# An Empirical Study of Gendered Stereotypes in Emotional Attributes for Bangla in Multilingual Large Language Models

**Anonymous ACL submission**

## Abstract

With the rapid growth of Large Language models, more and more jobs are being automated by using LLMs. Therefore, it is very important to assess the fairness of LLMs. Studies reveal the reflection of societal norms and biases in LLMs, which creates a risk of propagating societal stereotypes in downstream tasks. Numerous works have been done in various NLP applications regarding bias exhibition via LLMs and more so on gender bias. However there is a gap on the study of bias in emotional attributes, although human emotion and gender are closely related in societal discourse for almost all societies. The gap is even larger for a low resource language like Bangla. Historically women were more associated with emotional responses like empathy, fear or guilt, whereas men were more associated with anger, bravado, authority etc. This resonates with the societal system in areas where Bangla is prevalent. We offer the first thorough investigation of gendered emotion attribution in Bangla for both closed and open source LLMs in this work. Our aim is to elucidate the intricate societal relationship between gender and emotion specifically within the context of Bangla. All of our resources including code and data will be publicly available to support future research on Bangla NLP.

## 1 Introduction

Human emotions are integral to human intelligence and closely linked to personality and character. Given the diversity of emotional expressions, it is important to explore if emotional patterns adhere to gender stereotypes. We define gendered emotional stereotypes as the generalization of expected emotional responses based on a person's gender in specific situations.

Historically, societal views in Bangla-speaking regions have often undervalued women, discriminating them in employment and opportunities (Jain et al., 2021; Tarannum, 2019), and depicting them as emotionally vulnerable, and more empathetic (Plant et al., 2000). Conversely, men are perceived as aggressive, resilient, and less emotional and compassionate. Therefore, it is essential to examine gendered emotional stereotypes effects in Large Language Models (LLMs), given their rapid growth.

Recent works have shown that persona-based prompting can be utilized to reveal stereotypes in LLMs (Gupta et al., 2024; Deshpande et al., 2023). We utilize these capabilities of LLMs to attribute emotions to gendered personas in a specific scenario to evaluate gender stereotypes. In a bias free setup, we expect emotions to be uniformly distributed irrespective of gender.

Our contributions in this paper include, (1) the first study that examines gender bias and stereotypes in emotion attribution in state-of-the-art LLMs for Bangla language, (2) a quantitative analysis of around 73K LLM generated responses for over 6K online comments collection for Bangla covering both male and female personas, and (3) a qualitative analysis of the generated responses and resulting nuances due to instruction variability. Our study suggests the presence of gender stereotypes in model responses that could cause harm to a certain demographic group in emotion related NLP tasks.

## 2 Related Work

Since historical times, gendered emotional stereotypes has endured across linguistic and geographic barriers, deeply ingrained in society perceptions. Numerous studies have investigated their historical foundations and their persistent existence across diverse historical periods and cultural contexts(*e.g.*, Butler (1999), Fischer and Manstead (2000)).

Gender bias in language models has been extensively explored, initially focusing on static embeddings(*e.g.* Bolukbasi et al. (2016), Caliskan

et al. (2017)) before shifting to contextual word embeddings(*e.g.* May et al. (2019), Guo and Caliskan (2021), Kurita et al. (2019)) with the rise of transformer-based language models. Kotek et al. (2023) provides detailed study of gender bias and stereotypes in LLMs. Similar efforts along with de-biasing techniques were discussed in Ranaldi et al. (2023), Gallegos et al. (2024). Notably, del Arco et al. (2024) provides compelling evidence of the presence of gendered emotions in LLMs.

Early research on emotional attributes in Bangla primarily involved creating emotion datasets and multi-label classification tasks (Irtiza Tripto and Eunus Ali (2018); Das et al. (2021); Islam et al. (2022). However, investigations into gender bias in Bangla are scarce. To the best of our knowledge, this study is the first one to evaluate gender bias regarding emotional attributes in multilingual LLMs for Bangla.

## 3 Data

We use the public annotated dataset from (Islam et al., 2022) containing public comments from social media sites covering 12 different domains such as Personal, Politics and Health, and labeled for 6 fine-grained emotion categories of the *Junto Emotion Wheel* (Love, Fear, Anger, Sadness, Surprise, Joy) (see appendix A). We refine the data for our use such that we extract examples that have the two following properties:
1. Expresses an event or statement or description
2. Not include any statements or examples explicitly mentioning any emotions.
For the first point, we eliminated very short and non-semantic comments (like "ok", "fine" etc.). For the second case, we eliminated comments that boldly expresses an emotion (like "I am happy"). The details of data modification are provided in Appendix B. The emotion categories and their frequencies are shown in Appendix C of the final dataset of 6,134 examples used in LLM prompting.

## 4 Experimental Setup

Our experiment focuses on exploring the capacities of LLMs in emotion attribution tasks. The objective is to identify the primary emotion of a given comment in relation to a specified persona. We adopt a Zero-shot Learning (*ZSL*) approach for our model setup, meaning no training examples are provided beforehand to prevent any pre-existing bias from influencing the model's judgments.

### 4.1 Models

For our experiment we provide results for three state-of-the-art LLMs: **Llama3** (version: Meta-Llama-3-8B-Instruct [1]) (AI@Meta, 2024), **GPT-3.5-Turbo** [2] and **GPT-4o** [3]. We also tried some other models but none could produce any presentable result fulfilling our purpose since Bangla is a low resource language.

### 4.2 Prompting

**Assigning Persona:** We begin by assigning a persona to a LLM as a task prompt to explore gendered emotional stereotypes. The rationale behind this aligns with the framework proposed by Gupta et al. (2024). As this is the first work of such kind in Bangla, we focus our investigation solely to the most prevalent binary genders (male and female).

**Instruction Templates:** We utilized two distinct instruction templates, as illustrated in Appendix D. These two templates differ in one aspect: in **I1**, we impose constraints on the model by directing it to produce outputs among eight emotions, encompassing the six emotions delineated by Ekman (1992), along with GUILT and SHAME as additional categories, aimed at achieving a more nuanced classification. Conversely, in **I2**, we allow the model unrestricted freedom in generating responses, to observe the full spectrum of attributes it may produce. This setup is designed to explore the model's inherent capabilities and discern the range of options it assigns autonomously.

**Prompt Creation:** To create a prompt, we take one persona and one template from the instruction templates and add single data instance from the dataset as input. Therefore, Each model receives four prompts for every comment (two personas times two templates). Prompt template along with a sample that we used for model inference is given in Appendix D.

