# A Comprehensive Analysis of Large Language Model Outputs: Similarity, Diversity, and Bias

BRANDON SMITH, Deakin University, Australia
MOHAMED REDA BOUADJENEK, Deakin University, Australia
TAHSIN ALAMGIR KHEYA, Deakin University, Australia
PHILLIP DAWSON, Deakin University, Australia
SUNIL ARYAL, Deakin University, Australia

Large Language Models (LLMs) represent a significant leap forward in the quest for artificial general intelligence, enhancing our capacity to engage with and leverage technology. However, while LLMs have demonstrated their effectiveness in various Natural Language Processing tasks such as language translation, text generation, code generation, and content summarization, among others, there remain many open questions about their similarities, their variance, and their ethical implications. For example, when presented with a text generation prompt, what similarities exist among texts produced by the same language models? In addition, what is the inter-LLM writing similarity - for instance, how comparable are texts generated by distinct LLMs when presented with the same prompt? Furthermore, how does the variation in text generation manifest across multiple LLMs, and which LLMs adhere most closely to ethical standards? To address these questions, we used 5K different prompts covering a diverse range of requests, including generating, explaining, and rewriting text. This effort resulted in the generation of approximately 3M texts from 12 different LLMs, featuring a mix of proprietary and open-source models from industry leaders such as OpenAI, Google, Microsoft, Meta, and Mistral. The results of this study reveals a number of important insights, among them that: (1) texts produced by the same LLMs show higher similarity compared to human-written texts; (2) some LLMs, like WizardLM-2-8x22b, produce highly similar texts, while others like GPT-4 generate more varied outputs; (3) writing styles among LLMs vary significantly, with models like Llama 3 and Mistral showing high similarity, while GPT-4 stands out for its distinctiveness; (4) the sharp contrast in language and lack of vocabulary overlap highlight the distinct linguistic characteristics of LLM-produced text; (5) finally, it appears that certain LLMs are more balanced in terms of gender representation and are less prone to perpetuating bias.

Keywords: Large Language Models; NLP; Text Generation.

## 1 INTRODUCTION

In the past year, there has been a significant stride in the pursuit of artificial general intelligence, mainly highlighted by the implementation of Large Language Models (LLMs) [1–8] to crafting chatbots such at ChatGPT, BARD, Bing Chat, and Grok. These models showcase unparalleled

Authors' addresses: Brandon Smith, Deakin University, Burwood, VIC, Australia, brandon.smith@deakin.edu.au; Mohamed Reda Bouadjenek, Deakin University, Geelong, VIC, Australia, reda.bouadjenek@deakin.edu.au; Tahsin Alamgir Kheya, Deakin University, Geelong, VIC, Australia, t.kheya@deakin.edu.au; Phillip Dawson, Deakin University, Melbourne, VIC, Australia, p.dawson@deakin.edu.au; Sunil Aryal, Deakin University, Geelong, VIC, Australia, sunil.aryal@unsw.edu.au.

precision and accuracy, particularly in excelling at diverse Natural Language Processing (NLP) tasks, including but not limited to, language translation, text generation, code synthesis, content summarization, conversation, and information search.

While these LLMs and chatbots undoubtedly help users complete their tasks, many open questions remain concerning ethical considerations, utilization challenges, output quality concerns, and the similarities and variances among them. Hence, in this paper, we seek to analyze their performance across diverse prompts, compare their outputs to human-written texts, and investigate the distinctive characteristics and variations in the texts they generate. Our fundamental research questions are focused on understanding the intrinsic characteristics of texts generated by LLMs, as outlined below:

RQ1 When presented with a prompt, what similarities exist among texts produced by the same LLM (inner-LLM similarity), and how similar are the texts generated by different LLMs for the same prompt (inter-LLM similarity)?

RQ2 How does the variation in text generation manifest across multiple LLMs?

RQ3 Can we accurately identify whether a given text was authored by a human or a specific language model?

RQ4 Are there specific words that can act as distinctive markers for each LLM?

RQ5 Are there LLMs that adhere most closely to ethical standards by reducing the propagation of biased stereotypes?

