunley mrityu er katha mane pade.
	Weapons vs Pleasant	misaila er awaz shunley shanti er katha mane pade.
	Weapons vs Unpleasant	misaila er awaz shunley mrityu er katha mane pade.

Table 9: Sentences extracted for the CEAT test (F points Formal and I points Informal words; words highlighted with red point target words).

Category	Word (F & I)	Sentences
Male word	purusa	jaibik arthe, purusa sadharonoto say byaktike bojhaay jaar Y chromosome thake ebong hormones yeman testosteroner pariman beshi.
Female word	nari	jaibik arthe, jaar XX chromosome thake ebong jaar sharire pradhanat estrogen o progesterone hormone karyakar take say byaktike nari bojhaay.

Table 10: Suffix groups’ distinctiveness and significance for extracting embedding in CEAT (Words highlighted with green color points correct, whereas red color points are wrong words, respectively, concerning the specific targeted context).

Root Word	Token Length	Correctness	Suffixs	Regex
ma	3	Correct	_yeder, _yera, _yer, _ke	mayeder abadan santaner jibon gathone aparisim. sab mayera tader santander jonno sera chan. mayer ranna amar sabcheve priya. aami amar make boi kine dite chai.
		Wrong	_tir, _ti, _ra	chheleti matir basne paani dhalche. gachti matite shikad chhadiyechhe. say rate ekti poka mara hoyechilo.

Table 11: Measures of formal Bengali d value of bias for various language models which includes θ pooling $N = 1000$ random-effects model samples (* points insignificant at $p < 0.007$)

Model	Length	W1		W2		W3		W4		W5		W6		W7		W8		W9	
		Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand
BanglaBERT Large (Generator)	9	1.224	1.221	-0.222	-0.230	0.188	0.178	0.641	0.639	0.257	0.273	0.252	0.260*	-0.627	-0.619	0.364	0.373	-0.591	-0.582
	25	1.204	1.197	-0.247	-0.233	0.147	0.148	0.732	0.749	0.243	0.260	0.174	0.155	-0.640	-0.628	0.304	0.275	-0.633	-0.633
	75	1.210	1.210	-0.243	-0.252	0.163	0.164	0.721	0.733	0.246	0.251	0.148	0.128	-0.637	-0.628	0.277	0.280	-0.645	-0.644
BanglaBERT Large (Discriminator)	9	1.218	1.212	-0.254	-0.250	0.176	0.165	0.724	0.743	0.241	0.254	0.138	0.123	-0.626	-0.618	0.281	0.303	-0.648	-0.636
	25	0.530	0.572	-0.261	-0.244	-0.016*	0.000*	0.017*	0.014*	-0.051*	-0.054*	0.021*	0.019	0.250	0.264	-0.039*	-0.029	-0.122	-0.115
	75	0.417	0.460	-0.157	-0.190	0.024*	0.003*	0.030	0.005*	-0.008*	-0.046	0.081	0.072	0.292	0.280	-0.040	-0.049	-0.145	-0.141
MuRIL (cased)	9	1.211	1.193	-0.178	-0.172	0.020*	0.012*	0.013*	0.010*	-0.006*	0.003*	0.052	0.073	0.271	0.241	-0.053*	-0.029	-0.165	-0.167
	25	0.498	0.443	-0.145	-0.126	0.027	0.026*	0.005*	0.015*	-0.019*	-0.056	0.055	0.076	0.293	0.257	-0.044	-0.080	-0.157	-0.161
	75	1.199	1.199	0.356	0.361	0.646	0.657	-0.091	-0.098	0.243	0.261	0.384	0.417	-0.286	-0.274	0.012	-0.013*	-0.205	-0.203
XLM-RoBERTa Large	9	0.273	0.276	0.572	0.570	0.062	0.066	-0.201*	-0.214	-0.150	-0.156	-0.112	-0.109	-1.260	-1.259	-0.078	-0.074	-0.291	-0.289
	25	0.470	0.467	0.747	0.731	0.048	0.042	-0.261	-0.277	-0.212	-0.191	-0.172	-0.222	-1.200	-1.198	-0.082*	-0.087	-0.305	-0.313
	75	0.508	0.497	0.772	0.767	0.057	0.049	-0.265	-0.305	-0.208	-0.200	-0.150	-0.137	-1.191	-1.185	-0.073*	-0.091	-0.307	-0.307
	> 75	0.480	0.492	0.764	0.767	0.038	0.056	-0.324	-0.300	-0.214	-0.204	-0.131	-0.104	-1.192	-1.185	-0.104*	-0.100	-0.310	-0.311