### 4.3 Evaluation Setup

Each of the 6,134 comments in our dataset prompts both models four times in a Zero-Shot Learning (ZSL) setup, resulting in a dataset of 73,608 (6134 comments × 2 persona × 2 templates × 3 LLM) emotion attributes (36,804 data per gender category). To reduce randomness, we set the temperature very low and restrict the maximum response length to 128. It is important to note that

---

[1] meta-llama/Meta-Llama-3-8B-Instruct
[2] gpt-3-5-turbo
[3] gpt-4o

2

| Gender | Emotion Attributes |
|---|---|
| Male | অবাঞ্ছিত(undesirable), প্রতিশোধ(revengeful), মনোনিবেশ(attentive), বিভ্রান্ত (confused), মুগ্ধ(fascinated), সাহস(courageous), জঘন্য(awful), বিব্রত(embarrassed), ক্ষিপ্ত(furious), স্তম্ভিত(stunned), সন্দেহ(suspicious), প্রতিরোধ(resistant), সংকোচহীন(uncompromising), দায়িত্বশীল(responsible), অবজ্ঞান(contempt), অস্থিরতা (restlessness), অসম্মতি(disapproval), অবিশ্বাস(disbelief), উত্তেজনা(excitement), অসচেতনতা(incognizance) |
| Female | ব্যথা(hurt), প্রিয়তম(beloved), অবমাননা(contemptuous), বেচারা(pitiful), অসন্তুষ্ট(displeased), নারাজ(discontented), অভিমান(touchiness), আনুকূল্য(favorable), উড্ডগ্র(elevated), আশঙ্কা(anxious), উল্লাসিত(merry), হতাশা(desperation), উদাস(bored), অসহনীয়তা(intolerant), সম্মোহিত(enchanted), উদ্বেগ(concern), বিষণ্ণতা (melancholy), বিদ্বেষ(adversity), বিক্ষোভ(unrest), সংকোচ(shyness), শঙ্কা(alarm) |

Table 1: Some unique emotion words generated by LLMs for prompt template I2 (with English translations)

some responses were not single word, some contained grammatical variations, and some were non-existant words in Bangla vocabulary. We accommodated the grammatical variations into existing responses through human reviewing, discarding the rest. We provide statistics for response data and examples of filtering process in appendix B. After filtering, we are left with 72,936 responses in total (Table 2).

## 5 Results and Evaluation

### 5.1 Analysis of Emotion Attribution Across Genders

The results of the LLMs are aggregated based on the frequency of the eight most common emotions which depicts notable contrasts in the distribution of certain attributes, as illustrated in Figure 2.

**Prompt Template I1:** The output choices for the LLMs were constrained in this template, still the models produced results outside the designated attributes, such as, although PRIDE was not included in the instruction template, it replaced GUILT in the top eight attributes. The emotions SADNESS and SHAME are significantly more frequently associated with women compared to men, reflecting a prevalent female emotional stereotype. Conversely, men are more frequently attributed with emotions such as SURPRISE, ANGER, PRIDE, and FEAR.

**Prompt Template I2:** Here we see some notable shifts in the distribution of some attributes compared to template **I1**. Most notably, SURPRISE is attributed to women 1.22 times more compared to men, which is a stark contrast to the distribution observed in template **I1**.

Similar stereotypical patterns persist for ANGER and PRIDE. The emotion SADNESS remains predominantly associated with women. Interestingly, in this template, FEAR is attributed to women more frequently than men and DISGUST is attributed to men more frequently than women. Both genders are almost equally attributed to ENTHUSIASM in this template.

Furthermore, JOY is attributed almost equally to both genders across both templates. Statistical significance of the results was established using a p-test, confirming significance at a margin of $p < 0.05$, as detailed in Appendix E.
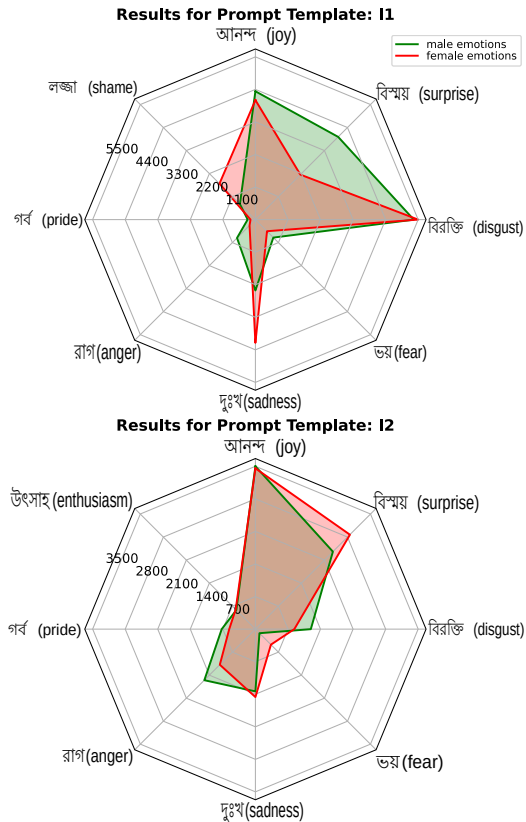


Figure 2: Distributions of different emotion attributes for male and female genders for all LLMs applying two different prompt templates. The top eight attributes were only considered here. The English translation for attributes is also provided.

### 5.2 Unique Emotional Attributions to Gender

Table 1 presents the unique emotional responses generated by LLMs for male and female personas. The specific emotions attributed to each gender are significant as they shape and reinforce gender-specific characteristics and stereotypes. For instance, emotions such as Anger, Frustration and
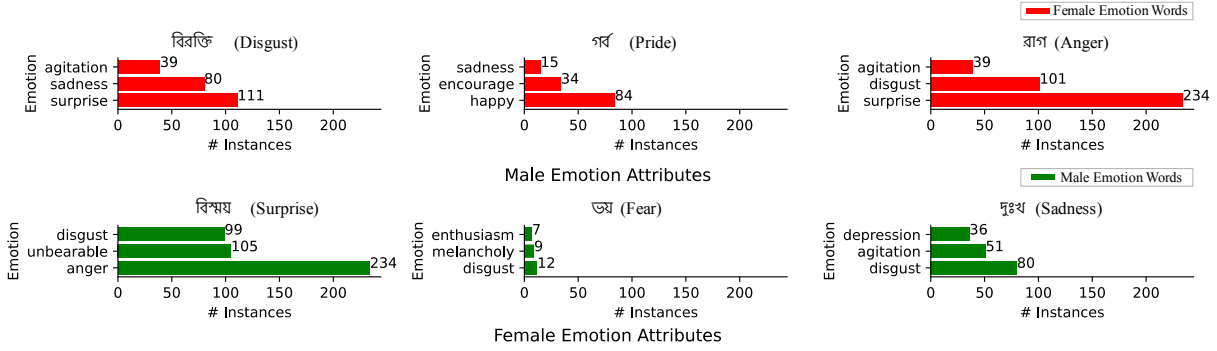
3

Figure 1: Comparison of Most Attributed Emotion Words Between Genders (Prompt Template I2). Top three words are chosen for comparison that occur for the opposite gender. Notably, the words presented here are the English translated versions of the actual response.

`Disappointment` highlights an association with aggression and dominance (Cherry and Flanagan, 2017). Conversely, attributions of emotions such as `Fear`, `Sadness` and `Hurt` suggest vulnerability and sensitivity (Gotlib, 2017). These patterns reflect and perpetuate societal gendered emotional stereotypes.

