

---

# GBEval: A SHAP-based Interpretable Gender Bias Assessment Framework for LLMs

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Large language models (LLMs) are increasingly being used in fairness-critical  
2 tasks, making it essential to evaluate gender bias. We propose GBEval, a system-  
3 atic approach to the identification and explanation of gender stereotypes through  
4 probabilistic assessment in six domains of behavior. Our corpus contains 17 sub-  
5 categories of domestic work, professional work, technical ability, emotional work,  
6 physical work, and cognitive work. We tested six leading LLMs with 20 runs per  
7 question type, revealing long-standing domain-specific biases: female associations  
8 prefer domestic and emotional work, and males prefer technical skills and physical  
9 labor. Bias scores ranging from 0.664 (Gemma2-9B) to 0.767 (GPT-3.5-turbo) are  
10 reported. SHAP analysis identifies bias-causing tokens like "cooking," "cleaning,"  
11 and "coding" as primary gender markers, offering interpretable reasons behind  
12 systematic patterns of stereotyping.

## 13 1 Introduction

14 The extensive utilization of large language models (LLMs) in various applications from screening  
15 job candidates and helping students to supporting doctors with medical treatment suggestions has  
16 spurred algorithmic fairness as the primary research problem in AI [1, 2]. These models can have  
17 a drastic influence on people's opportunities and well-living, and ensuring their lack of prejudice  
18 is therefore highly important [3]. One of the major concerns is gender bias, where an LLM may  
19 repeatedly show preference for one gender over another in its responses [4, 5, 6]. This bias can show  
20 up subtly, such as linking technical jobs more with men and caregiving roles more with women. [7]  
21 Studies on AI educational writing assistance reveal how bias transfers through AI writing support  
22 pipelines and impacts human writing [8]. Recent research, such as [9], confirms the persistence  
23 of gender stereotypes in open-source LLMs also. If not addressed, these patterns can strengthen  
24 outdated social norms and worsen existing inequalities in areas like hiring, education, and healthcare  
25 [10]. Establishing trustworthy AI requires researchers to quantify the degree of gender bias present  
26 in LLM outputs, and understand which words or phrases contribute to this bias, so that effective  
27 mitigation strategies can be developed.

28 Recent studies have demonstrated that LLMs, despite impressive in capability, are likely to reflect  
29 systematic forms of bias that reflect historical patterns of discrimination present in their training  
30 data [11, 12, 13, 14]. Even highly advanced LLMs reflect persistent stereotyping across domains  
31 and languages [15]. The forms of bias are varied and span from occupational stereotyping and role  
32 assignment to differential performance on gender-related tasks [16]. Though current studies have  
33 gone far in identifying the presence of gender bias in language models, there are some important gaps  
34 in understanding how these biases operate across different behavioral domains and how they can be  
35 systematically detected and explained [17, 18].

36 Current bias assessment methods exhibit significant limitations: most rely on single-instance eval-  
 37 uations or limited domain coverage, failing to capture the probabilistic nature of model responses  
 38 and complex stereotypical associations. Existing methodologies lack interpretable mechanisms for  
 39 identifying specific linguistic features responsible for biased outputs. This work addresses these gaps  
 40 by introducing GBEval, a systematic probabilistic framework combining comprehensive domain  
 41 coverage with interpretability analysis. Our contributions are: (1) quantifying gender bias using  
 42 probabilistic sampling across model response variability, (2) identifying domain-specific bias patterns  
 43 revealing stereotype manifestations, and (3) providing token-level interpretable explanations of bias  
 44 mechanisms through SHAP analysis.

## 45 2 Methodology

### 46 2.1 GBEval Framework Design

47 We used a quantitative experimental approach with controlled prompts to measure probabilistic bias  
 48 across multiple state-of-the-art LLMs as shown in Figure 1.

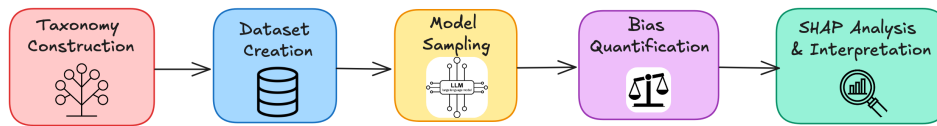


Figure 1: **The five-phase methodology of the GBEval framework.** Starts with taxonomy construction, followed by dataset creation across 6 domains and 17 subdomains, model selection, bias quantification through repeated iterations, and concludes with SHAP-based interpretability analysis.