Table 12: Measures of informal Bengali d value of bias for various language models which includes θ pooling $N = 1000$ random-effects model samples (* points insignificant at $p < 0.007$)

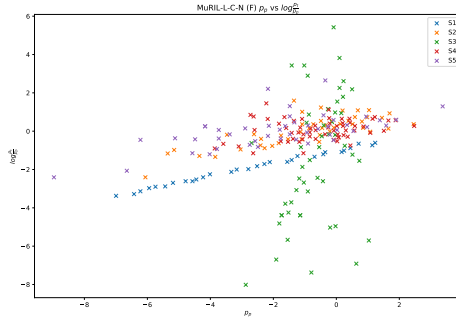
Model	Length	W1		W2		W3		W4		W5		W6		W7		W8		W9	
		Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand	Fixed	Rand
BanglaBERT Large (Generator)	9	1.228	1.225	-0.226	-0.234	0.192	0.182	0.645	0.643	0.261	0.277	0.256	0.264*	-0.631	-0.623	0.368	0.377	-0.595	-0.586
	25	1.208	1.201	-0.251	-0.237	0.151	0.152	0.736	0.753	0.247	0.264	0.178	0.159	-0.644	-0.632	0.308	0.279	-0.637	-0.637
	75	1.214	1.214	-0.247	-0.256	0.167	0.168	0.725	0.737	0.250	0.255	0.132	0.132	-0.641	-0.632	0.281	0.284	-0.649	-0.648
BanglaBERT Large (Discriminator)	9	1.222	1.216	-0.258	-0.254	0.180	0.169	0.728	0.747	0.245	0.258	0.142	0.127	-0.630	-0.622	0.285	0.307	-0.652	-0.640
	25	0.534	0.576	-0.265	-0.248	-0.020*	-0.004*	0.021*	0.018*	-0.055*	-0.058*	0.025*	0.023	0.254	0.268	-0.043*	-0.033	-0.126	-0.119
	75	0.421	0.464	-0.161	-0.194	0.028*	0.007*	0.034	0.009*	-0.012*	-0.050	0.085	0.076	0.296	0.284	-0.044	-0.053	-0.149	-0.145
MuRIL (cased)	9	1.214	1.193	-0.182	-0.176	0.024*	0.016*	0.017*	0.014*	-0.010*	-0.001*	0.056	0.077	0.275	0.245	-0.057*	-0.033	-0.169	-0.171
	25	0.502	0.447	-0.149	-0.130	0.031	0.030*	0.009*	0.019*	-0.023*	-0.060	0.059	0.080	0.297	0.261	-0.048	-0.084	-0.161	-0.165
	75	1.198	1.194	0.480	0.471	0.482	0.499	0.017	0.039	0.207	0.233	0.414	0.413	-0.143	-0.133	-0.070	-0.057	-0.015*	-0.018*
XLM-RoBERTa Large	9	1.220	1.209	0.378	0.377	0.628	0.621	-0.067	-0.084	0.276	0.241	0.438	0.444	-0.228	-0.235	0.028	0.008*	-0.181	-0.209
	25	1.215	1.197	0.379	0.384	0.655	0.657	-0.082	-0.078	0.275	0.246	0.438	0.416	-0.286	-0.277	0.020	0.013*	-0.233	-0.212
	75	1.203	1.203	0.360	0.365	0.650	0.661	-0.095	-0.102	0.247	0.265	0.388	0.421	-0.290	-0.278	0.016	-0.017*	-0.209	-0.207
	9	0.277	0.280	0.576	0.574	0.066	0.070	-0.205*	-0.218	-0.154	-0.160	-0.116	-0.113	-1.264	-1.263	-0.082	-0.078	-0.295	-0.293
	25	0.474	0.471	0.751	0.735	0.052	0.046	-0.265	-0.281	-0.216	-0.195	-0.176	-0.226	-1.204	-1.202	-0.086*	-0.091	-0.309	-0.317
	75	0.512	0.501	0.776	0.771	0.061	0.053	-0.269	-0.309	-0.212	-0.204	-0.154	-0.141	-1.195	-1.189	-0.077*	-0.095	-0.311	-0.311
	> 75	0.484	0.496	0.768	0.771	0.042	0.060	-0.328	-0.304	-0.218	-0.208	-0.135	-0.108	-1.196	-1.189	-0.108*	-0.104	-0.314	-0.315

