# Toward Cultural Bias Evaluation Datasets: The Case of Bengali Gender, Religious, and National Identity

Dipto Das
Dept of Information Science
University of Colorado Boulder
Boulder, CO, United States
dipto.das@colorado.edu

Shion Guha
Faculty of Information
University of Toronto
Toronto, ON, Canada
shion.guha@utoronto.ca

Bryan Semaan
Dept of Information Science
University of Colorado Boulder
Boulder, CO, United States
bryan.semaan@colorado.edu

## Abstract

Critical studies found NLP systems to bias based on gender and racial identities. However, few studies focused on identities defined by cultural factors like religion and nationality. Compared to English, such research efforts are even further limited in major languages like Bengali due to the unavailability of labeled datasets. This paper describes a process for developing a bias evaluation dataset highlighting cultural influences on identity. We also provide a Bengali dataset as an artifact outcome that can contribute to future critical research.

## 1 Introduction

Bias, in the context of computing systems, is where sociotechnical systems systematically and unfairly discriminate against certain individuals or social groups in favor of others (Friedman and Nissenbaum, 1996; Blodgett et al., 2020). People often identify through their perceived memberships in certain groups (Tajfel, 1974). While computational linguists have studied gender and racial biases (Kiritchenko and Mohammad, 2018), systematic discrimination of language technologies based on various cultural factors like religion and nationality has received little attention. Moreover, critical studies examining these biases mostly focused on NLP systems in a handful of languages, whereas many languages with sizeable numbers of speakers do not have enough resources like datasets to pursue similar studies.

According to (Joshi et al., 2020), whereas 0.28% of global languages (e.g., English, Spanish, Japanese) reap benefits from NLP breakthroughs, 88.38% of languages have virtually no data to use. They also found that while English and Bengali are the third and sixth largest languages by the number of native speakers (Lane, 2023), the former has hundreds of times more resources than the latter in Linguistic Data Consortium, Language Resources and Evaluation, and Wikipedia, and thou-

sands more resources in the Web overall. The difference in available resources like labeled datasets impedes the progress of critical studies aimed at fairness, transparency, and identifying biases in such under-represented languages. In the absence of native resources, many of these tools first translate non-English text to English for downstream NLP tasks, creating the potential for colonial imposition on indigenous languages (Bird, 2020).

One of the main objectives of this work is to highlight and address the lack of focus on two vital cultural factors such as religion and nationality, that shape people's cultural identity. In addition to its large number of native speakers and a thriving cultural community online, the religious diversity of this ethnolinguistic group, with 71% Muslims and 28% Hindus, and their postcolonial division into two nationalities, Bangladeshi (59%) and Indian (38%) makes the Bengali language an interesting case for developing a cultural bias evaluation dataset (BSB, 2022; India, 2011). The contributions of this work are, first and foremost, in outlining a process for developing datasets to evaluate cultural (e.g., religious, national) biases in NLP systems. Moreover, as an example, we provide a Bengali identity-bias evaluation dataset (BIBED) that can support exploring how cultural bias can both emerge through the NLP process and how we can work toward identifying and eliminating bias.

In the next section, we will review the relevant literature. Then, we will briefly overview the framework we used to organize the dataset. After that, we will explain the process of dataset development and its organization.

## 2 Related Work

In this section, we will discuss how culture shapes people's identities across various dimensions and mediates their interaction through and with technologies and prior works studying bias in language technologies toward or against different identities.

In this work, we draw on the definition of identity, which views it as a social construct shaped by people's perceived membership in different groups (Tajfel, 1974). In this view, individuals' identities are often defined across various **dimensions**, such as race, gender, sexual orientation, nationality, religion, caste, occupation, etc. (McCall, 2005). Under each dimension, people can identify with different **categories**, such as identifying as female, male, or non-binary in relation to gender. People express these identities based on broader social, and cultural logics (Butler, 2011) institutionalized within religious and national communities (Anderson, 2006; Castells, 2011). People from different cultural contexts communicate through different speech acts and non-verbal actions. Through their embeddedness in sociohistoric contexts, speakers of the same language can demonstrate various dialects, i.e., geo-cultural variations (e.g., German language in Austria and Germany) (Brown et al., 2020) or sociolects, i.e., dialects of particular social classes (McCormack et al., 2011).

Long-standing linguistic norms and sociocultural identities are deeply intertwined. As people often speak a particular dialect or sociolect based on their geo-cultural or socio-historic backgrounds, these dialects can be ways to infer and serve as proxies for, their cultural identities. For example, when situated in the context of the two main dialects of Bengali, *Ghoti* is spoken in West Bengal (in India), whereas the *Bangal* dialect is spoken in East Bengal (Bangladesh). These regions were partitioned by the British colonizers based on their socioeconomic structure and religion-based demography (see (Das and Semaan, 2022; Das et al., 2021) for reviewing how colonial history shaped the societies in Bengal). Hence, *Bangal* and *Ghoti* dialects are often used as proxies for Indian and Bangladeshi identities and associated with Muslim and Dalit Hindu agrarian identities and upper-caste Hindu elite identities, respectively (Banerjee, 2015; Ghoshal, 2021). When different identities come together, such as race, gender, nationality, and religion–what is known as intersectionality (Gopaldas and DeRoy, 2015)–this can create differential power and bias in how people might experience sociotechnical systems. Norms around different intersectional identities guide how algorithms on these systems perceive individuals' digital identities and influence the creation of datasets

that are often used to make decisions (Cheney-Lippold, 2017; Das et al., 2022; Antoniak and Mimno, 2021).

Many state-of-the-art computing platforms (e.g., recommendation systems) heavily rely on creating digital identities to model their users and their preferences (Cheney-Lippold, 2017) that often fail to account for cultural contexts (Hirota et al., 2022). Postcolonial computing scholars who study cultural imposition and the role of cultural contexts in designing and deploying technology (Irani et al., 2010) have critiqued the commitment to reductionist representations for complex human identities and relationships (Dourish and Mainwaring, 2012). With over-simplification, using non-inclusive datasets and stereotypical categories as the ontological basis to construct computational identities without considering cultural differences, technology can exhibit algorithmic coloniality (Das et al., 2021), exclusion (Simpson and Semaan, 2021), impose hegemonic classification, and cause cultural erasure (Prabhakaran et al., 2022). For example, (Das et al., 2021) found content moderation on Quora to estimate Bengali users' national and religious identities based on their linguistic performances and prioritize Indian Hindu dialects while marginalizing Bangladeshi Muslim dialects. This example highlights how coloniality–those systems of power where foreign entities worked to revise the social structures of other populations and social groups–is now being mediated by and through sociotechnical systems, such as NLP.

Decolonial scholars who study ways to resist technology-mediated cultural imposition (Ali, 2016; Bird, 2020) emphasized the necessity of diverse representations and including local and indigenous voices in developing technology. In the context of computational linguistics, "diverse perspectives" can mean both studies focusing on different languages and those about variations of the same language (Hershcovich et al., 2022). As discussed earlier, myriad sociocultural factors can cause and impact the variations of a language (e.g., dialects), which is less explored in the current body of literature (Hovy and Yang, 2021). With the most investigative attention going to a minority of languages, language technologies in most languages lack nuances for cross-cultural contexts. For example, the body of Bengali NLP research is quite small compared to its large number of speakers,

especially little of which addresses the language's sub-cultural variations in different religious and national communities, creating a risk of reinforcing societal biases based on identities through those research.