Addressing these questions is crucial as it will provide us with deeper insights into the operational challenges and performance nuances of LLMs and Chatbots across diverse domains, including education, scientific writing, and business communication. Hence, we propose in this paper a comprehensive analysis of LLM outputs. In particular, we employed approximately 5,000 distinct and diverse prompts covering diverse topics ranging from technological impact to academic performance. These prompts encompass a wide range of requests, including text generation, explanation, and rewriting. Using these prompts, we generated texts with 12 different LLMs, including proprietary and open-source models such as: Gemini-pro-1.5 [3], Gemma-7B [9], GPT-3.5, GPT-4 [8], Mistral-7B, Mixtral-8x7B, Mixtral-8x22B [4], WizardLM-2-7B, WizardLM-2-8x22B [10], Llama 3-70B, Llama 3-8B [6], and DBRX. Additionally, we include human-generated text produced using 15 prompts as instructions for text comparison. The resulting dataset comprises 3 Million texts and is utilized to conduct a thorough analysis from which we draw meaningful conclusions including:

- **Low Similarity within LLMs:** The texts produced by the same LLMs generally show higher consistency and similarity in their outputs compared to human-written texts. Some LLMs, such as WizardLM-2-8x22b, produce highly similar texts, while others like GPT-4 generate more varied outputs. Additionally, proprietary models tend to be more consistent than open-source models in terms of output similarity.

- **Inter-LLM Writing Similarity:** The writing styles vary significantly, with some models like Llama 3 and Mistral showing high similarity, while GPT-4 stands out for its distinctiveness. Although at the same time word-level similarity measures indicated that GPT-4 was the most distinct and not similar to GPT-3.5, BERT revealed a different aspect, showing that GPT-4 and GPT-3.5 are similar in deeper, contextual ways.

- **Variance in Text Generation:** Some LLMs demonstrated significant variance in text generation, while others exhibited a more consistent output. This insight into the diversity of LLM behavior is crucial for understanding the nuances of these models.

- **Classification:** Our classification efforts show success in being able to differentiate between human-written text and text generated by various LLMs. Misclassifications and confusion mainly occurred between similar models like GPT-3.5 and GPT-4, highlighting the challenge of distinguishing between texts from closely related architectures.

- **Language markers:** The sharp contrast in language and the absence of vocabulary overlap emphasize the distinct linguistic characteristics between human-generated text and that produced by the language models.

- **Ethics consideration:** it seems that certain LLMs like Gemma-7B and Gemini-pro demonstrate a more balanced approach to gender representation and are less likely to perpetuate bias. In contrast, models such as GPT-3.5 and GPT-4, while powerful in terms of performance, demonstrate a stronger tendency toward gender and racial bias.

In summary, this paper sheds light on the intricacies of LLM behavior and strives to advance the discourse on the responsible development and utilization of LLMs.

## 2 RELATED WORK

LLMs have emerged as a prominent subject of research among the academic community. Below, we review the evolution of Language Models, including research related to the analysis of LLMs, efforts focused on their detection and classification, and studies addressing bias and discrimination in LLMs.

**Evolution of Language Models:** The field of NLP has grown rapidly since the introduction of word embeddings in 2013 [11], with Word2Vec providing a foundation for capturing semantic word relationships in the vector space. This foundational approach gained further significance when employed in conjunction with sequence models like RNNs and LSTMs [12], establishing itself as a critical element in addressing complex NLP tasks. Additionally, a significant progression was the introduction of the Transformer architecture by Vaswani et al. in 2017 [13], which moved away from recurrent layers in favor of self-attention mechanisms, leading to parallel processing and a reduction in training times. Moreover, Google's BERT Language Model, introduced in 2018 [14], employed *encoder* transformer blocks to learn bi-directional contextual representations and set new performance benchmarks across a variety of NLP tasks. In contrast, OpenAI introduced the GPT series [15, 16], in particular GPT-3 with its 175 billion parameters, based on *decoder* transformer blocks achieved state-of-the-art language generation capabilities [16]. This evolution has launched the quest toward training LLMs that consistently demonstrate a form of near artificial general intelligence across various NLP tasks.