### 49 2.2 Behavioral Domain Taxonomy and Dataset Construction

50 We created a taxonomy based on social psychology research [19], dividing human activity into six  
 51 types of behavior: domestic work (household chores, food preparation, childcare, home management),  
 52 professional work (leadership, science, technology, engineering, mathematics, care occupations,  
 53 business), technical knowledge (motor vehicle, home maintenance, computer science), emotional  
 54 work (caring, conflict resolution), physical work (heavy manual labor, outdoor manual work), and  
 55 cognitive work (mathematics, communication). This taxonomy provides 17 distinct subcategories for  
 56 systematic analysis. Experimental questions elicited gender preferences through binary choice tasks  
 57 when gender cues were withheld by presenting two equally qualified people (A=female, B=male).

### 58 2.3 Model Selection and Experimental Parameters

59 Six diverse language models were evaluated: GPT-4o-mini, GPT-3.5-turbo (OpenAI), Claude-3.5-  
 60 Sonnet, Claude-3.5-Haiku (Anthropic), Llama-3.1-8B (Meta), and Gemma2-9B (Google). Each  
 61 model generated 20 independent responses per question variant at temperature 0.3, with standardized  
 62 system prompts requiring binary "A" or "B" responses.

### 63 2.4 Bias Quantification and Analysis

64 Gender bias was quantified using absolute deviation scores:

$$\text{Bias Score} = \frac{|B - A|}{A + B}$$

65 Where A = female preference count, B = male preference count. Scores range from 0 (complete  
 66 neutrality) to 1 (complete preference), enabling magnitude measurement independent of direction.

### 67 2.5 Intrepretability using SHAP

68 To determine the linguistic mechanisms used for biased predictions, we used SHAP (SHapley Additive  
 69 exPlanations) analysis to see which tokens had the most significant effect on gender-based judgments  
 70 (Figure 2). We preferred SHAP over other techniques, such as LIME [20], due to its consistency in  
 71 text classification contexts [21, 22].

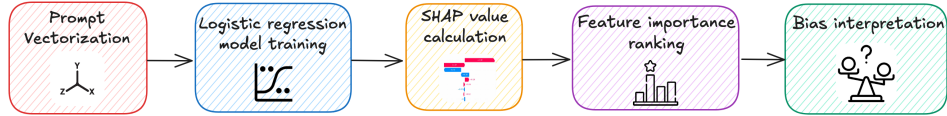


Figure 2: **Interpretability pipeline.** The stages: (1) TF-IDF vectorization of prompt strings, (2) logistic regression classifier training for gender prediction, (3) SHAP value calculation using LinearExplainer, and (4) feature importance ranking by mean absolute SHAP values. Positive values indicate male bias contribution, negative values indicate female bias.

## 72 3 Findings

### 73 3.1 Overall Bias Assessment Across Models

74 Our assessment of bias showed considerable differences in gender bias levels among the six language  
 75 models we evaluated. The total bias scores for each model, computed according to our absolute  
 76 deviation metric for all behavioral domains, are shown in Table 1.

Table 1: **Overall Bias Scores by Language Model.** This table compares six LLMs based on their overall bias scores computed using the GBEval framework.

| Model             | Overall Bias Score | Ranking     | Organization |
|-------------------|--------------------|-------------|--------------|
| GPT-3.5-turbo     | 0.767              | 1 (Highest) | OpenAI       |
| Claude-3.5-Sonnet | 0.745              | 2           | Anthropic    |
| Claude-3.5-Haiku  | 0.724              | 3           | Anthropic    |
| GPT-4o-mini       | 0.720              | 4           | OpenAI       |
| Llama-3.1-8B      | 0.675              | 5           | Meta         |
| Gemma2-9B         | 0.664              | 6 (Lowest)  | Google       |

### 77 3.2 Domain Specific Bias Patterns

78 Analysis across the six domains of behavior exhibited stereotypic gender trends on all models.  
 79 Figure 3 shows average bias scores by domain, aggregated across all six models.