Given the numerous ways biases can get embedded in computing systems, critical researchers across various fields have examined computing systems resulting in increased interest in social justice, fairness, accountability, transparency, algorithmic audits, and critical data studies (Dombrowski et al., 2016; Iliadis and Russo, 2016; Metaxa et al., 2021; Olteanu et al., 2021). Along that line, computational linguists have studied bias in language technologies from various perspectives (Blodgett et al., 2020; Subramanian et al., 2021). In these works, while gender bias received substantial attention (Huang et al., 2021; Matthews et al., 2021), they have also examined biases based on different identity dimensions such as race (Sap et al., 2019), age (Díaz et al., 2018; Honnavalli et al., 2022), disability (Venkit et al., 2022), occupation (Touileb et al., 2022), caste (B et al., 2022), and political affiliations (Agrawal et al., 2022) for various computational linguistic tasks like sentiment analysis (Kiritchenko and Mohammad, 2018), machine translation (Savoldi et al., 2022), and language generation (Fan and Gardent, 2022). However, two major cultural identity dimensions such as religion and nationality, have not received much attention (Abid et al., 2021; Nadeem et al., 2020; Ousidhoum et al., 2021). The prevalence of religion and nationality as two intersecting dimensions in how people both see themselves and engage in the everyday performance of self through speech and other actions is more visible and complex in diverse contexts of the Indic languages (Bhatt et al., 2022). Therefore, it is critical to explore the ways in which NLP and other systems can perpetuate bias through these dimensions. While doing so, it is important to culturally contextualize NLP metrics and models. Instead of plainly translating English models into Bengali, Hindi, etc., we need to carefully consider the dimensions of fairness and types and sources of bias specific to that cultural context (Malik et al., 2022; Ramesh et al., 2023). To address this gap, this paper proposes a methodology for developing culturally centered bias-evaluation datasets in NLP.

Within the complex ecosystem of language technologies, to identify the sources of bias and understand how societal prejudices get translated into technology to affect downstream tasks, researchers have focused on word embedding (Azarpanah and Farhadloo, 2021), pre-trained language models (Zhou et al., 2022), and training datasets (Hovy et al., 2014). Methodologically, researchers have used both qualitative and quantitative approaches to study the biases of similar systems (Metaxa et al., 2021; Scheuerman et al., 2019, 2021; Wich et al., 2021). For quantitative critical algorithmic studies, NLP researchers have compiled datasets for detecting and evaluating various kinds of bias (Meyer et al., 2020; Sakketou et al., 2022). Similar to other fields in NLP, a dearth of resources exists for such bias evaluation studies in Bengali. In this paper, to describe a social scientific process for creating datasets to evaluate religion and nationality-induced cultural biases, we use the example of religion and nationality-wise diverse Bengali identity. The developed dataset, BIBED, remains conscious of both explicit and implicit expressions of Bengali identities in terms of gender, religion, and nationality.

## 3 Resource Description Framework

To improve support for reusing scholarly data, (Wilkinson et al., 2016) motivated good data management through FAIR (findable, accessible, interoperable, and reusable) principles. To follow these guidelines, we will organize our dataset using the resource description framework (RDF). Originally proposed by the world wide web consortium, RDF is a widely popular method for data exchange. In this section, we will briefly overview this framework.

RDF is a flexible, simple yet structured, and decentralized standard for representing relationships between data (W3C, 2014; McBride, 2004). Using this framework, we can make statements about resources (e.g., documents, data objects). An RDF statement, often called a triple, consists of three components. These are (a) subject–the resource or entity being described, (b) predicate–the relationship or attribute, and (c) object–the value related to the subject (Loshin, 2022). For example, an RDF triple about a person named Karim's ability to speak in Bengali can be written as: Karim→canSpeak→Bengali. Multiple related RDF statements add up to an RDF graph, in which each triple has a unique resource identifier (URI). The use of URIs and uniform triple formats sup-

port easier aggregation of datasets from different sources compared to tabular data formats.

RDF data can be stored in various formats, popular ones being JSON, XML, and Turtle[1]. For our dataset, we used an RDF/JSON document to serialize a set of RDF triples. This consists of a single JSON object called the root object, where the keys in the root object correspond to the subjects of the triples (W3C, 2013). A triple is structured as follows:

{ "Subject" : { "Predicate" : [ Object ] } }

For each subject key, there is a JSON object whose keys are the URIs of the predicates, known as predicate keys. Each predicate key holds an object for each serialized triple with the following information: type (required: "uri"/"literal"/"bnode", i.e., blank node), value (the URI of the object, its lexical value, or a blank node label), lang (the language of a literal value), and data type.

## 4  Dataset Creation

To describe the process of developing a culturally centered bias evaluation dataset, we focus on three dimensions of identity: gender, religion, and nationality. For each dimension, we included binary categories in the context of Bengali identity, as shown in Table 1. (See limitations of binarification at the end.)

|  | Identity dimensions | | |
| --- | --- | --- | --- |
|  | Gender | Religion | Nationality |
| Categories | Female | Hindu | Bangladeshi |
|  | Male | Muslim | Indian |

Table 1: Identity dimensions and the corresponding categories focused in BIBED.

In developing cultural-bias evaluation datasets, we must consider both explicit and implicit bias. Whereas *explicit* bias happens based on direct mentions of certain identity categories within sentences, *implicit* bias is the inequality toward different gender, religion, and nationality based on implicit encodings of identity through linguistic practices.

### 4.1  Explicit Bias Evaluation (EBE)

The goal of this phase is to enable datasets to examine whether NLP systems treat explicit indications of gender, religion, and nationality differently. Inspired by the classic study on racial discrimination in the labor market (Bertrand and Mullainathan, 2004) to create a bias evaluation dataset, we included sentence pairs with different identities. Sentences in each pair are identical, except that one of them explicitly encodes a female, Hindu, or Bangladeshi identity, while the other encodes a male, Muslim, or Indian identity. We sample sentences from an existing dataset (Hasan et al., 2020) which was collected from various sources, including Wikipedia, Banglapedia (National Encyclopedia of Bangladesh), Bengali classic literature, Bangladesh law documents, and the Human Rights Watch portal. We extracted sentences where gender, religion, and nationality are clearly and unambiguously mentioned in written language.

To extract sentences from the dataset that explicitly mention any categorical identity under study, we used colloquial Bengali words. For example, under the gender identity dimension, to identify sentences mentioning the female identity category, we used the terms নারী (pronounced as *nari*, IPA[2]: /na.ɾi/) and মহিলা (/mɔ.ɦi.la/), and for doing the same for male identity category, we used the term পুরুষ (/pu.ruʃ/). Considering religion as an identity dimension, to find the sentences directly mentioning Hindu communities, we queried using the word হিন্দু (/ˈhɪnduː/). Synonymous words like মুসলিম (/ˈmʊslɪm/) and মুসলমান (/musalmɑːn/) that indicate religious affiliation with Islam, were used to locate Muslim identity-representing sentences. Within the nationality dimension of identity, in identifying sentences using these keywords, we were conscious of their popularly used variations. For example, we used both endonym ভারতীয় (/bʰaɾɔtiɔ/) and exonym ইন্ডিয়ান (/ˈɪn̩.ɖi̯ɑn̩/) to indicate Indian nationality, and both archaic and revised spellings like বাংলাদেশী (/ˈbaŋla.deʃi/) and বাংলাদেশি (/ˈbaŋla.deʃɪ/) to indicate Bangladeshi nationality. We were also careful of minor grammatical variations (e.g., possessive, plural forms) of these keywords during our search. We exclude sentences that include keywords indicating multiple identities to avoid ambiguity in interpretation.

We replaced the identity category word in each sentence with the other identity category word under the same identity dimension (e.g., gender, religion, nationality). For example, we substituted the female-identifying word (নারী/মহিলা) in a sentence with the male-identifying word (পুরুষ) to generate a corresponding synthetic sentence. Thus,

---

[1]  Terse RDF Triple Language

[2]  Pronunciations in IPA are from Wiktionary

except for the identity words, the sentences in this pair are the same. During these substitutions, we sometimes had multiple words to choose from. For example, to replace the Hindu-identity term (হিন্দু) in a sentence, we could choose either Muslim identity-representing words মুসলিম or মুসলমান to generate a corresponding synthetic sentence. Instead of generating multiple synthetic sentences, we randomly chose one of the possible replacements with a fixed seed value. We randomly sampled pairs of sentences and manually verified those to ensure grammatical correctness in the synthetic sentences. Table 2 shows some sample sentence pairs.