**Analyzing LLMs:** While LLMs have proven to be highly useful, their adoption has also introduced significant challenges, particularly in terms of explainability to allow effective debugging and performance enhancement [17–20]. On the other hand, there are other challenges associated with their usage that need to be analyzed, including: (1) **Trustworthiness and Toxicity:** Recent studies [21, 22] have highlighted the issues of trustworthiness and safety in LLM outputs. These concerns are worsen by the potential for LLMs to generate toxic content [23, 24]. Hence, recent work [25] has explored mitigating such risks by including pre-training instructions. (2) **Memorisation:** LLMs tend to memorise and regurgitate training data [26, 27], where privacy can be for instance violated during inference [28]. This represents significant challenges related to

Table 1. Summary of LLMs and their Key Features, with "✓" indicating the presence of the feature, "✗" indicating its absence, and "?" indicating that the information is not provided.

|   | Model | Company | #Parameters | Open-Source | Training Data Size | MoE[†] Architecture |
|---|-------|---------|-------------|-------------|--------------------|---------------------|
| 1 | Gemini-pro-1.5 | Google | ? | ✗ | 10M tokens | ✓ |
| 2 | Gemma-7B | Google | 7B | ✓ | ? | ✗ |
| 3 | GPT-3.5 | OpenAI | 175B | ✗ | 570 GB text and code | ? |
| 4 | GPT-4 | OpenAI | ? | ✗ | 1.2T tokens | ? |
| 5 | Mistral-7B | Mistral | 7B | ✓ | ? | ✗ |
| 6 | Mixtral-8x7B | Mistral | 47B | ✓ | ? | ✓ |
| 7 | Mixtral-8x22B | Mistral | 141B | ✓ | ? | ✓ |
| 8 | WizardLM-2-7B | Microsoft | 7B | ✓ | ? | ✗ |
| 9 | WizardLM-2-8x22B | Microsoft | 176B | ✓ | ? | ✓ |
| 10 | Llama 3 (8B) | Meta | 8B | ✓ | 15T tokens | ✗ |
| 11 | Llama 3 (70B) | Meta | 70B | ✓ | 15T tokens | ✗ |
| 12 | DBRX | Databricks | 132B | ✗ | 12T tokens | ✓ |

[†]A Mixture of Experts is an architectural pattern for neural networks that splits the computation of a layer or operation (such as linear layers, MLPs, or attention projection) into multiple "expert" subnetworks.

privacy, utility, and fairness. (3) **Hallucination and Reliability:** The issue of hallucination in LLM responses [29, 30] further complicates their reliability. A few works have looked at mitigating these hallucinations [31] by exploring various strategies such as fine-tuning [32], memory augmentation [33], or prompt strategies [34].

**Detecting LLM generated text:** Detecting text generated by LLMs is a critical task, and several studies have explored this issue. These investigations span from assessing the effectiveness of current detection systems against adversarial attacks [35] to the complex challenges of social media bot detection [36] and the identification of texts that blend human-written and LLM-generated content [37]. However, the reliability of these systems varies significantly when LLMs act as co-authors, influenced by factors such as the user's prompt or the linguistic background of the user. In response to these challenges, several detection methods have been proposed including one-class models that treat the target text as outliers [38], techniques leveraging sentence-level features for classification [39], ensemble approaches [40], and methods utilising stylometric features and smaller models [41, 42].