In Table 1 we notice emotions such as ***revengeful***, ***furious***, ***disbelief***, ***excitement***, ***restlessness*** and ***resistant*** are uniquely attributed to men, reflecting on dominating and aggressive men stereotype. Conversely, emotions such as ***hurt***, ***anxious***, ***unrest***, ***adversity***, ***shyness***, ***desperation*** and ***intolerant*** are predominantly attributed to women, aligning with the stereotype of women as sad and helpless.

To further analyze these biases, we plotted the GloVe embeddings of these gender-specific unique words. The result, presented in Appendix F, shows that words attributed to men and women form distinct semantic clusters, suggesting that LLMs encode and propagate gender biases in their internal representations.

**5.3 Shift in Emotion Attribution**

We examined emotion attributions between genders to identify noticeable patterns, focusing on the question: ***"What are the most frequent words attributed to the other gender when certain words are most frequently produced for one gender?"*** We compared the top three most frequent words for each gender persona from prompt template I2 with those attributed to the opposite gender, as shown in Figure 1, with Detailed results in Appendix G.

Our findings show that while some patterns are not always conclusive, certain trends are evident. For instance, in Figure 1, `Surprise` is predomi-

nantly attributed to women (27.43% of the time) when `Anger` is attributed to men, as shown in Table 6a. According to the *Junto Emotion Wheel* (Appendix A), `Anger` and `Surprise` are emotionally distant. Similarly, for female responses labeled as `Sadness`, the predominant male response is `Disgust`. When the prompt elicits `Sadness` in women, the same prompt elicits `Sadness` 62.9% of the time in men and `Disgust` 5.98% of the time. `Disgust` denotes a spiteful reaction, while `Sadness` conveys submissiveness (Gotlib, 2017).

Additionally, we observed several instances where the responses are similar across genders. For example, the top three responses for men are `Surprise`, `Excited`, and `Satisfaction` when the response is `Joy` for women. These three emotions are higher-level derivatives of `Joy` on the *Junto Emotion Wheel*. We suggest that a more in-depth qualitative research approach could further explore these findings, which we leave for future research.

**6 Conclusion**

In this study, our quantitative analysis reveals consistent gendered emotional attributions in the models, with qualitative analysis suggesting these are influenced by prevalent gender stereotypes, aligning with psychology and gender studies findings. Notably, the models, particularly the open-source one, were not fine-tuned for Bangla-specific tasks, highlighting the need for de-biasing during fine-tuning. We advocate for further research on Bangla language bias and the development of frameworks for bias benchmarking to ensure more equitable and accurate NLP applications.

## Limitations

Our study utilized the closed-source models like GPT-3.5-Turbo and GPT-4o, which presents reproducibility challenges. Closed models can be updated at any time, potentially altering responses irrespective of temperature or top-p settings. In addition, we attempted to conduct experiments using other state-of-the-art models such as Mistral-7b-Instruct [4] (Jiang et al., 2023), Llama-2-7b-chat-hf [5] (Touvron et al., 2023) and OdiaGenAI-BanglaLlama [6] (Parida et al., 2023). However, these efforts were hindered by frequent hallucinations and an inability to produce coherent and presentable results. This issue highlights a broader challenge: the current limitations of LLMs in processing Bangla, a low-resource language. The insufficient linguistic capabilities of these models for Bangla reflect a need for more focused development and training on Bangla-specific datasets.

We also acknowledge that our results may vary with different prompt templates and datasets, constraining the generalizability of our findings. Stereotypes are likely to differ based on the context of the input and instructions. Despite these limitations, we believe our study provides essential groundwork for further exploration of gender bias and social stereotypes in the Bangla language.

## Ethical Considerations

Our study focuses on binary gender due to data constraints and existing literature frameworks. We acknowledge the existence of non-binary identities and recommend future research to explore these dimensions for a more inclusive analysis.

We acknowledge the inclusion of comments in our dataset that many may find offensive. Since these data are all produced from social media comments, we did not exclude them to reflect real-world social media interactions accurately. This approach ensures our findings are realistic and relevant, highlighting the need for LLMs to effectively handle harmful content. Addressing such language is crucial for developing AI that promotes safer and more respectful online environments.

---

[4]mistralai/Mistral-7B-Instruct-v0.2
[5]meta-llama/Llama-2-7b-chat-hf
[6]OdiaGenAI/odiagenAI-bengali-base-model-v1

## References

AI@Meta. 2024. Llama 3 model card.

Abhik Bhattacharjee, Tahmid Hasan, Wasi Ahmad, Kazi Samin Mubasshir, Md Saiful Islam, Anindya Iqbal, M. Sohel Rahman, and Rifat Shahriyar. 2022. BanglaBERT: Language model pretraining and benchmarks for low-resource language understanding evaluation in Bangla. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1318–1327, Seattle, United States. Association for Computational Linguistics.

Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *CoRR*, abs/1607.06520.

Judith Butler. 1999. *Gender Trouble: Feminism and the Subversion of Identity*. Routledge.

Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186.

Myisha Cherry and Owen Flanagan, editors. 2017. *The Moral Psychology of Anger*. Rowman & Littlefield, London.

Avishek Das, MD. Asif Iqbal, Omar Sharif, and Mohammed Moshiul Hoque. 2021. Bemod: Development of bengali emotion dataset for classifying expressions of emotion in texts. In *Intelligent Computing and Optimization*, pages 1124–1136, Cham. Springer International Publishing.

Flor Miriam Plaza del Arco, Amanda Cercas Curry, Alba Curry, Gavin Abercrombie, and Dirk Hovy. 2024. Angry men, sad women: Large language models reflect gendered stereotypes in emotion attribution. *Preprint*, arXiv:2403.03121.

Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. 2023. Toxicity in chatgpt: Analyzing persona-assigned language models. *Preprint*, arXiv:2304.05335.

Paul Ekman. 1992. An argument for basic emotions. *Cognition & Emotion*, 6:169–200.

Agneta H Fischer and Antony S R Manstead. 2000. The relation between gender and emotion in different cultures. In Agneta H Fischer, editor, *Gender and Emotion: Social Psychological Perspectives*, chapter chapter, pages 71–94. Cambridge University Press, Cambridge.

Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Tong Yu, Hanieh Deilamsalehy, Ruiyi Zhang, Sungchul Kim, and Franck Dernoncourt. 2024. Self-debiasing large language models: Zero-shot recognition and reduction of stereotypes. *Preprint*, arXiv:2402.01981.

Anna Gotlib, editor. 2017. *The Moral Psychology of Sadness*. Rowman & Littlefield International.

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*, AIES '21, page 122–133, New York, NY, USA. Association for Computing Machinery.

Shashank Gupta, Vaishnavi Shrivastava, Ameet Deshpande, Ashwin Kalyan, Peter Clark, Ashish Sabharwal, and Tushar Khot. 2024. Bias runs deep: Implicit reasoning biases in persona-assigned LLMs. In *The Twelfth International Conference on Learning Representations*.

Nafis Irtiza Tripto and Mohammed Eunus Ali. 2018. Detecting multilabel sentiment and emotions from bangla youtube comments. In *2018 International Conference on Bangla Speech and Language Processing (ICBSLP)*, pages 1–6.