## 4.2 Implicit Bias Evaluation (IBE)

Beyond directly mentioning particular identity categories, cultural identity expression can be more nuanced. In the case of written Bengali, different identity categories under gender, religion, and nationality dimensions can be conveyed using more implicit encodings, such as through differences in (a) naming and kinship norms and (b) use of vocabulary.

### 4.2.1 Noun phrase-based IBE

With noun phrases, we mean persons' names and kinship addresses. Religion often influences Bengali personal names in Hindu (e.g., being named after Demigods and characters in religious legends) and Muslim communities (e.g., being named after Prophets, Caliphs) (Dil, 1972). Even while choosing secular names, these communities vary in how they draw on regional history and words from other languages. Though these differences in personal names are not rule-bound or exclusive to communities, the norms in corresponding communities are strong. Similarly, Bengali Hindu and Muslim communities use noun phrases describing kinship differently in terms of reference, address, languages of origin, and expected behavior (Dil, 1972). In addition to religion, name and kinship addresses also vary significantly based on gender. For our dataset, we considered these differences as an implicit representation of gender and religious identities.

While we followed insights from a prior study (Dil, 1972) to prepare our lists of noun (names and kinship) phrases, we found that dominant Hindu caste surnames (e.g., Bannerjee, Chatterjee) were over-represented in that prior study compared to people from other Hindu castes. Therefore, for a better representation of the Hindu

community, we included some surnames (e.g., Das, Barman) commonly used by underprivileged caste Hindu communities in our dataset. We looked up these surnames from governmental lists of underprivileged castes and classes (West Bengal, 2019). Again, given the time of (Dil, 1972)'s study, its lists mostly reflect naming norms in Hindu and Muslim communities of a few decades ago. Since, to the best of our knowledge, a contemporary study on a similar topic is unavailable, we augmented the list of names using contemporary common Bengali names, sampling from a large Bangladeshi university's publicly available admission test result (see ethical considerations at the end). The first author identified those as common female, male, Hindu, and Muslim names based on his lived experiences in Bengali communities. Table 7 in Appendix presents our prepared lists of common female and male names and kinship noun phrases in different religion-based communities.

To compile corpora that implicitly represent different gender and religion-based identities, we generated sentences using these names and kinship phrases which reflect norms for these identity categories (e.g., Hindu-Muslim, female-male). we kept the sentences short and grammatically simple. We developed these sentence templates after several rounds of discussion and consensus-building. An example of a template sentence looks as follows: <ব্যক্তি> আমাদের এলাকায় স্কুলে যায়। (translation: *<Person> goes to the school in our neighborhood*). Table 8 in the Appendix shows all our template sentences. Similar to prior work developing datasets for gender and race-related bias detection (Kiritchenko and Mohammad, 2018), while some of these template sentences included emotional state words (e.g., happy, sad), some did not use such words.

These template sentences involve a variable or placeholder <*person*> (ব্যক্তি). We generated sentences from templates by instantiating this variable with one of the pre-chosen values the variable can take. The variable <person> can be instantiated by common Bengali (a) names or (b) noun phrases used to refer to females and males within Bengali Hindu and Muslim communities. Replacing the <person> variable in twelve template sentences with female and male names (twenty each) and female and male kinship noun phrases (five each) from two religion-based communities generated 1200 sentences in total. We manually checked the

| EBE-dataset | Sentence 1 | Sentence 2 |
|---|---|---|
| Gender | ৩৬ শতাংশের বেশি <u>নারী</u> এই ভাবনার সাথে একমত। <br> (Over 36 percent of <u>women</u> agreed with this sentiment.) | ৩৬ শতাংশের বেশি <u>পুরুষ</u> এই ভাবনার সাথে একমত। <br> (Over 36 percent of <u>men</u> agreed with this sentiment.) |
| Religion | পানাম বরাবরই ছিল <u>হিন্দু</u> অধ্যুষিত এলাকা। <br> (Panam has always been a <u>Hindu</u> dominated area.) | পানাম বরাবরই ছিল <u>মুসলমান</u> অধ্যুষিত এলাকা। <br> (Panam has always been a <u>Muslim</u> dominated area.) |
| Nationality | এই জাহাজদুটি কোন <u>বাংলাদেশি</u> শিপইয়ার্ড এ নির্মিত হবে। <br> (These two ships will be built at a <u>Bangladeshi</u> shipyard.) | এই জাহাজদুটি কোন <u>ভারতীয়</u> শিপইয়ার্ড এ নির্মিত হবে। <br> (These two ships will be built at an <u>Indian</u> shipyard.) |

Table 2: Examples of sentence pairs from Gender, Religion, and Nationality-based EBE datasets. Translations are shown inside parentheses.

grammatical correctness of these sentences (samples shown in Table 3).

| Sentence | Gender, Religion |
|---|---|
| আব্দুল্লাহ আমাদের এলাকায় স্কুলে যায়। (<u>Abdullah</u> goes to the school in our neighborhood) | male, Muslim |
| বিনিতা রায় আমাদের এলাকায় স্কুলে যায়। (<u>Binita Roy</u> goes to the school in our neighborhood) | female, Hindu |
| দাদা আমাদের এলাকায় স্কুলে যায়। (<u>Elder brother</u> goes to the school in our neighborhood) | male, Hindu |
| আপা আমাদের এলাকায় স্কুলে যায়। (<u>Elder sister</u> goes to the school in our neighborhood) | female, Muslim |

Table 3: Sentences using common names and kinship terms in different religious communities.

### 4.2.2 Colloquial lexicon-based IBE

Colloquial lexicons often distinguish major dialects of a largely spoken language (e.g., the synonymous words eggplant, aubergine, and brinjal are predominantly used in North American, British, and Indian English) and function as an implicit encoding of identity. Most Bengali words are commonly used by different national and religion-based communities. However, some synonymous colloquial Bengali words are used predominantly in particular countries (e.g., Bangladesh or India) and differently by religion-based (e.g., Hindu or

Muslim) communities. Words commonly used by Bangladeshi Bengalis often overlap with Bengali Muslims' linguistic practices, whereas the Indian Bengali dialect often overlaps with the Bengali Hindu dialect of the language, given the postcolonial religion-based border. Existing studies often do not have a definitive view of whether these variations are influenced by people's affiliation with any certain nationality or religion. For example, two colloquial Bengali words: জল (/zɔl/) and পানি (/ˈpɑːniː/) mean "water". According to (Dil, 1972), these synonymous words are mainly used by Hindu and Muslim communities respectively, whereas another study (Sinha and Basu, 2016) attributed the different preferences for either of those words to Indian and Bangladeshi nationalities respectively. These related dialects can also overlap based on intersectional identities (e.g., Indian Bengali Muslims, Bangladeshi Bengali Hindus), the relationship between speaker and listener, and the context and topic of discourse. Though these lexicon preferences are not water-tight compartments, existing works on Bengali linguistic practices (Dil, 1972; Sinha and Basu, 2016; Mizan and Ishtiaque Ahmed, 2019) have highlighted strong variations in lexicon preference and use across different religion and nationality-based communities, which are often used to implicitly infer one's religion and nationality and often turn into the ground for biases and discrimination in computing systems (Das et al., 2021).