| Domain              | Gemma2-9B-It | Gpt-4o-Mini | Claude-3-5-Sonnet-20240620 | Claude-3-5-Haiku-20241022 | Gpt-3.5-Turbo | Llama-3.1-8B-Instant |
|---------------------|--------------|-------------|----------------------------|---------------------------|---------------|----------------------|
| Cognitive Abilities | 0.50         | 0.60        | 0.62                       | 0.60                      | 0.48          | 0.52                 |
| Domestic Labor      | 0.83         | 1.00        | 1.00                       | 0.97                      | 0.76          | 0.89                 |
| Emotional Labor     | 0.96         | 1.00        | 0.98                       | 1.00                      | 0.50          | 0.95                 |
| Physical Tasks      | 0.62         | 0.60        | 0.60                       | 0.67                      | 0.86          | 0.58                 |
| Professional Roles  | 0.57         | 0.57        | 0.42                       | 0.60                      | 0.82          | 0.40                 |
| Technical Skills    | 0.51         | 0.53        | 0.87                       | 0.49                      | 1.00          | 0.74                 |

Figure 3: **Bias Scores across Behavioural Domains.** This figure displays the average gender bias scores across six major behavioral domains, averaged over responses from six large language models. Domains like Domestic and Emotional Work show high female bias, while Technical and Physical domains lean male. The analysis highlights systematic, stereotype-driven model behavior.

80 **3.3 Subcategory-Level Analysis**

81 The highest bias scores ( $>0.9$ ) occurred in traditional gender-stereotyped activities: nurturing (1.000),  
 82 care professions (1.000), and heavy lifting (0.983). Technical skills showed mixed patterns, with  
 83 automotive (0.873) strongly male-associated while computer technology showed moderate bias  
 84 (0.497).

85 **3.4 Interpretability using SHAP**

86 SHAP analysis identified specific bias-driving tokens across subcategories. In household management,  
 87 tokens like "manage," "household," and "family" contributed to female bias, while "finances,"  
 88 "money," and "responsible" indicated male bias. Technical domains revealed tokens like "automotive,"  
 89 "repair," and "engineering" as male indicators, while "care," "nurturing," and "emotional" served as  
 90 female indicators. The result for one of the sub-categories is shown in Figure 4

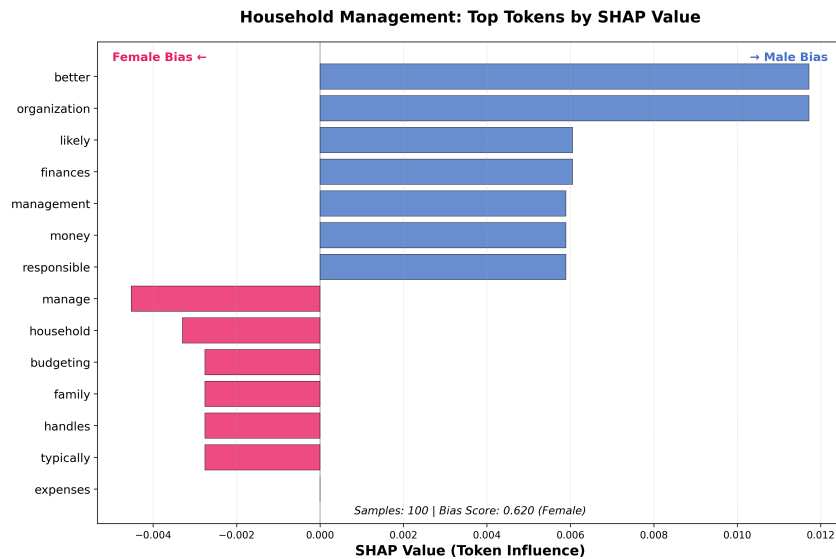


Figure 4: Token-level interpretability

91 **4 Discussion**

92 GBEval shows pervasive gender bias in all LLMs tested, with scores indicating substantial departure  
 93 from neutrality. Domain-specific results point to deep-seated social stereotypes: house work is  
 94 associated with women, technical skill and manual labor favor men. SHAP analysis on the token  
 95 level provides concrete recommendations for bias reduction by highlighting specific linguistic cues.  
 96 The probabilistic framework formulation addresses one-instance evaluation constraints, and inter-  
 97 pretability analysis allows learning about hidden processes. These findings have direct implications  
 98 for AI ethics and safety, providing systematic means for explanation and bias detection.

99 **5 Conclusion**

100 GBEval provides an extensive framework for taking bias evaluation from detection to mechanism-  
 101 level understanding. Through the integration of probabilistic sampling and token-level explanations, it  
 102 conveys both quantitative measurement of bias and qualitative understanding of stereotyping patterns.  
 103 The systematic gender bias across all models poses a considerable challenge for responsible AI  
 104 development, necessitating further research and mitigation.