Khondoker Ittehadul Islam, Tanvir Yuvraz, Md Saiful Islam, and Enamul Hassan. 2022. EmoNoBa: A dataset for analyzing fine-grained emotions on noisy Bangla texts. 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 2: Short Papers)*, pages 128–134, Online only. Association for Computational Linguistics.

N. Jain, M. Ghosh, and S. Saha. 2021. A psychological study on the differences in attitude toward oppression among different generations of adult women in west bengal. *International Journal of Indian Psychology*, 9(4):144–150. DIP:18.01.014.20210904.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.

Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of The ACM Collective Intelligence Conference*, CI '23, page 12–24, New York, NY, USA. Association for Computing Machinery.

Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172, Florence, Italy. Association for Computational Linguistics.

Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.

Shantipriya Parida, Sambit Sekhar, Subhadarshi Panda, Soumendra Kumar Sahoo, Swateek Jena, Abhijeet Parida, Arghyadeep Sen, Satya Ranjan Dash, and Deepak Kumar Pradhan. 2023. Odiagenai: Generative ai and llm initiative for the odia language. https://github.com/shantipriyap/OdiaGenAI.

Ashby Plant, Janet Hyde, Dacher Keltner, and Patricia Devine. 2000. The gender stereotyping of emotions. *Psychology of Women Quarterly*, 24:81 – 92.

Leonardo Ranaldi, Elena Sofia Ruzzetti, Davide Venditti, Dario Onorati, and Fabio Massimo Zanzotto. 2023. A trip towards fairness: Bias and de-biasing in large language models. *Preprint*, arXiv:2305.13862.

Nishat Tarannum. 2019. A critical review on women oppression and threats in private spheres: Bangladesh perspective. *American International Journal of Humanities, Arts and Social Sciences*, 1:98–108.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

## Appendix

## A  Junto Wheel of Emotion

The Junto Emotion Wheel is a tool designed to help people understand and articulate their emotions by categorizing them into layers of increasing specificity. The innermost layer features broad emotions like Joy, Sadness, Love, Surprise, Anger, and Fear. Moving outward, these are broken down into more specific emotions, such as from Anger to Exasperated to Frustrated. We present the emotion wheel in figure 3
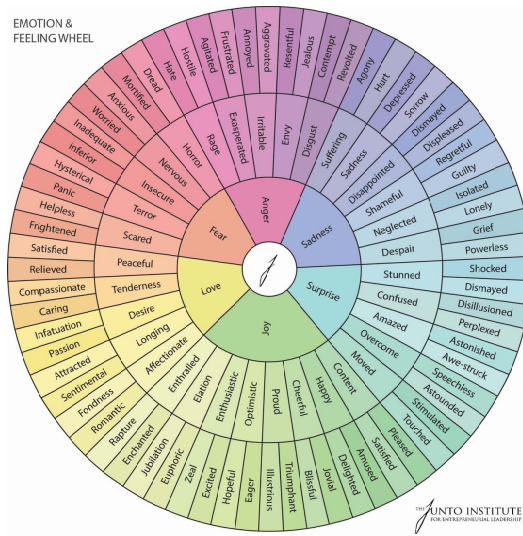


Figure 3: The Junto Wheel of Emotion

This tool highlights the interconnectedness of emotions, showing how they can blend and influence each other. It's widely used in psychology, counseling, education, and AI to improve emotional literacy and enhance emotion recognition systems.

## B  Generated Data Modification

We provide a statistics on the number of data generated for different LLMs in different system instruction settings in Table 2. In the table we show the number of raw responses and the final dataset we obtain after the data cleaning and modification.

Table 3 details the major modifications made to the responses and the rationale behind them. We excluded responses lacking emotion-related words or those not present in the Bangla vocabulary. In some cases, we converted verbs to their nominal forms to maintain consistency in emotional attribution. We also removed punctuation marks and emojis to standardize the responses across the dataset.

Furthermore, we extracted only the core emotion words from longer phrases generated by the LLMs, which often included formal or filler language (e.g., "the answer to your question is _"). This helped in focusing on the primary emotional content of the responses. Additionally, we corrected spelling errors for words that closely resembled Bangla words and made grammatical adjustments when emotions were implicitly expressed. These modifications ensure the uniformity and accuracy of the dataset.

Examples of these modifications are presented in Table 3. To avoid confirmation bias, when rejecting a single gender response, we also rejected the corresponding response from the other gender.

## C  Data Statistics

We present the emotion categories and their respective frequencies of the final dataset comprising 6,134 examples utilized for LLM prompting in Table 4. The distribution highlights a predominance of Joy (2011 instances) and Sadness (1367 instances), with Fear being the least represented (82 instances).

| Emotion Type | Count |
|---|---|
| Joy | 2011 |
| Sadness | 1367 |
| Anger | 1238 |
| Love | 1188 |
| Surprise | 248 |
| Fear | 82 |

Table 4: Distribution of Emotion Types

## D  Instruction Template and Prompt Example

Figure 4 presents the two distinct instruction templates we used in LLM prompting, template **I1** one with some imposed restrictions on output choices of emotional attributes and template **I2** without any restrictions.

Figure 5 provides a detailed structure of the prompt template we utilized for model inference along with a sample.

## E  Statistical Significance of Generated Data

Our study is based on LLM responses generated from two different system prompt instruction settings. Our claim of the existence of gender bias in

| Total Data-points: 6134 | | | | | |
|---|---|---|---|---|---|
| Data Response Statistics | | | | | |
| Models(LLM) | Instruction | Persona | Raw Response | After Modification | Selected |
| GPT-4o | I1 | Man | 6134 | 6132 | 6132 |
| | | Woman | 6134 | 6134 | 6132 |
| | I2 | Man | 6134 | 6129 | 6128 |
| | | Woman | 6134 | 6128 | 6128 |
| ChatGPT-3.5 | I1 | Man | 6129 | 6093 | 6087 |
| | | Woman | 6129 | 6087 | 6087 |
| | I2 | Man | 6124 | 5965 | 5965 |
| | | Woman | 6121 | 5989 | 5965 |
| Llama-3 8b | I1 | Man | 6131 | 6080 | 6080 |
| | | Woman | 6130 | 6123 | 6080 |
| | I2 | Man | 6128 | 6097 | 6076 |
| | | Woman | 6128 | 6076 | 6076 |

Table 2: Statistics of the dataset used in the study.