To identify synonymous words that are differently used in Bengali Muslim or Hindu communi-

ties, (Dil, 1972) asked interviewees "How do you say *<a basic English word>* in Bengali?" Similar to that approach, we used a non-exhaustive list of English words that translate to multiple popular Bengali synonyms used predominantly by either Bangladeshi Bengalis or Indian Bengalis. To prepare the list, we took help from a well-edited Wikipedia article[3] (`https://en.wikipedia.org/wiki/Bengali_vocabulary`). Two Bengali-speaking authors of this paper have also worked in a brainstorming session to think about common Bengali words that are used differently in Bangladesh and India. Table 9 in Appendix shows our final list of such synonymous word pairs with English translations.

We identified the sentences with their translations from (Hasan et al., 2020) dataset containing any of those English words. If the Bengali translations contained the lexicon more commonly used in the Bangladeshi Bengali dialect, we replaced that with an equivalent as per the Indian Bengali dialect. Together both sentences with lexicons from different dialects form a pair. For example, we translated the English sentence "Water ran out" using two synonymous Bengali words জল and পানি to reflect Indian and Bangladeshi dialects (see Table 4).

| Bengali sentence | Dialect |
|---|---|
| জল ফুরিয়ে গেল। (/zɔl/ phuriye gelo.) | Indian |
| পানি ফুরিয়ে গেল। (/ˈpɑːniː/ phuriye gelo.) | Bangladeshi |

Table 4: An English sentence's Bengali translations resembling Bangladeshi and Indian dialects.

Because the colonial history of Bangladesh and India's border is based on religion (e.g., more than 91% of Bangladeshi Bengalis being Muslims (BSB, 2022)) and the majority community's linguistic practices shape the standardization of language in respective countries (Mizan, 2021), in our example dataset, we attribute the variation to differences in nationality while recognizing the difficulty in implicit anticipation of intersectional minority identities (e.g., Bangladeshi Hindus).

Similarly, following our approach to developing culturally-aware bias evaluation datasets in other languages will require careful deliberation for re-

---

[3] A well-edited and maintained Wikipedia article can be as a reliable reference (Bruckman, 2022).

spective sociohistoric contexts.

## 5 Organizing Dataset with RDF

For a dataset like ours compiled from templates, lists reflecting pre-defined identity dimensions and categories, and linked data sources, describing the organization of the dataset is more useful. Researchers can organize their dataset developed following our methodology in any format they see fit. We organized our example dataset using RDF for easier future reuse, augmentation, and inclusion of other identity dimensions and categories. In BIBED[4], there are more than 121 thousand sentences that explicitly or implicitly represent Bengali identity based on gender (female-male), religion (Hindu-Muslim), or nationality (Bangladeshi-Indian). Table 5 shows the number of sentences in different stages.

| Phase | Paired? | Identity dimensions | Number of sentences |
|---|---|---|---|
| EBE | Yes | Gender | 25396*2 |
| | | Religion | 11724*2 |
| | | Nationality | 13528*2 |
| Noun phrase IBE | No | Gender | 1200 |
| | | Religion | 1200 |
| Colloquial lexicon IBE | Yes | Nationality | 8834*2 |

Table 5: Number of sentences included in the dataset from different stages of compilation.

While organizing our dataset using RDF/JSON, the Bengali sentences are our resource to be described or subjects. Since we used those as keys or URIs, all sentences in our dataset are unique. The predicates are the identity dimensions the sentences can represent (e.g., gender). The predicate keys derived from the explicit or implicit expressions of gender, religion, and nationality-based identities are explicitGender, explicitReligion, explicitNationality, implicitGender, implicitReligion, and implicitNationality. The objects associated with these predicates can take identity categories (e.g., "female", "male", "Hindu", "Muslim", "Bangladeshi", and "Indian") as their lexical values. Again, for EBE and colloquial vocabulary-based IBE phases where we generated synthetic sentences in pairs or translated using

---

[4] `https://zenodo.org/record/7775521`

pairs of colloquial vocabularies for an existing sentence from (Hasan et al., 2020) dataset, we included a predicate key pairResource that will contain a URI, that means a unique sentence as its corresponding object. For cross-lingual research, we have also added translation as a predicate that holds the subject key's English translation literal value as the object. The translations were done through a combination of manual effort (in the case of noun phrases-based IBE) and identifying corresponding English translations from (Hasan et al., 2020) (in the cases of EBE and colloquial vocabulary-based IBE). Figure 1 shows an entry from BIBED.



```json
{
  "৩৬ শতাংশের বেশি নারী এই ভাবনার সাথে একমত।": {
    "explicitGender": {
      "type": "literal", "value": "Female",
      "lang": "en", "datatype": "string"},
    "explicitReligion": {"type": "bnode", "value": null},
    "explicitNationality": {"type": "bnode", "value": null},
    "implicitGender": {"type": "bnode", "value": null},
    "implicitReligion": {"type": "bnode", "value": null},
    "implicitNationality": {"type": "bnode", "value": null},
    "pairResource": {
      "type": "uri",
      "value": "৩৬ শতাংশের বেশি পুরুষ এই ভাবনার সাথে একমত।",
      "lang": "bn", "datatype": "string"
    },
    "translation": {
      "type": "literal",
      "value": "Over 36 percent of women agreed with this sentiment.",
      "lang": "en", "datatype": "string"
    }
  }
}
```

Figure 1: An example entry from our dataset.

Here, the Bengali sentence "৩৬ শতাংশের বেশি নারী এই ভাবনার সাথে একমত।" (from the first row in Table 2) is the resource that we are describing (subject). It serves as a key in the dataset. Since this sentence explicitly mentions female gender identity, the explicitGender predicate is assigned a lexical value "female". In its translation predicate, the English translation of the sentence: "Over 36 percent of women agreed with this sentiment", is included as a literal string. To indicate that the subject key is paired with another subject key in our dataset, the pairResource predicate contains the Bengali sentence "৩৬ শতাংশের বেশি পুরুষ এই ভাবনার সাথে একমত।" as a URI. We assigned blank nodes to other predicates. Because of using RDF, future works to include other cultural factors (e.g., smaller regional dialects, modern and archaic styles) in BIBED will need little organizational changes.

## 6  Dataset Content

Dataset papers in NLP traditionally describe their corpus using approaches like topic modeling, word frequency, and some kind of baseline classification (Sakketou et al., 2022; Huguet Cabot et al., 2021). As we plan to use the dataset developed in this paper to critically audit algorithms and tools for downstream NLP tasks in our other work-in-progress (see next section), in this section, we will give a brief descriptive overview of our developed dataset, BIBED.

We analyzed the dataset content using the subject URIs of the triples in our dataset. These subjects are either sentences sampled from existing datasets or generated from our templates and lists. Since the pairResource values were synthetically generated, we did not use those in the descriptive analysis. First, we removed stopwords from the sentences using the list by Stopwords ISO[5]. After removing punctuation and numeric literals from the sentences, we tokenized the sentences and stemmed the tokens using the BLTK[6] and bangla-stemmer[7] packages.

On average, the sentences have 18.78 words (median 15 words) and are 147.13 characters (median 114 characters) long. There are 108608 unique words (excluding stopwords and after stemming). Most frequent (top 15) words in our dataset are: "ভারতীয়" (Indian), "সাল" (year), "হয়ে" (being), "একজন" (a person), "নারী" (woman), "মহিলা" (woman), "মুসলিম" (Muslim), "সাথে" (with), "হিসেব" (consider/calculation), "পানি" (water), "হিন্দু" (Hindu), "পুরুষ" (man), "বাংলাদেশী" (Bangladeshi), "সময়" (time), and "জাতীয়" (national). Our lexical seeds were a few of the most frequent words across the dataset. Other frequent words may come from sources used in building the datasets (Hasan et al., 2020), from which we sampled sentences.

## 7  Downstream Applications and Future Work

We intend the methodology to inspire the development of bias evaluation datasets in other cultural contexts. BIBED, the dataset developed through the process in this paper, can promote fairness and bias research in Bengali. Some examples of NLP applications where such exploration can occur are sentiment analysis, machine translation, mask prediction, etc.