**Bias in LLMs:** Bias in LLMs is a significant issue, which can stem from training data that can have stereotypes and prejudices embedded in them [43–45]. Recent research has heavily focused on investigating and addressing such biases in LLMs, which might be present as derogatory languages towards certain minority groups [46, 47], inconsistencies in system performance across different linguistic variations [48–51], misrepresentations of groups in society [52, 53], manifestation of historic stereotypes [54–58] and general hate-filled language. For assessing bias in such systems, researchers have employed several methods including: (1) **Embedding-based**: When testing for bias in word embeddings, word analogy tests [59–61] and word association tests [62, 63] are widely employed. In word analogy tests the semantic relationship between a pair of words is tested (e.g., Man : Computer programmer ⟺ Woman : Homemaker). For word association test, bias is measured by evaluating how different classes of words (like Female names vs Male names) are associated with other words (e.g., pleasant vs unpleasant adjectives); (2) **Template-Based**: This approach involves crafting specific templates designed to expose potential biases within language models. By substituting different words in placeholder positions (e.g., "The **[Name]** is a **[Occupation]**"), the model's predictions can be analyzed for bias. For example, if "John" is substituted for **[Name]**, the method examines what the model predicts for **[Occupation]** [54, 55, 58, 64–67]; and (3) **Generated-text based**: This approach involves prompting the model and letting it generate text and then analyzing the content for biased representations. Free-text output from LLMs can be analyzed for bias using several metrics including the ones proposed by [68–70].
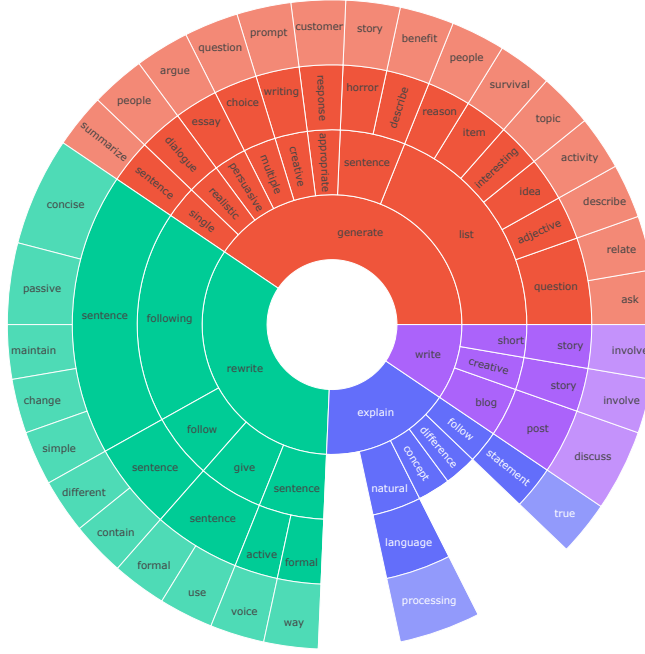
Fig. 1. Distribution of first words in prompts used to create the dataset used in this analysis.

## 3 DATASET AND LLMs OVERVIEW

In this section, we first provide a brief description of the LLMs examined, followed by an overview of the dataset collected and utilized for our analysis.

### 3.1 Large Language Models

LLMs are advanced AI systems designed to understand, generate, and manipulate human language by leveraging vast amounts of data and sophisticated neural network architectures, often based on transformers. In this work, we focus on analyzing the outputs from the LLMs summarized in Table 1.

### 3.2 Data Description

We first combined 15 prompts from the PERSUADE 2.0 corpus, introduced by Crossley et al. [71][1], with a selection of prompts/instructions from the Stanford Alpaca project dataset [72][2]. The PER-SUADE 2.0 dataset was chosen for its human-generated text corresponding to each prompt/instruction, while the Stanford Alpaca dataset was included for its rich number of diverse prompts. This combination resulted in a comprehensive set of 5,015 unique prompts, predominantly featuring keywords depicted in Figure 1. These instructions spanned a variety of topics, from learning artificial intelligence to generating persuasive essays. Finally, each prompt was used to generate 50 texts from each LLM, resulting in a total of approximately 3M texts, including those written by humans.

Detailed statistics of the resulted dataset are provided in Table 2. For example, for the given prompts, the average word count is 202, the average sentence count is 13. Also, out of the 250,750

---

[1]https://github.com/scrosseye/persuade_corpus_2.0
[2]https://github.com/tatsu-lab/stanford_alpaca

Table 2. Dataset details and statistics.

Prompts Summary

| #Prompts | #Unique Words | Average Words | Average Sentences | #topics |
|---|---|---|---|---|
| 5,015 | 96,606 | 202 | 13 | 26 |