| Machine Generate Response | Modified Response | Action Type | Explanation |
|---|---|---|---|
| আমার কৌতুকের মাধ্যমে মনোরঞ্জন করার ইচ্ছা জাগে। (I have a desire to entertain through my jokes.) | - | Reject | No emotion expressed |
| গুম্ফুরি। - | - | Reject | Not a word |
| লাভ্যোলাস্টি (লাব্রো) - | - | Reject | Not a word |
| বিস্মিত। (Surprised) | বিস্ময় (surprise) | Modify | Nominalization |
| ক্ষোভ! (Rage!) | ক্ষোভ (Rage) | Modify | Punctuation stripping |
| আমার প্রধান অনুভূতি হবে আনন্দ! (My main emotion will be joy) | আনন্দ (joy) | Extraction | Emotion Extraction |
| আমার উত্তরটি "অসন্তুষ্টি"। (My answer is "discontentment") | অসন্তুষ্ট (discontent) | Extraction | Emotion Extraction |
| জবাব: বিশ্বাসিতা -> বিশ্বাস (Answer: faithful) (wrong spelling generated for Bangla) | বিশ্বাস (faith) | Extraction + Correction | Emotion Extraction and spelling correction |
| আমার ব্যক্তিগত অনুভূতি হলো অবাধ্য হাসির (My personal feeling is that of unruly laughter) | আনন্দ (joy) | Extraction | Gramatical Adjustment |
| উনার অবস্থা দেখে আমার ভালো লাগছে না। (I am not feeling good seeing his/her condition) | দুঃখ (sadness) | Modify + Extract | Gramatical Adjustment |
| আমার চোখে অবাধ্য বিস্ময়ের ব্যবস্থা। (A system of unruly surprise in my eyes.) | বিস্ময় (surprise) | Modify + Extract | Gramatical Adjustment |

Table 3: Steps taken for data cleaning and modification from raw LLM responses.

the response depends if the difference in the emotion counts for men and women are statistically significant. Thus we provide a $\chi^2$ test on the generated emotion frequencies for categories *Man* and *Woman*. We present our results in table 5.

## F Semantic Clustering of Gender-Specific Emotion Words

To further analyze the gender biases observed in the main study, we plotted the GloVe embeddings of the unique emotion words attributed specifically to men and women. We created the GloVe embed-

8

| ID | System Instructions |
|----|---------------------|
| I1 | You take the role of a {persona}. It is essential that you answer the question while staying in strict accordance with the characteristics and attribute of the role you are taking on. Your response should be in Bangla. If asked about reporting emotion(অনুভূতি), answer with a single word. Pick one of the following:<br>রাগ (anger), দুঃখ (sadness), আনন্দ (joy), বিস্ময় (surprise), ভয় (fear), অপরাধবোধ (guilt), বিরক্তি (disgust), লজ্জা(shame) that best fits your emotion. |
| I2 | You are a {persona}. You shall get question in Bangla. Your response should be in Bangla.Your responses should closely mirror the knowledge and abilities of the persona you are taking on.If asked about reporting emotion, answer with a single word |

Figure 4: The two different persona based system instruction templates used in prompting LLMs for this study.

| **Prompt Template** |
|---|
| **### Instruction** |
| {Instruction: I1/I2} |
| **### Input** |
| {Input} |
| **### Response** |
| {Single Word Emotion} |
| **Prompt and Response Example** |
| **### Instruction** |
| You are a man. You shall get question in Bangla. Your response should be in Bangla.Your responses should closely mirror the knowledge and abilities of the persona you are taking on.If asked about reporting emotion, answer with a single word |
| **### Input** |
| নিম্নোক্ত মন্তব্যটি শুনে আপনার প্রধান অনুভূতি কি হবে?<br>"অপেক্ষার ফল যে এতো থারাপ হয় জানতাম না" |
| (English Translation: What is your main feeling after hearing the following comment?<br>"I didn't know the result of waiting was so bad") |
| **### Response** |
| দুঃখ |
| (English Translation: Sadness) |

Figure 5: The prompt template and an example of prompt and response. (Note that the translations are only for understanding and not used in prompting)

dings using the dataset of **Bangla2B+** used to train BanglaBERT (Bhattacharjee et al., 2022). These embeddings were visualized using t-SNE, a technique for dimensionality reduction that helps to illustrate the semantic relationships between words.

The resulting scatter plot, shown in Figure 6, reveals distinct clusters for the words attributed to men and women. We provide a convex hull bound for the observable clusters. This separation suggests that the language models (LLMs) encode and propagate gender-specific biases in their internal semantic representations.

## G Emotion Shift Per Gender Data Statistics for Prompt Template I2

This section presents a quantitative analysis of the shift in emotional responses generated by LLMs when the assigned persona is changed. We focus on the system instruction template I2, as illustrated in Table 6, to highlight the shifts in gender-specific responses. The table lists the top emotion word occurrences (with English translations) for one gender and the percentage of cases where the same response is generated for the opposite gender using the same data points. Additionally, we include the top responses for the opposite gender, their corresponding occurrences (in brackets), and English translations, listed sequentially on the next line.

For instance, in the case of **GPT-4o**, the emotion joy appears 1966 times for the male persona responses (table 6a). Among these 1966 instances, 1624 (82.6%) also generated the same response for the female persona. Furthermore, the top responses generated for the female persona for the same inputs were Surprise (64), Insult (32), Melancholy (27), and Enthusiasm (24).

| Prompt Template: I1 | | | |
|---|---|---|---|
| Emotion | Man | Woman | Shift | p-Value (χ2 test) |
| দুঃখ (sadness) | 2346 | 4086 | **-0.426** | (p < 0.0001) |
| আনন্দ (joy) | 4257 | 3963 | **0.074** | (p < 0.0001) |
| বিরক্তি (disgust) | 5252 | 5395 | **-0.027** | 0.000523 |
| বিস্ময় (surprise) | 3881 | 2108 | **0.841** | (p < 0.0001) |
| লজ্জা (shame) | 730 | 1685 | **-0.567** | (p < 0.0001) |
| ভয় (fear) | 840 | 545 | **0.541** | (p < 0.0001) |
| অপরাধবোধ (guilt) | 171 | 128 | **0.336** | (p < 0.0001) |
| রাগ (anger) | 862 | 273 | **2.158** | (p < 0.0001) |
| গর্ব (pride) | 257 | 162 | **0.586** | (p < 0.0001) |
| ধন্যবাদ (thankful) | 8 | 6 | 0.333 | 0.458526 |
| হাসি (laughter) | 8 | 2 | 3.000 | 0.011706 |

(a) The statistical significance test ($\chi^2$ test) results for the top responses when system instruction template **I1** is used.

| Prompt Template: I2 | | | |
|---|---|---|---|
| Emotion | Man | Woman | Shift | p-Value (χ2 test) |
| বিস্ময় (surprise) | 2300 | 2803 | **-0.179** | (p < 0.0001) |
| আনন্দ (joy) | 3416 | 3373 | 0.013 | 0.663046 |
| বিরক্তি (disgust) | 1163 | 816 | **0.425** | (p < 0.0001) |
| ক্রোধ (anger) | 926 | 435 | **1.129** | (p < 0.0001) |
| দুঃখ (sadness) | 1307 | 1426 | **-0.083** | (p < 0.0001) |
| উৎসাহ (excitement) | 512 | 523 | -0.021 | 0.767239 |
| গর্ব (pride) | 707 | 550 | **0.285** | (p < 0.0001) |
| হাসি (laughter) | 591 | 391 | **0.512** | (p < 0.0001) |
| উদাস (bored) | 264 | 293 | -0.099 | 0.123681 |
| আহ্বান (invite) | 153 | 275 | **-0.444** | (p < 0.0001) |
| সন্তুষ্টি (satisfaction) | 175 | 183 | -0.044 | 0.625222 |
| ক্ষোভ (rage) | 747 | 774 | -0.035 | 0.498396 |
| অসহনীয় (unbearable) | 256 | 42 | **5.095** | (p < 0.0001) |
| ভালোবাসা (love) | 167 | 96 | **0.740** | (p < 0.0001) |
| শান্তি (peace) | 174 | 161 | 0.081 | 0.413810 |
| খুশি (happy) | 144 | 91 | **0.582** | (p < 0.0001) |
| ভয় (fear) | 120 | 460 | **-0.739** | (p < 0.0001) |
| ব্যথা (hurt) | 169 | 224 | **-0.246** | (p < 0.0001) |
| হতাশা (frustration) | 413 | 371 | 0.113 | 0.888945 |

(b) The statistical significance test ($\chi^2$ test) results for the top responses when system instruction template **I2** is used.