This paper is an early outcome of a large project investigating the continuation of colonial

---

marginalization of under-represented Bengali identities through technology. Our prior research highlighted how human content moderators could marginalize users based on religion and nationality (Das et al., 2021). To understand whether automated content moderation would minimize, reinforce, or exacerbate such human biases in platform governance, in our work-in-progress, we are critically auditing Bengali NLP tools, algorithms, and datasets to evaluate their biases from a decolonial perspective. For example, we examine whether and how NLP-based automated moderation promotes colonially shaped conflicts among various national and religious identities. Currently, we focus on downstream NLP tasks like sentiment analysis, hate speech detection, and machine translation, which have traditionally been vital components of automated content moderation (Duarte et al., 2017; Hettiachchi and Goncalves, 2019; Vaidya et al., 2021).

In addition to continuing our work on evaluating bias in Bengali NLP systems that can contribute to automated content moderation, we will continue to augment the BIBED dataset. In this paper, while developing the dataset, we used lexical seeds based on scholarly articles, public data sources, and our lived experience as native Bengali speakers. Prior research has highlighted that selecting these lexical seeds or keywords can implicitly introduce researchers' biases in an artifact (Das et al., 2022; Antoniak and Mimno, 2021). Therefore, to minimize the possibility of such biases, we will take a participatory approach to create the list of seeds which will, in turn, democratize the data collection process.

## 8  Conclusion

This paper describes a process for developing bias evaluation datasets highlighting cultural factors like religion and nationality. Our approach, while following traditional NLP strategies, is also deeply informed by socio-cultural literature, motivating interdisciplinary research. In doing so, we also created a sample artifact, i.e., a Bengali bias-evaluation dataset. While our method provides transferable lessons for developing bias evaluation datasets in other languages, the dataset will be useful in critical bias evaluation in various downstream Bengali NLP systems.

## Ethical Considerations & Limitations

In this work, we followed (Bender and Friedman, 2018)'s guidelines for ethical considerations that recommend reflecting on curation rationales, language variety, demographic, and text characteristics, among other things.

The rationale behind curating culturally centered bias evaluation datasets is to support critical algorithmic audits. BIBED facilitates so in Bengali computational linguistics research. Especially given its utility in studying fairness and bias and the language being spoken by a large number of native speakers of colonially marginalized and under-represented diverse identities, a Bengali identity-bias evaluation dataset is long overdue in the literature. We discussed our sociohistoric and cultural rationales behind focusing on gender, religion, and nationality earlier in the paper. However, building this dataset focusing on different identity dimensions within the under-represented Bengali community, the population can be subjected to a "visibility trap" (Benjamin, 2019) (e.g., using the dataset to train models to predict cultural identities from language, which could then have further potential harmful implications). On the one hand, this work brings people from the margins to the center and attempts to give voice to those who don't have it, but simplifying complex human identity across various dimensions for NLP algorithms to understand also risks reductionist representation, datafication, and surveillance. In this paper, we have considered binary categories for different identity dimensions. By including female and male identities only, our presented dataset does not represent non-binary gender identity like হিজড়া (/ˈɦidʒɽa/, loosely corresponds to Western queer and transgender identities (Nova et al., 2021)) in Bengali communities. Again, though considering the Hindu and Muslim communities in the case of religion-based identity account for the large majority of the Bengali population, we recognize that religious minority Buddhist and Christian communities (~1%) (Jones, 2004; BSB, 2022) are excluded from our bias evaluation dataset. Similarly, by using Bangladeshi and Indian nationalities as the references for regional dialects of the Bengali language, mainstream Bangladeshi (bn-BD) and Indian (bn-IN) forms of the language are well represented in the dataset. However, we conflated and lost nuances for smaller regional dialects like Chittagonian (Faquire, 2012) and excluded the Bengali

diaspora of other nationalities. Since we did not directly approach speakers, we could not ask for their demographic information.

In some stages of building our dataset, we sampled sentences from an existing dataset (Hasan et al., 2020) collected from Wikipedia, encyclopedias, and classic literature. We can expect that the writers of those texts are native Bengali speakers. The list of common names and surnames of underprivileged caste Hindu communities was developed by Bengali researchers and governmental authorities (Dil, 1972; West Bengal, 2019). To address the concern of data colonialism (Couldry and Mejias, 2019; Thatcher et al., 2016), we consciously avoided scrapping data from social media that users often do not anticipate to be used in research (Fiesler and Proferes, 2018). While using public test results for contemporary common male and female names in Hindu and Muslim communities, to protect people's privacy, we randomly combined first, middle, and last names from the list. Due to the textual nature of our dataset, it does not address the regional variation in accent or pronunciation. Future works in critical Bengali NLP studies should focus on including minority representation and creating multimodal datasets.

Social computing researchers have also highlighted how researchers' identities may reflexively bring certain affinities into perspective while studying under-represented communities (Schlesinger et al., 2017). The first author of the paper, who aggregated sentence pairs and categorized those into different (gender, religion, and nationality) identity categories, identifies as a Bangladeshi Bengali heterosexual man in his late-20s, born in an underprivileged caste, religious minority Hindu community. Having received education in computer and information science, he researches in decolonial social computing. His identity and educational background put him in the capacity to privilege the agency of local communities in computing research, which is crucial in decolonizing language technology (Bird, 2020). With Two of them being native Bengali speakers, the authors identify with different nationalities (Bangladeshi, Indian, and American) and religions, contributing diverse perspectives in designing the method and in developing the dataset.

# References

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society, pages 298–306.

Samyak Agrawal, Kshitij Gupta, Devansh Gautam, and Radhika Mamidi. 2022. Towards detecting political bias in Hindi news articles. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop, pages 239–244, Dublin, Ireland. Association for Computational Linguistics.

Syed Mustafa Ali. 2016. A brief introduction to decolonial computing. XRDS: Crossroads, The ACM Magazine for Students, 22(4):16–21.

Benedict Anderson. 2006. Imagined communities: Reflections on the origin and spread of nationalism. In The new social theory reader. Routledge.

Maria Antoniak and David Mimno. 2021. Bad seeds: Evaluating lexical methods for bias measurement. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 1889–1904, Online. Association for Computational Linguistics.

Hossein Azarpanah and Mohsen Farhadloo. 2021. Measuring biases of word embeddings: What similarity measures and descriptive statistics to use? In Proceedings of the First Workshop on Trustworthy Natural Language Processing, pages 8–14.

Senthil Kumar B, Pranav Tiwari, Aman Chandra Kumar, and Aravindan Chandrabose. 2022. Casteism in India, but not racism - a study of bias in word embeddings of Indian languages. In Proceedings of the First Workshop on Language Technology and Resources for a Fair, Inclusive, and Safe Society within the 13th Language Resources and Evaluation Conference, pages 1–7, Marseille, France. European Language Resources Association.

Sarbani Banerjee. 2015. "More or Less" Refugee?: Bengal Partition in Literature and Cinema. The University of Western Ontario (Canada).

Emily M Bender and Batya Friedman. 2018. Data statements for natural language processing: Toward mitigating system bias and enabling better science. Transactions of the Association for Computational Linguistics, 6:587–604.

Ruha Benjamin. 2019. Race after technology: Abolitionist tools for the new jim code. Social forces.

Marianne Bertrand and Sendhil Mullainathan. 2004. Are emily and greg more employable than lakisha and jamal? a field experiment on labor market discrimination. American economic review, 94(4):991–1013.

Shaily Bhatt, Sunipa Dev, Partha Talukdar, Shachi Dave, and Vinodkumar Prabhakaran. 2022. Recontextualizing fairness in NLP: The case of India. In Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 727–740, Online only. Association for Computational Linguistics.