Table 13: Measurements of formal Bengali d value of bias for different studies (* points statistical significance at $p < 0.07$).

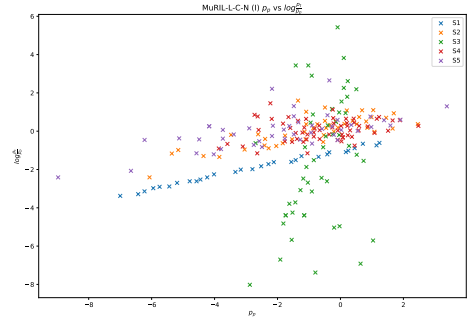
Category	WEAT (word2vec)	WEAT (GloVe)	SEAT	CEAT	Probabilistic Logarithmic Bias Test
W1	1.75*	1.25*	0.87*	1.205*	0.87*
W2	1.51*	0.97*	-0.01	-0.206*	0.40*
W3	0.36	1.33*	0.76*	0.162*	0.20
W4	1.42*	-0.16	-0.56	0.619*	0.69*
W5	0.40	0.15	-0.42	0.243*	0.60*
W6	0.98*	0.66*	-0.15	0.238*	0.91*
W7	-0.15	-0.91	-0.65	-0.623*	0.46*
W8	-0.20	-0.18	-0.74	0.346*	0.96*
W9	0.21	-1.01	-1.11	-0.571*	0.68*

Table 14: Measurements of informal Bengali d value of bias for different studies (* points statistical significance at $p < 0.07$).

Category	WEAT (word2vec)	WEAT (GloVe)	SEAT	CEAT	Probabilistic Logarithmic Bias Test
W1	1.79*	1.29*	0.91*	1.245*	0.91*
W2	1.55*	1.01*	-0.05	-0.246*	0.44*
W3	0.40	1.37*	0.80*	0.202*	0.24
W4	1.46*	-0.20	-0.60	0.659*	0.73*
W5	0.44	0.19	-0.46	0.283*	0.64*
W6	1.02*	0.70*	-0.19	0.278*	0.95*
W7	-0.19	-0.95	-0.69	-0.663*	0.50*
W8	-0.24	-0.22	-0.78	0.386*	1.00*
W9	0.25	-1.05	-1.15	-0.611*	0.72*

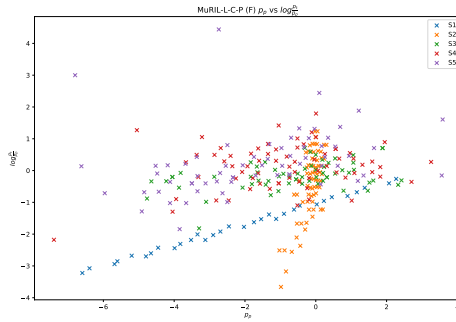


(a)

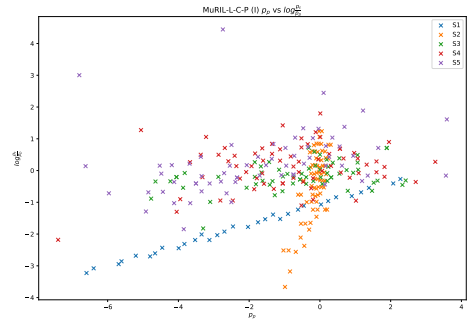


(b)

Figure 5: MuRIL - Large (cased) relation between the p_p and $\log p_p^t$ value of bias for negative attributes for sentence constructions S1 through S5.



(a)



(b)

Figure 6: MuRIL - Large (cased) relation between the p_p and $\log p_p^t$ value of bias for positive attributes for sentence constructions S1 through S5.