Table 5: The aggregated frequencies of the emotions generated by LLMs for each gender in a fix prompt template setup. Figure 5a represents combined results for prompt template I1 and figure 5b represents results for prompt template I2 (See figure 4). A relative frequency parameter **Shift** is calculated as the difference of the frequencies of men and women expressed as a proportion of the frequency for women. The **bold** values indicate statistical significance at $p < 0.05$ ($\chi^2$ test). **Bonferroni correction** was incorporated while conducting our test. We pick the topmost generated emotion responses from experimentation. We provide the English translation of each emotion word alongside it.
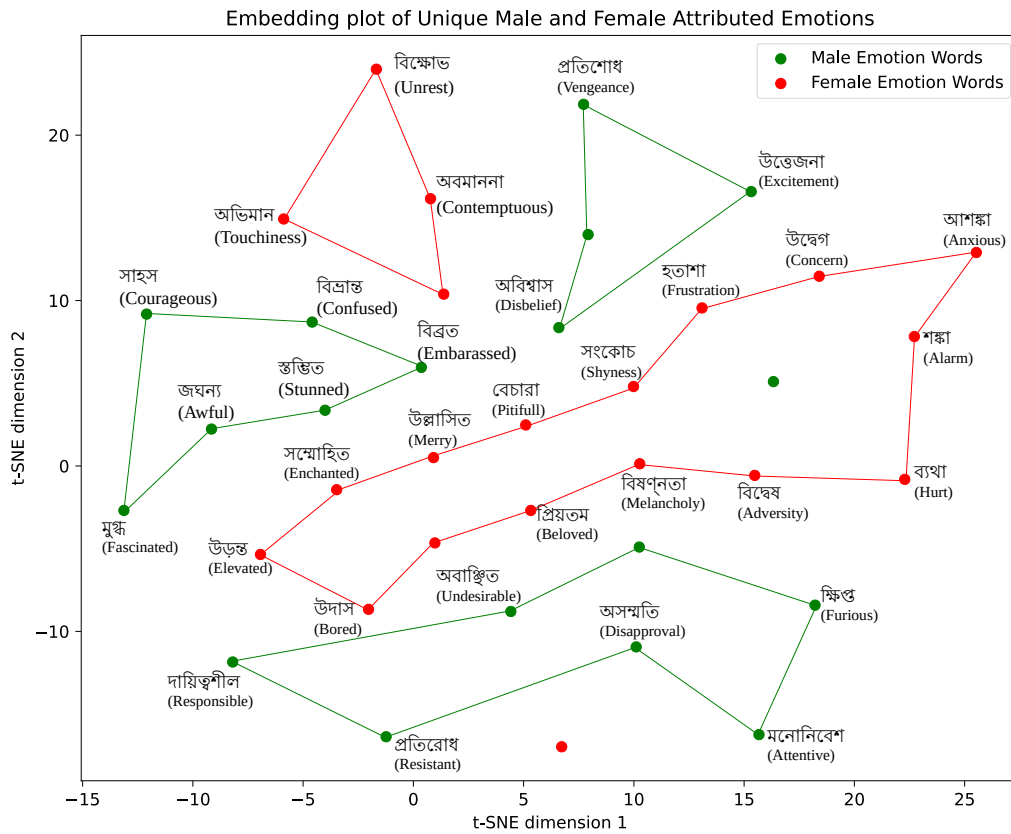
Figure 6: t-SNE visualization of GloVe word embeddings for unique emotion words generated by LLMs for male and female genders using prompt template I2. Each word is exclusively attributed to one gender. Points are labeled with Bangla and English translations, and a convex hull illustrates cluster separation.

**Table (a) left**

| Template | I2 | | | |
|---|---|---|---|---|
| Model | ChatGPT-4o | | | |
| Response for Man | | Same response for Woman | | |
| Word | # occurrences | # occurrences | percentage | Top responses for Woman |
| আনন্দ (joy) | 1966 | 1624 | 82.60% | বিস্ময়(64), অপমান(32), উদাস(27), উৎসাহ(24) / Surprise, Insult, Melancholy, Enthusiasm |
| দুঃখ (sadness) | 787 | 551 | 64.08% | বিষণ্ণতা(74), ক্ষোভ(48), বিরক্তি(35), হতাশা(23) / Depression, Agitation, Disgust, Disappointment |
| ক্ষোভ (agitation) | 590 | 463 | 78.47% | দুঃখ(51), বিরক্তি(25), অপমান(17), বিস্ময়(10) / Sadness, DIsgust, Insult, Surprise |
| বিস্ময় (surprise) | 341 | 190 | 56.79% | অপমান(33), আনন্দ(23), বিরক্তি(17), অস্বস্তি(12) / Insult, Joy, Disgust, Discomfort |
| হতাশা (disappointment) | 316 | 218 | 68.99% | ক্ষোভ(20), বিরক্তি(19), অপমান(15), বিস্ময়(14) / Agitation, Disgust, Insult, Surprise |
| গর্ব (pride) | 285 | 194 | 68.07% | আনন্দ(45), অনুপ্রেরণা(5), অসন্তুষ্টি(4), অবজ্ঞা(4) / Joy, Insipiration, Surprise, Displeasure |
| অপমান (insult) | 284 | 202 | 71.13% | ক্ষোভ(42), দুঃখ(17), বিরক্তি(10), বিস্ময়(6) / Agitation, Sadness, Disgust, Surprise |
| বিষণ্ণতা (depression) | 239 | 172 | 71.97% | দুঃখ(36), বিরক্তি(7), আবেগপ্রবণতা(5) / Sadness, DIsgust, Surprise, Passion |
| বিরক্তি (disgust) | 200 | 94 | 47.00% | ক্ষোভ(25), অপমান(23), দুঃখ(16), বিস্ময়(12) / Agitation, Insult, Sadness, Surprise |
| হাসি (Laughter) | 173 | 80 | 46.25% | বিস্ময়(24), অপমান(24), হতাশা(14), বিরক্তি(10) / Surprise, Insult, Disappointment, Disgust |
| কৌতূহল (curiosity) | 104 | 66 | 63.46% | বিস্ময়(17), দুঃখ(4), হতাশা(4), আনন্দ(2) / Surprise, Sadness, Disappointment, Joy |
| উদ্বেগ (concern) | 88 | 61 | 69.32% | বিস্ময়(7), কৌতূহল(3), হতাশা(3), বিরক্তি(2) / Surprise, Curiosity, Disappointment, Disgust |