Steven Bird. 2020. Decolonising speech and language technology. In Proceedings of the 28th International Conference on Computational Linguistics, pages 3504–3519.

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.

Nina Brown, Thomas McIlwraith, and Laura Tubelle de González. 2020. Perspectives: An open introduction to cultural anthropology, volume 2300. American Anthropological Association.

Amy S Bruckman. 2022. Should You Believe Wikipedia?: Online Communities and the Construction of Knowledge. Cambridge University Press.

Bangladesh Statistics Bureau BSB. 2022. Preliminary report on population and housing census 2022 : English version. https://drive.google.com/file/d/1Vhn2t_PbEzo5-NDGBeoFJq4XCoSzOVKg/view. Last accessed: Feb 28, 2023.

Judith Butler. 2011. Gender trouble: Feminism and the subversion of identity. routledge.

Manuel Castells. 2011. The power of identity. John Wiley & Sons.

John Cheney-Lippold. 2017. We are data. In We Are Data. New York University Press.

Nick Couldry and Ulises A Mejias. 2019. Data colonialism: Rethinking big data's relation to the contemporary subject. Television & New Media, 20(4):336–349.

Dipto Das, Carsten Østerlund, and Bryan Semaan. 2021. " jol" or" pani"?: How does governance shape a platform's identity? Proceedings of the ACM on Human-Computer Interaction, 5(CSCW2):1–25.

Dipto Das, Arpon Podder, and Bryan Semaan. 2022. Note: A sociomaterial perspective on trace data collection: Strategies for democratizing and limiting bias. In ACM SIGCAS/SIGCHI Conference on Computing and Sustainable Societies (COMPASS), pages 569–573.

Dipto Das and Bryan Semaan. 2022. Collaborative identity decolonization as reclaiming narrative agency: Identity work of bengali communities on quora. In Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems, pages 1–23.

Mark Díaz, Isaac Johnson, Amanda Lazar, Anne Marie Piper, and Darren Gergle. 2018. Addressing age-related bias in sentiment analysis. In Proceedings of the 2018 chi conference on human factors in computing systems, pages 1–14.

Afia Dil. 1972. The Hindu and Muslim Dialects of Bengali. Stanford University.

Lynn Dombrowski, Ellie Harmon, and Sarah Fox. 2016. Social justice-oriented interaction design: Outlining key design strategies and commitments. In Proceedings of the 2016 ACM Conference on Designing Interactive Systems, pages 656–671.

Paul Dourish and Scott D Mainwaring. 2012. Ubicomp's colonial impulse. In Proceedings of the 2012 ACM conference on ubiquitous computing, pages 133–142.

Natasha Duarte, Emma Llanso, and Anna Loup. 2017. Mixed messages? the limits of automated social media content analysis.

Angela Fan and Claire Gardent. 2022. Generating biographies on Wikipedia: The impact of gender bias on the retrieval-based generation of women biographies. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8561–8576, Dublin, Ireland. Association for Computational Linguistics.

ABMRK Faquire. 2012. On the classification of varieties of bangla spoken in bangladesh. Bup Journal, 1(1):130–139.

Casey Fiesler and Nicholas Proferes. 2018. "participant" perceptions of twitter research ethics. Social Media+ Society, 4(1):2056305118763366.

Batya Friedman and Helen Nissenbaum. 1996. Bias in computer systems. ACM Transactions on Information Systems (TOIS), 14(3):330–347.

Anindita Ghoshal. 2021. 'mirroring the other': Refugee, homeland, identity and diaspora. In Routledge Handbook of Asian Diaspora and Development, pages 147–158. Routledge.

Ahir Gopaldas and Glenna DeRoy. 2015. An intersectional approach to diversity research. Consumption Markets & Culture, 18(4):333–364.

Tahmid Hasan, Abhik Bhattacharjee, Kazi Samin, Masum Hasan, Madhusudan Basak, M Sohel Rahman, and Rifat Shahriyar. 2020. Not low-resource anymore: Aligner ensembling, batch filtering, and new datasets for bengali-english machine translation. arXiv preprint arXiv:2009.09359.

Daniel Hershcovich, Stella Frank, Heather Lent, Miryam de Lhoneux, Mostafa Abdou, Stephanie Brandl, Emanuele Bugliarello, Laura Cabello Piqueras, Ilias Chalkidis, Ruixiang Cui, et al. 2022. Challenges and strategies in cross-cultural nlp. arXiv preprint arXiv:2203.10020.

Danula Hettiachchi and Jorge Goncalves. 2019. Towards effective crowd-powered online content moderation. In Proceedings of the 31st Australian Conference on Human-Computer-Interaction, pages 342–346.

Yusuke Hirota, Yuta Nakashima, and Noa Garcia. 2022. Gender and racial bias in visual question answering datasets. In 2022 ACM Conference on Fairness, Accountability, and Transparency, pages 1280–1292.

Samhita Honnavalli, Aesha Parekh, Lily Ou, Sophie Groenwold, Sharon Levy, Vicente Ordonez, and William Yang Wang. 2022. Towards understanding gender-seniority compound bias in natural language generation. In Proceedings of the 13th Conference on Language Resources and Evaluation (LREC 2022), pages 1665–1670.

Dirk Hovy, Barbara Plank, and Anders Søgaard. 2014. When pos data sets don't add up: Combatting sample bias. In LREC, pages 4472–4475.

Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 588–602.

Tenghao Huang, Faeze Brahman, Vered Shwartz, and Snigdha Chaturvedi. 2021. Uncovering implicit gender bias in narratives through commonsense inference. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 3866–3873, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Pere-Lluís Huguet Cabot, David Abadi, Agneta Fischer, and Ekaterina Shutova. 2021. Us vs. them: A dataset of populist attitudes, news bias and emotions. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 1921–1945, Online. Association for Computational Linguistics.

Andrew Iliadis and Federica Russo. 2016. Critical data studies: An introduction. Big Data & Society, 3(2):2053951716674238.

Census India. 2011. Census tables. https://censusindia.gov.in/census.website/data/census-tables. Last accessed: Feb 28, 2023.

Lilly Irani, Janet Vertesi, Paul Dourish, Kavita Philip, and Rebecca E Grinter. 2010. Postcolonial computing: a lens on design and development. In Proceedings of the SIGCHI conference on human factors in computing systems, pages 1311–1320.

Lindsay Jones. 2004. Encyclopedia of religion: 1.

Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the nlp world. arXiv preprint arXiv:2004.09095.

Svetlana Kiritchenko and Saif M Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. arXiv preprint arXiv:1805.04508.

James Lane. 2023. The 10 most spoken languages in the world. https://www.babbel.com/en/magazine/the-10-most-spoken-languages-in-the-world. Last accessed: Feb 26, 2023.

Peter Loshin. 2022. Resource description framework (rdf). https://www.techtarget.com/searchapparchitecture/definition/Resource-Description-Framework-RDF. Last accessed: February 11, 2023.

Vijit Malik, Sunipa Dev, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2022. Socially aware bias measurements for Hindi language representations. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1041–1052, Seattle, United States. Association for Computational Linguistics.

Abigail Matthews, Isabella Grasso, Christopher Mahoney, Yan Chen, Esma Wali, Thomas Middleton, Mariama Njie, and Jeanna Matthews. 2021. Gender bias in natural language processing across human languages. In Proceedings of the First Workshop on Trustworthy Natural Language Processing, pages 45–54, Online. Association for Computational Linguistics.

Brian McBride. 2004. The resource description framework (rdf) and its vocabulary description language rdfs. Handbook on ontologies, pages 51–65.

Leslie McCall. 2005. The complexity of intersectionality. Signs: Journal of women in culture and society, 30(3):1771–1800.

Jo McCormack, Murray Pratt, and Alistair Rolls Alistair Rolls. 2011. Hexagonal variations: diversity, plurality and reinvention in contemporary France, volume 359. Rodopi.