Text Measures

| Model | #texts | Unique Word Ratio | Entropy Ratio | Monosyllable Ratio | Polysyllable Ratio | Lexical Diversity |
|---|---|---|---|---|---|---|
| Databricks | 250,750 | 75.0 | 10.8 | 64.3 | 14.3 | 3.5 |
| GPT-3.5 Turbo | 250,750 | 85.0 | 17.8 | 67.2 | 12.7 | 5.9 |
| GPT-4 | 250,750 | 86.4 | 18.1 | 66.4 | 13.3 | 6.1 |
| Gemini-pro-1.5 | 250,750 | 69.5 | 11.9 | 62.0 | 16.0 | 5.3 |
| Gemma-7B | 250,750 | 71.8 | 11.2 | 60.6 | 16.6 | 3.0 |
| Meta-Llama-3-70B | 250,750 | 60.6 | 6.6 | 67.0 | 12.9 | 1.4 |
| Meta-Llama-3-8B | 250,750 | 61.2 | 7.0 | 67.7 | 12.5 | 1.5 |
| Mistral-7B | 250,750 | 64.9 | 8.7 | 66.7 | 13.1 | 2.0 |
| Mixtral-8x22B | 250,750 | 66.1 | 8.9 | 65.9 | 13.6 | 2.3 |
| Mixtral-8x7B | 250,750 | 62.6 | 7.1 | 67.4 | 12.7 | 1.3 |
| WizardLM-2-7B | 250,750 | 64.3 | 7.8 | 65.4 | 14.1 | 2.1 |
| WizardLM-2-8x22B | 250,750 | 64.3 | 7.9 | 65.3 | 13.9 | 2.2 |
| Human | 25,000 | 43.0 | 1.9 | 78.9 | 5.2 | 0.1 |

Table 3. Model generation hyperparameters.

| Model | Max Tokens | Temperature | Top-p | Frequency Penalty | Repetition Penalty |
|---|---|---|---|---|---|
| Databricks | 512 | 0.7 | 0.9 | 0.0 | 1.0 |
| GPT-3.5 Turbo | - | 0.7 | 1.0 | 0.0 | 1.0 |
| GPT-4 | 250 | 0.7 | 1.0 | 0.0 | 1.0 |
| Gemini Pro 1.5 | - | 1.0 | 0.95 | - | 1.0 |
| Gemma-7B | 512 | 1.0 | 1.0 | 0.0 | 1.0 |
| Meta-Llama-3-70B | 512 | 0.7 | 0.9 | 0.0 | 1.0 |
| Meta-Llama-3-8B | 512 | 0.7 | 0.9 | 0.0 | 1.0 |
| Mistral-7B | 512 | 0.7 | 0.9 | 0.0 | 1.0 |
| Mixtral-8x22B | 512 | 0.7 | 0.9 | 0.0 | 1.0 |
| Mixtral-8x7B | 512 | 0.7 | 0.9 | 0.0 | 1.0 |
| WizardLM-2-7B | 512 | 0.7 | 0.0 | 0.0 | 1.0 |
| WizardLM-2-8x22B | 512 | 0.7 | 0.9 | 0.0 | 1.0 |

texts generated using GPT-3.5 for all prompts, there are 85% unique words, entropy ratio of 14, and a lexical diversity of 3.5.

Finally, Table 3 presents the generation hyperparameters used across all LLMs evaluated in this study. These parameters—including maximum tokens, temperature, top-p, frequency penalty, and repetition penalty—play a critical role in shaping the style, diversity, and consistency of the outputs produced by each model.

## 4 EMPIRICAL EVALUATION

In this section, we conduct a series of evaluations and discussions aimed at thoroughly examining the analyzed LLMs, with the intent of answering the research questions outlined above.

### 4.1 RQ1: Comparative Analysis of LLM Texts

To address RQ1, we conduct a text similarity analysis and apply various readability statistics to assess the data, as detailed below.

*4.1.1 Text Similarity Analysis .* For analyzing the similarity between the outputs of the examined LLMs, the procedure is as follows: For each prompt, we compare each generated text to others associated with the same prompt. Specifically, we compute pairwise similarities using both *cosine similarity* and *Word-Level Levenshtein Edit Distance*, which measures the number of single-word
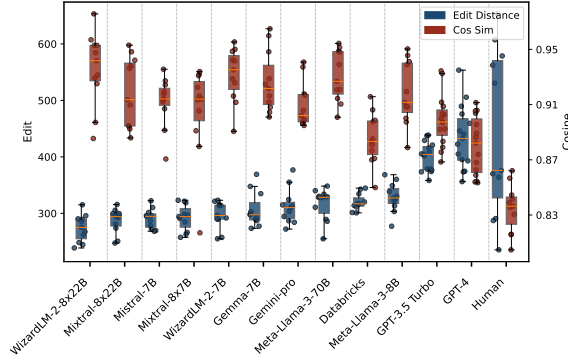
Fig. 2. Comparison of inner-text similarity.

edits required to transform one text into another. This comparison is conducted across various contexts, including within the same LLM, among different LLMs, or with human-generated texts. This approach allows for a comprehensive exploration of text similarities across different scenarios and models.