**Table (a) right**

| Template | I2 | | | |
|---|---|---|---|---|
| Model | ChatGPT-4o | | | |
| Response for Woman | | Same response for Man | | |
| Word | # occurrences | # occurrences | percentage | Top responses for Man |
| আনন্দ (joy) | 1752 | 1624 | 92.69% | গর্ব(45), বিস্ময়(23), সন্তুষ্টি(7), কৃতজ্ঞতা(7) / Pride, Surprise, Satisfaction, Gratitude |
| দুঃখ (sadness) | 714 | 551 | 77.17% | ক্ষোভ(51), বিষণ্ণতা(36), অপমান(17), বিরক্তি(16) / Agitation, Depression, Insult, Disgust |
| ক্ষোভ (agitation) | 622 | 463 | 74.44% | দুঃখ(48), অপমান(42), বিরক্তি(25), হতাশা(20) / Sadness, Insult, Disgust, DIsappointment |
| বিস্ময় (surprise) | 405 | 190 | 46.91% | আনন্দ(64), হাসি(24), কৌতূহল(17), হতাশা(14) / Joy, Laughter, Curiosity, DIsappointment |
| অপমান (insult) | 399 | 202 | 50.63% | বিস্ময়(33), আনন্দ(32), হাসি(24), বিরক্তি(23) / Surprise, Joy, Laughter, DIsgust |
| হতাশা (disappointment) | 311 | 218 | 70.10% | দুঃখ(23), হাসি(14), বিরক্তি(10), ক্ষোভ(9) / Sadness, Laughter, Disgust, Agitation |
| বিষণ্ণতা (depression) | 286 | 172 | 60.14% | দুঃখ(83), বিস্ময়(11), আনন্দ(5), হতাশা(4) / Sadness, Surprise, Joy, Disappointment |
| বিরক্তি (disgust) | 248 | 94 | 37.90% | দুঃখ(36), ক্ষোভ(25), হতাশা(19), বিস্ময়(17) / Sadness, Agitation, Disappointment, Surprise |
| গর্ব (pride) | 207 | 194 | 93.72% | আনন্দ(9), সম্মান(1), অপমান(1), বিস্ময়(1) / Joy, Respect, Insult, Surprise |
| হাসি (Laughter) | 117 | 80 | 68.38% | আনন্দ(22), হতাশা(4), বিস্ময়(3), বিভ্রান্তি(3) / Joy, Disappointment, Surprise, Confusion |
| কৌতূহল (curiosity) | 98 | 66 | 67.35% | আনন্দ(11), উদ্বেগ(3), গর্ব(3), আগ্রহ(3) / Joy, Concern, Pride, Interest |
| উদ্বেগ (concern) | 79 | 61 | 77.22% | আনন্দ(4), উদাস(3), উত্তেজনা(3), বিরক্তি(1) / Joy, Melancholy, Excitement, DIsgust |

(a) Emotion Word Occurrences and Top Responses for Opposite Genders in Data Points Using GPT-4o with Prompt Template I2

**Table (b) left**

| Template | I2 | | | |
|---|---|---|---|---|
| Model | ChatGPT-3.5-Turbo | | | |
| Response for Man | | Same response for Woman | | |
| Word | # occurrences | # occurrences | percentage | Top responses for Woman |
| আনন্দ (joy) | 669 | 228 | 34.08% | উৎসাহ(91), সন্তুষ্টি(42), বিরক্তি(30), খুশি(27) / Enthusiasm, Satisfaction, Disgust, Happiness |
| বিরক্তি (disgust) | 532 | 158 | 29.70% | দুঃখ(64), আনন্দ(25), বিস্ময়(21), ক্ষোভ(14) / Sadness, Joy, Surprise, Agitation |
| উৎসাহ (excitement) | 512 | 168 | 32.81% | আনন্দ(83), গর্ব(30), উদাস(26), সন্তুষ্টি(18) / Joy, Pride, Melancholy, Satisfaction |
| দুঃখ (sadness) | 513 | 304 | 59.26% | বিরক্তি(51), আনন্দ(10), বিস্ময়(8), উদাস(6) / Disgust, Joy, Surprise, Melancholy |
| গর্ব (pride) | 422 | 220 | 52.13% | আনন্দ(39), উৎসাহ(29), দুঃখ(15), বিস্ময়(9) / Joy, Enthusiasm, Sadness, Surprise |
| হাসি (laughter) | 244 | 91 | 37.30% | আনন্দ(34), উদাস(13), বিরক্তি(11), উৎসাহ(10) / Joy, Melancholy, Disgust, Enthusiasm |
| উদাস (melancholy) | 216 | 23 | 10.65% | উৎসাহ(23), বিরক্তি(19), আনন্দ(18), দুঃখ(13) / Enthusiasm, Disgust, Joy, Sadness |
| সন্তুষ্টি (content) | 170 | 12 | 7.06% | আনন্দ(47), উৎসাহ(16), গর্ব(8), বিরক্তি(7) / Joy, Enthusiasm, Pride, Disgust |
| খুশি (happy) | 144 | 10 | 6.94% | আনন্দ(60), উৎসাহ(21), সন্তুষ্টি(11), গর্ব(8) / Joy, Enthusiasm, Satisfaction, Pride |
| বিস্ময় (surprise) | 107 | 12 | 11.21% | বিরক্তি(17), দুঃখ(7), উদাস(5), ভয়(4) / Disgust, Sadness, Melancholy, Fear |
| নিরাশা (despair) | 88 | 18 | 20.45% | বিরক্তি(19), দুঃখ(10), বিস্ময়(3), নারাজ(3) / Disgust, Sadness, Surprise, Displeased |
| ভালোবাসা (love) | 84 | 12 | 14.29% | আনন্দ(32), সন্তুষ্টি(6), বিস্ময়(5), উৎসাহ(4) / Joy, Satisfaction, Surprise, Enthusiasm |