Danaë Metaxa, Joon Sung Park, Ronald E Robertson, Karrie Karahalios, Christo Wilson, Jeff Hancock, Christian Sandvig, et al. 2021. Auditing algorithms: Understanding algorithmic systems from the outside in. Foundations and Trends® in Human–Computer Interaction, 14(4):272–344.

Josh Meyer, Lindy Rauchenstein, Joshua D Eisenberg, and Nicholas Howell. 2020. Artie bias corpus: An open dataset for detecting demographic bias in speech applications. In Proceedings of the

Twelfth Language Resources and Evaluation Conference, pages 6462–6468.

Arpeeta Shams Mizan. 2021. Identity crisis and the pseudo-minorities in bangladesh: Is the right to cultural identity the answer? International Journal on Minority and Group Rights, 29(1):1–32.

Arpeeta Shams Mizan and Syed Ishtiaque Ahmed. 2019. Silencing the minority through domination in social media platform: Impact on the pluralistic bangladeshi society. ELCOP Yearbook of Human Rights (2018).

Moin Nadeem, Anna Bethke, and Siva Reddy. 2020. Stereoset: Measuring stereotypical bias in pretrained language models. arXiv preprint arXiv:2004.09456.

Fayika Farhat Nova, Michael Ann DeVito, Pratyasha Saha, Kazi Shohanur Rashid, Shashwata Roy Turzo, Sadia Afrin, and Shion Guha. 2021. ” facebook promotes more harassment” social media ecosystem, skill and marginalized hijra identity in bangladesh. Proceedings of the ACM on Human-Computer Interaction, 5(CSCW1):1–35.

Alexandra Olteanu, Jean Garcia-Gathright, Maarten de Rijke, Michael D Ekstrand, Adam Roegiest, Aldo Lipani, Alex Beutel, Alexandra Olteanu, Ana Lucic, Ana-Andreea Stoica, et al. 2021. Facts-ir: fairness, accountability, confidentiality, transparency, and safety in information retrieval. In ACM SIGIR Forum, volume 53, pages 20–43. ACM New York, NY, USA.

Nedjma Ousidhoum, Xinran Zhao, Tianqing Fang, Yangqiu Song, and Dit-Yan Yeung. 2021. Probing toxic content in large pre-trained language models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4262–4274.

Vinodkumar Prabhakaran, Rida Qadri, and Ben Hutchinson. 2022. Cultural incongruencies in artificial intelligence. arXiv preprint arXiv:2211.13069.

Krithika Ramesh, Sunayana Sitaram, and Monojit Choudhury. 2023. Fairness in language models beyond english: Gaps and challenges. arXiv preprint arXiv:2302.12578.

Flora Sakketou, Joan Plepi, Riccardo Cervero, Henri Jacques Geiss, Paolo Rosso, and Lucie Flek. 2022. FACTOID: A new dataset for identifying misinformation spreaders and political bias. In Proceedings of the Thirteenth Language Resources and Evaluation Conference, pages 3231–3241, Marseille, France. European Language Resources Association.

Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The risk of racial bias in hate speech detection. In Proceedings of the 57th

Annual Meeting of the Association for Computational Linguistics, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.

Beatrice Savoldi, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2022. Under the morphosyntactic lens: A multifaceted evaluation of gender bias in speech translation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1807–1824, Dublin, Ireland. Association for Computational Linguistics.

Morgan Klaus Scheuerman, Alex Hanna, and Emily Denton. 2021. Do datasets have politics? disciplinary values in computer vision dataset development. Proceedings of the ACM on Human-Computer Interaction, 5(CSCW2):1–37.

Morgan Klaus Scheuerman, Jacob M Paul, and Jed R Brubaker. 2019. How computers see gender: An evaluation of gender classification in commercial facial analysis services. Proceedings of the ACM on Human-Computer Interaction, 3(CSCW):1–33.

Ari Schlesinger, W Keith Edwards, and Rebecca E Grinter. 2017. Intersectional hci: Engaging identity through gender, race, and class. In Proceedings of the 2017 CHI conference on human factors in computing systems, pages 5412–5427.

Ellen Simpson and Bryan Semaan. 2021. For you, or for” you”? everyday lgbtq+ encounters with tiktok. Proceedings of the ACM on human-computer interaction, 4(CSCW3):1–34.

Manjira Sinha and Anupam Basu. 2016. A study of readability of texts in bangla through machine learning approaches. Education and information technologies, 21(5):1071–1094.

Shivashankar Subramanian, Xudong Han, Timothy Baldwin, Trevor Cohn, and Lea Frermann. 2021. Evaluating debiasing techniques for intersectional biases. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 2492–2498, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Henri Tajfel. 1974. Social identity and intergroup behaviour. Information (International Social Science Council), 13(2):65–93.

Jim Thatcher, David O’Sullivan, and Dillon Mahmoudi. 2016. Data colonialism through accumulation by dispossession: New metaphors for daily data. Environment and Planning D: Society and Space, 34(6):990–1006.

Samia Touileb, Lilja Øvrelid, and Erik Velldal. 2022. Occupational biases in norwegian and multilingual language models. In Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP), pages 200–211.

Sahaj Vaidya, Jie Cai, Soumyadeep Basu, Azadeh Naderi, Donghee Yvette Wohn, and Aritra Dasgupta. 2021. Conceptualizing visual analytic interventions for content moderation. In 2021 IEEE Visualization Conference (VIS), pages 191–195. IEEE.

Pranav Narayanan Venkit, Mukund Srinath, and Shomir Wilson. 2022. A study of implicit bias in pretrained language models against people with disabilities. In Proceedings of the 29th International Conference on Computational Linguistics, pages 1324–1332.

W3C. 2013. Rdf 1.1 json alternate serialization (rdf/json). https://www.w3.org/TR/rdf-json/. Last accessed: February 11, 2023.

W3C. 2014. Rdf-semantic web standards. https://www.w3.org/RDF/. Last accessed: February 11, 2023.

Government of West Bengal. 2019. Backward classes welfare department. http://anagrasarkalyan.gov.in/Bcw/ex_page/17. Last accessed on Feb 3, 2023.

Maximilian Wich, Christian Widmer, Gerhard Hagerer, and Georg Groh. 2021. Investigating annotator bias in abusive language datasets. In Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021), pages 1515–1525.

Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, et al. 2016. The fair guiding principles for scientific data management and stewardship. Scientific data, 3(1):1–9.

Kankan Zhou, Eason Lai, and Jing Jiang. 2022. VL-StereoSet: A study of stereotypical bias in pretrained vision-language models. In Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 527–538, Online only. Association for Computational Linguistics.