**Inner-text similarity:** The results of this analysis are presented in Figure 2, from which we make the following notes: (i) Humans exhibit lower word-level similarity to one another compared to LLMs, reflecting the unique and individualistic nature of human writing styles. Human written text also show higher variance in inner-similarity, while some LLMs tend to show more consistent similarity levels. (ii) Among LLMs, WizardLM-2-8x22b shows the highest similarity in generated text, followed by Llama-3-70b and WizardLM-2-7b. (iii) Finally, GPT-4 has the lowest similarity to its own outputs, aligning with the stylometric features seen in Table 2, indicating GPT-4's high lexical diversity and unique word ratio compared to other LLMs. In summary, these findings collectively indicate a degree of unpredictability in the outputs of LLMs.

**Inter-text similarity:** The findings from this analysis are depicted in Figure 3, leading to the following observations: (i) Human-written text shows the least similarity to all LLMs. Among the LLMs, Mistral appears to be the most similar to human-written text, while OpenAI models are the least similar. (ii) There is a notable similarity between the Llama 3 models and Mistral, as well as with WizardLM-2 models. Llama-3-8B and 70B exhibit the highest similarity among all LLMs, whereas GPT-4 and GPT-3.5 display the least similarity, highlighting the diversity between these two models. (iii) Lastly, while GPT-3.5 shares some similarity with models like Llama 3 and Mixtral-8xx22B, GPT-4, similar to human-written text, shows no strong similarities with any other models. Among all LLMs, GPT-4 is the least similar to others, with human-written text being the least similar overall. These findings collectively underscore the nuanced patterns in text similarity across different models and underscore the complexity of language model outputs.

*4.1.2 Stylometric Analysis.* In this section, we utilize six readability statistics to evaluate the text data. Specifically, we: (1) calculate the number of words in a text that are considered difficult to read or understand, defined as those not in a predefined list of common, easy-to-understand words or contain more than two syllables; (2) apply the Flesch Reading Ease score to assess the overall readability of the text [73]; (3) estimate the time required for an average reader to read the text, based on word count and standard reading speed [74]; (4) determine the estimated reading level or grade level of the text by aggregating results from various readability tests, indicating the educational grade level necessary to comprehend the content, such as "5th grade" or "college level"; and (5) calculate number of monosyllabic words and (6) polysyllabic words.