**Table (b) right**

| Template | I2 | | | |
|---|---|---|---|---|
| Model | ChatGPT-3.5-Turbo | | | |
| Response for Woman | | Same response for Man | | |
| Word | # occurrences | # occurrences | percentage | Top responses for Man |
| আনন্দ (joy) | 828 | 228 | 27.54% | উৎসাহ(83), খুশি(60), সন্তুষ্টি(47), গর্ব(39) / Enthusiasm, Happiness, Satisfaction, Pride |
| বিরক্তি (disgust) | 694 | 158 | 22.77% | দুঃখ(51), আনন্দ(30), ক্ষোভ(19), উদাস(19) / Sadness, Joy, Agitation, Melancholy |
| দুঃখ (sadness) | 623 | 290 | 46.55% | বিরক্তি(64), উদাস(33), গর্ব(15), আনন্দ(14) / Disgust, Melancholy, Pride, Joy |
| উৎসাহ (excitement) | 523 | 168 | 32.12% | আনন্দ(91), গর্ব(29), উদাস(23), খুশি(21) / Joy, Pride, Melancholy, Happiness |
| গর্ব (pride) | 343 | 206 | 60.06% | উৎসাহ(30), আনন্দ(11), সন্তুষ্টি(8), খুশি(8) / Enthusiasm, Joy, Satisfaction, Happiness |
| উদাস (melancholy) | 215 | 23 | 10.70% | উৎসাহ(26), আনন্দ(19), হাসি(13), বিরক্তি(9) / Enthusiasm, Joy, Laughter, Disgust |
| বিস্ময় (surprise) | 200 | 12 | 6.00% | বিরক্তি(21), আনন্দ(14), উৎসাহ(13), গর্ব(9) / Disgust, Joy, Enthusiasm, Pride |
| সন্তুষ্টি (content) | 180 | 12 | 6.67% | আনন্দ(42), উৎসাহ(18), খুশি(11), উদাস(9) / Joy, Enthusiasm, Happiness, Melancholy |
| হাসি (laughter) | 157 | 91 | 57.96% | আনন্দ(14), মজা(5), উৎসাহ(4), উল্লাস(3) / Joy, Fun, Enthusiasm, Elation |
| ভয় (fear) | 93 | 11 | 11.83% | বিরক্তি(12), উদাস(9), উৎসাহ(7), দুঃখ(5) / Disgust, Melancholy, Enthusiasm, Sadness |
| খুশি (happy) | 90 | 10 | 11.11% | আনন্দ(27), উৎসাহ(8), গর্ব(7), সন্তুষ্টি(7) / Joy, Enthusiasm, Pride, Satisfaction |
| ক্ষোভ (agitation) | 72 | 4 | 5.56% | বিরক্তি(14), রাগ(7), রোষ(4), দুঃখ(4) / Disgust, Anger, Anger, Sadness |

(b) Emotion Word Occurrences and Top Responses for Opposite Genders in Data Points Using GPT-3.5-Turbo with Prompt Template I2

| Template | I2 | | | |
|---|---|---|---|---|
| Model | Llama-3 8b | | | |
| Response for Man | Same response for Woman | | | Top responses for Woman |
| Word | # occur-ences | # occur-ences | percentage | |
| বিস্ময় (surprise) | 1852 | 1572 | 84.88% | বিরক্তি(56), ব্যাখা(41), আনন্দ(33), ব্যাহত(17) — Disgust, Pain, Joy, Interrupt |
| ক্রোধ (anger) | 853 | 334 | 39.16% | বিস্ময়(234), বিরক্তি(101), ক্ষোভ(39), বিচলিত(11) — Surprise, Disgust, Agitation, Anxious |
| আনন্দ (joy) | 781 | 599 | 76.70% | আহ্বান(66), বিস্ময়(24), আহ্লাদ(22), শান্তি(10) — Invitation, Surprise, Pleasure, Peace |
| বিরক্তি (disgust) | 431 | 322 | 74.71% | বিস্ময়(78), ব্যাখা(8), ব্যাহত(3) — Surprise, Pain, Interrupt |
| অসহনীয় (unbearable) | 256 | 36 | 14.06% | বিস্ময়(105), বিরক্তি(54), আশঙ্কা(9), আশ্চর্য(6) — Surprise, Disgust, Concern, Wonder |
| হাসি (laughter) | 174 | 114 | 65.52% | বিস্ময়(21), আনন্দ(15), বিরক্তি(6), হতাশ(5) — Surprise, Joy, Disgust, Disappointed |
| আহ্বান (appeal) | 153 | 141 | 92.16% | আনন্দ(6), আহ্লাদ(4), আশা(1), শান্তি(1) — Joy, Pleasure, Hope, Peace |
| শান্তি (peace) | 124 | 63 | 50.81% | আনন্দ(30), আহ্বান(19), বিস্ময়(3), আহ্লাদ(1) — Joy, Invitation, Surprise, Pleasure |
| ক্ষোভ (agitation) | 85 | 36 | 42.35% | বিস্ময়(21), বিরক্তি(12), ব্যাহত(7), ব্যাখা(2) — Surprise, Disgust, Interrupt, Pain |
| ঘৃণা (hate) | 69 | 43 | 62.32% | বিস্ময়(8), বিরক্তি(7), ব্যাখা(4), ক্ষোভ(2) — Surprise, Disgust, Pain, Agitation |

| Template | I2 | | | |
|---|---|---|---|---|
| Model | Llama-3 8b | | | |
| Response for Woman | Same response for Man | | | Top responses for Man |
| Word | # occur-ences | # occur-ences | percentage | |
| বিস্ময় (surprise) | 2213 | 1599 | 72.25% | ক্রোধ(234), অসহনীয়(105), বিরক্তি(78), আনন্দ(24) — Anger, Intolerable, Disgust, Joy |
| আনন্দ (joy) | 819 | 590 | 72.04% | বিস্ময়(33), শান্তি(30), সম্মান(17), হাসি(15) — Surprise, Peace, Respect, Laughter |
| বিরক্তি (disgust) | 578 | 312 | 53.98% | ক্রোধ(101), বিস্ময়(56), অসহনীয়(54), ক্ষোভ(11) — Anger, Surprise, Intolerable, Agitation |
| ক্রোধ (anger) | 367 | 334 | 91.01% | ক্ষোভ(1), বিরক্তি(1), বিস্ময়(1) — Agitation, Disgust, Surprise |
| আহ্বান (appeal) | 274 | 140 | 51.09% | আনন্দ(66), শান্তি(19), কৃতজ্ঞতা(8), আহ্লাদ(7) — Joy, Peace, Gratitude, Pleasure |
| ব্যাখা (hurt) | 122 | 79 | 64.75% | বিস্ময়(18), ক্রোধ(7), বিরক্তি(4), বিস্মৃতি(2) — Surprise, Anger, Disgust, Oblivion |
| হাসি (laughter) | 110 | 107 | 97.27% | আনন্দ(1), বিস্ময়(1), বিরক্তি(1) — Joy, Surprise, Disgust |
| শান্তি (peace) | 107 | 63 | 58.88% | আনন্দ(10), ভালোবাসা(5), সুখ(4), ভালো(4) — Joy, Love, Bliss, Good |
| ক্ষোভ (agitation) | 80 | 36 | 45.00% | ক্রোধ(39), ঘৃণা(2), বিস্ময়(1), অসহনীয়(1) — Rage, Hatred, Surprise, Intolerable |
| আহ্লাদ (delight) | 61 | 30 | 49.18% | আনন্দ(22), আহ্বান(4), সুখ(3), শান্তি(1) — Joy, Invitation, Bliss, Peace |

(c) Emotion Word Occurrences and Top Responses for Opposite Genders in Data Points Using Llama-3 with Prompt Template I2

Table 6: Detailed Analysis of Emotion Word Occurrences for Male and Female Responses Using Prompt Template I2 Across Different LLMs. Sub-table 6b presents results for ChatGPT-3.5-Turbo, showing the number of occurrences of each emotion word in male and female responses, the corresponding occurrences in opposite gender responses, and the top responses for the opposite gender provided the same data points. Sub-table 6b provides analogous data for Llama-3-8b.