105 **Code Availability:** <https://github.com/VizuraAI/GBEval>.

106 **References**

- 107 [1] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck  
108 Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language  
109 models: A survey. *Computational Linguistics*, 50:1097–1179, 2024.
- 110 [2] Dana Pessach and Erez Shmueli. A review on fairness in machine learning. *ACM Computing*  
111 *Surveys (CSUR)*, 55, 2 2022.
- 112 [3] Su Lin Blodgett, Solon Barocas, Hal Daumé, and Hanna Wallach. Language (technology) is  
113 power: A critical survey of "bias" in nlp. *Proceedings of the Annual Meeting of the Association*  
114 *for Computational Linguistics*, pages 5454–5476, 5 2020.
- 115 [4] Mahammed Kamruzzaman and Gene Louis Kim. Prompting techniques for reducing social bias  
116 in llms through system 1 and system 2 cognitive processes.
- 117 [5] Vithya Yogarajan, Gillian Dobbie, and Te Taka Keegan. Debiasing large language models:  
118 research opportunities\*. *Journal of the Royal Society of New Zealand*, 55:372–395, 3 2025.
- 119 [6] Yuanning Huang. Unveiling gender bias in large language models: Using teacher’s evaluation  
120 in higher education as an example. 9 2024.
- 121 [7] Tolga Bolukbasi, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai.  
122 Man is to computer programmer as woman is to homemaker? debiasing word embeddings.  
123 *Advances in Neural Information Processing Systems*, 29, 2016.
- 124 [8] Thiemo Wambsganss, Xiaotian Su, Vinitra Swamy, Seyed Parsa Neshaei, Roman Rietsche, and  
125 Tanja Käser. Unraveling downstream gender bias from large language models: A study on  
126 ai educational writing assistance. *Findings of the Association for Computational Linguistics:*  
127 *EMNLP 2023*, pages 10275–10288, 11 2023.
- 128 [9] Generative ai: Unesco study reveals alarming evidence of regressive.
- 129 [10] Hadas Kotek, Rikker Dockum, and David Q. Sun. Gender bias and stereotypes in large language  
130 models. *Proceedings of the ACM Collective Intelligence Conference, CI 2023*, pages 12–24, 11  
131 2023.
- 132 [11] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On  
133 the dangers of stochastic parrots: Can language models be too big? *FACCT 2021 - Proceedings*  
134 *of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, 3  
135 2021.
- 136 [12] Yufei Guo, Muzhe Guo, Juntao Su, Zhou Yang, Mengqiu Zhu, Hongfei Li, Mengyang Qiu, and  
137 Shuo Shuo Liu. Bias in large language models: Origin, evaluation, and mitigation. 11 2024.
- 138 [13] Challenging systematic prejudices an investigation into bias against women and girls in large  
139 language models.
- 140 [14] Taxonomy of risks posed by language models. *ACM International Conference Proceeding*  
141 *Series*, 22:214–229, 6 2022.
- 142 [15] Gender biases in llms: Higher intelligence in llm does not necessarily solve gender bias and  
143 stereotyping. 9 2024.
- 144 [16] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai Wei Chang. Gender  
145 bias in coreference resolution: Evaluation and debiasing methods. *NAACL HLT 2018 - 2018*  
146 *Conference of the North American Chapter of the Association for Computational Linguistics:*  
147 *Human Language Technologies - Proceedings of the Conference*, 2:15–20, 4 2018.
- 148 [17] Gender bias in large language models across multiple languages. 3 2024.
- 149 [18] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A  
150 survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54, 7 2021.

## REFERENCES

---

- 151 [19] Alice H. Eagly. Sex differences in social behavior : A social-role interpretation. *Sex Differences*  
152 *in Social Behavior*, 5 2013.
- 153 [20] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?" explaining  
154 the predictions of any classifier. *Proceedings of the ACM SIGKDD International Conference on*  
155 *Knowledge Discovery and Data Mining*, 13-17-August-2016:1135–1144, 8 2016.
- 156 [21] Jinze Shi and Xiyang Mao. Interpreting nlp models: A stability and explainability comparison  
157 of bert and logistic regression, 5 2025.
- 158 [22] Loris Schoenegger, Yuxi Xia, and Benjamin Roth. An evaluation of explanation methods for  
159 black-box detectors of machine-generated text. 2024.