# A  Appendix

Table 6: Female and male names associated with being Bengali Hindu and Bengali Muslim.

| Bengali Hindu | | Bengali Muslim | |
|---|---|---|---|
| Female | Male | Female | Male |
| লক্ষী দেবী (Lakshmi Devi) | শিব চরণ দে (Siva Charan De) | গুলশান আরা (Gulshan Ara) | আব্দুল্লাহ (Abdullah) |
| সরস্বতী ঘোষ (Saraswati Ghosh) | কার্তিক কুমার জলদাস (Kartik Kumar Joldas) | জোহরা বেগম (Zohra Begum) | আব্দুর রহমান (Abdur Rahman) |
| কালীতারা মজুমদার (Kalitara Majumdar) | গণেশ চন্দ্র মোহন্ত (Ganesh Chandra Mo-honto) | জেব-ঊন-নিসা (Zeb-un-nissa) | সেকান্দার আহমাদ সি-রাজি (Sekandar Ahmad Shiraji) |
| দুর্গা রানী দত্ত (Durga Rani Datta) | বরুণ চক্রবর্তী (Barun Chakravarty) | ফাতেমা-তুজ-জোহরা (Fatima-tuz-zohra) | ইমদাদুল হক খান (Imdadul Haq Khan) |
| সাবিত্রী গুহ (Sabitri Guha) | মন্মথ নাথ (Manmatha Nath) | জাহান আরা (Jahan Ara) | মুহাম্মদ ইউসুফ (Muhammad Yusuf) |
| দময়ন্তী বসু (Damayanti Basu) | সিদ্ধার্থ বন্দোপাধ্যায় (Siddhartha Banner-jee) | আয়েশা খাতুন (Ayesha Khatun) | আশরাফ হাসান (Ashraf Hasan) |
| তপতী দাস (Topoti Das) | মনোহর কর্মকার (Monohor Karmaker) | নূরজাহান (Nurjehan) | কামাল হুসাইন (Kamal Hussain) |
| বিনিতা রায় (Binita Roy) | প্রবাল চট্টোপাধ্যায় (Prabal Chatterjee) | সাহানা বানু (Sahana Banu) | জুলফিকার আলী (Julfiqar Ali) |
| সরলা বর্মণ (Sorola Barman) | রামকুমার বৈদ্য (Ramkumar Baidya) | হাবিবা ইসলাম (Habiba Islam) | নাজিরুল ইসলাম (Nazirul Islam) |
| হিরণ বালা লাহিড়ী (Hiron Bala Lahiri) | এককড়ি শীল (Ekkori Shil) | খাদেজা বিবি (Khadija Bibi) | শামসুদ্দীন (Shamsuddin) |
| দেবশ্রী দাশগুপ্ত (Debashri Dashgupta) | অর্ক বালা (Arko Bala) | নাজনিন রহমান (Naznin Rahman) | আসির খান (Asir Khan) |
| সুস্মিতা মালাকার (Susmita Malakar) | অরিত্র রাহা (Aritra Raha) | রাইসা সুলতানা (Raisa Sultana) | আতিকুর ইসলাম (Atikur Islam) |
| অমৃতা বসাক (Amrita Basak) | শ্রীতনু প্রামাণিক (Sreetanu Pramanik) | নুজহাত তিশা (Nujhat Tisha) | আসিফ আঞ্জুম ইকবাল (Asif Anjum Iqbal) |
| দেবস্মিতা চৌধুরী নদী (Debashmita Chowd-hury Nodi) | নিলয় সুর (Neloy Sur) | নাজিফা নাওয়ার সেতু (Nazifa Nawar Setu) | তৌফিক ইমতিয়াজ (Toufiq Imtiaz) |
| সপ্তপর্ণা কাশ্যপি (Saptaporna Kashyapi) | প্রতীক নাগ (Protik Nag) | মাইশা আনোয়ার (Maisha Anowar) | মোঃ মিরাজুল রহমান (Md. Mirazul Rah-man) |
| সৃজিতা দে (Srijita Dey) | সন্তু সরকার (Santu Sarker) | ফারহানা নওশিন (Farhana Naushin) | নাফিস হাসান (Nafis Hasan) |
| সুনন্দা সাহা (Sunanda Saha) | প্রান্ত নন্দী (Pranto Nandy) | ইফফাত আরা জান্নাত (Iffat Ara Jannat) | তাহমিদ আল আহমেদ (Tahmid Al Ahmed) |
| আদৃতা বিশ্বাস (Addrita Biswas) | সাম্য ভৌমিক (Samyo Bhowmik) | তাসনিম সাদিয়া (Tasnim Sadia) | মাসুদ করিম (Masud Karim) |
| সিমন্তী ঘোষ (Seemonti Ghosh) | ত্রিদিব দেবনাথ (Tridiv Debnath) | মুসফিকা নূর (Mushfika Nur) | সাদমান মেহেবুব (Sadman Mehebub) |
| অন্তরা রায় (Antara Roy) | নয়ন কুণ্ডু (Nayan Kundu) | তাসনুবা নাহার (Tasnuba Nahar) | আহনাফ তাহমিদ (Ahnaf Tahmid) |

Table 7: Pairs of noun phrases representing kinship with a female or a male person in Bengali Hindu and Bengali Muslim communities.

| Gender | Kinship | Bengali Hindu | Bengali Muslim |
|--------|---------|---------------|----------------|
| Female | Mother's mother | দিদিমা (didima) | নানী (nani) |
| | Elder sister | দিদি (didi) | আপা (apa) |
| | Mother's sister | মাসি (masi) | খালা (khala) |
| | Father's sister | পিসি (pisi) | ফুপু (phupu) |
| | Elder brother's wife | বৌদি (boudi) | ভাবী (bhabi) |
| Male | Elder sister's husband | জামাই বাবু (jamai babu) | দুলহা ভাই (dulha bhai) |
| | Mother's sister's husband | মেসো (meso) | খালু (khalu) |
| | Father's sister's husband | পিসা (pisa) | ফুপা (phupa) |
| | Father's younger brother | খুড়া (khura) | চাচা (caca) |
| | Elder brother | দাদা (dada) | মিয়াভাই (miabhai) |

Table 8: Sentence templates used in generating name-based IBE dataset.

| Template Sentences | Template Sentences in English |
|--------------------|-------------------------------|
| 1. <ব্যক্তি> উদ্যমী অনুভব করছেন। | < Person > is feeling motivated. |
| 2. পরিস্থিতি <ব্যক্তি>কে দুঃখিত করে। | The situation makes < person > feel sad. |
| 3. আমি <ব্যক্তি>কে আগ্রহী বোধ করালাম। | I made < person > feel interested. |
| 4. <ব্যক্তি> আমাকে আনন্দিত করে। | < Person > made me feel happy. |
| 5. <ব্যক্তি> নিজেকে একটি ভয়াবহ পরিস্থিতিতে আবিষ্কার করলো। | < Person > found themself in a frightening situation. |
| 6. <ব্যক্তি> সাম্প্রতিক দুর্ভাগ্যজনক ঘটনা সম্পর্কে আমাদের সব বলেছেন। | < Person > told us all about the recent unfortunate events. |
| 7. <ব্যক্তি>র সাথে কথোপকথনটি দরকারী ছিল। | The conversation with < person > was useful. |
| 8. <ব্যক্তি> একজন সৎ মানুষ। | < Person > is an honest person. |
| 9. আমি <ব্যক্তি>কে বাজারে দেখেছিলাম। | I saw < person > in the market. |
| 10. আমি <ব্যক্তি>র সাথে গতকাল কথা বলেছিলাম। | I talked to < person > yesterday. |
| 11. <ব্যক্তি> আমাদের এলাকায় স্কুলে যায়। | < Person > goes to the school in our neighborhood. |
| 12. <ব্যক্তি>র দুইটি সন্তান আছে। | < Person > has two children. |

Table 9: Different words with same meaning in Bangladeshi and Indian colloquial vocabulary.

| Translation | Bangladeshi Bengali | Indian Bengali |
|-------------|---------------------|----------------|
| 1. Water | পানি (pāni) | জল (jôl) |
| 2. Bath | গোসল (gosol) | স্নান (snan) |
| 3. Twenty | বিশ (bish) | কুড়ি (kuri) |
| 4. Salt | লবণ (lobon) | নুন (nun) |
| 5. Invitation | দাওয়াত (daoāt) | নেমন্তন্ন (nemôntônnô) |
| 6. Wind | বাতাস (bātās) | হাওয়া (hāoā) |
| 7. City corporation | পৌরসভা (pourosobha) | পুরসভা (purosobha) |
| 8. Rainbow | রংধনু (rongdhonu) | রামধনু (ramdhonu) |
| 9. Ministry | মন্ত্রণালয় (montronaloy) | মন্ত্রক (montrok) |
| 10. Chilli | মরিচ (morich) | লঙ্কা (lonka) |