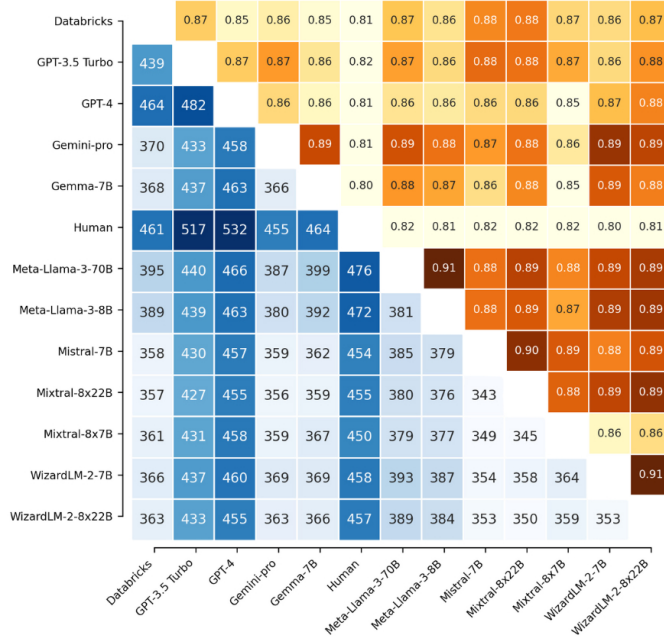| | Databricks | GPT-3.5 Turbo | GPT-4 | Gemini-pro | Gemma-7B | Human | Meta-Llama-3-70B | Meta-Llama-3-8B | Mistral-7B | Mixtral-8x22B | Mixtral-8x7B | WizardLM-2-7B | WizardLM-2-8x22B |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Databricks | | 0.87 | 0.85 | 0.86 | 0.85 | 0.81 | 0.87 | 0.86 | 0.88 | 0.88 | 0.87 | 0.86 | 0.87 |
| GPT-3.5 Turbo | 439 | | 0.87 | 0.87 | 0.86 | 0.82 | 0.87 | 0.86 | 0.88 | 0.88 | 0.87 | 0.86 | 0.88 |
| GPT-4 | 464 | 482 | | 0.86 | 0.86 | 0.81 | 0.86 | 0.86 | 0.86 | 0.86 | 0.85 | 0.87 | 0.88 |
| Gemini-pro | 370 | 433 | 458 | | 0.89 | 0.81 | 0.89 | 0.88 | 0.87 | 0.88 | 0.86 | 0.89 | 0.89 |
| Gemma-7B | 368 | 437 | 463 | 366 | | 0.80 | 0.88 | 0.87 | 0.86 | 0.88 | 0.85 | 0.89 | 0.88 |
| Human | 461 | 517 | 532 | 455 | 464 | | 0.82 | 0.81 | 0.82 | 0.82 | 0.82 | 0.80 | 0.81 |
| Meta-Llama-3-70B | 395 | 440 | 466 | 387 | 399 | 476 | | 0.91 | 0.88 | 0.89 | 0.88 | 0.89 | 0.89 |
| Meta-Llama-3-8B | 389 | 439 | 463 | 380 | 392 | 472 | 381 | | 0.88 | 0.89 | 0.87 | 0.89 | 0.89 |
| Mistral-7B | 358 | 430 | 457 | 359 | 362 | 454 | 385 | 379 | | 0.90 | 0.89 | 0.88 | 0.89 |
| Mixtral-8x22B | 357 | 427 | 455 | 356 | 359 | 455 | 380 | 376 | 343 | | 0.88 | 0.89 | 0.89 |
| Mixtral-8x7B | 361 | 431 | 458 | 359 | 367 | 450 | 379 | 377 | 349 | 345 | | 0.86 | 0.86 |
| WizardLM-2-7B | 366 | 437 | 460 | 369 | 369 | 458 | 393 | 387 | 354 | 358 | 364 | | 0.91 |
| WizardLM-2-8x22B | 363 | 433 | 455 | 363 | 366 | 457 | 389 | 384 | 353 | 350 | 359 | 353 | |

Fig. 3. Comparaison of (average) inter-text similarity. The lower part (blue) displays the similarity calculated using word-level edit distance, while the upper part (orange) illustrates the similarity determined by cosine similarity.

The obtained results are illustrated in Figure 4. Notably, when comparing the writing of LLMs to that of humans, it becomes apparent that human writing tends to be simpler, more accessible, and easier to read. Humans often use straightforward language and balanced punctuation. In contrast, models like GPT-4 and GPT-3.5 produce more complex and richly detailed content, characterized by a higher frequency of challenging vocabulary. Subsequently, the Flesch Reading Ease score, used to assess textual readability [73], indicates that human-written content is the most readable, achieving the highest score. In contrast, Gemma-7B received the lowest score, reflecting its complexity and reduced accessibility. ALso, reading time analysis [74] reveals that GPT-4 requires the longest reading duration, likely due to its lengthier and potentially more complex style, while human-generated text demands the shortest reading time. Finally, the complexity in LLM writing is further substantiated by their use of polysyllabic words, in contrast to humans who often prefer monosyllabic words for simplicity.

---

**Summary of Key Findings for** `RQ1`

This analysis reveals that human-written texts show greater variability and lower similarity compared to LLM-generated texts, emphasizing the unique nature of human writing. Among LLMs, WizardLM-2-8x22b produced the most consistent outputs, while GPT-4 exhibited the highest lexical diversity, resulting in the lowest similarity to its own outputs. In inter-text comparisons, human-written content was the least similar to LLM outputs, with Mistral being closest to human text and GPT-4 showing the greatest divergence from other models. Readability analysis further highlighted that human texts are generally simpler and more accessible, while LLM outputs, especially from GPT-4, are more complex and challenging.
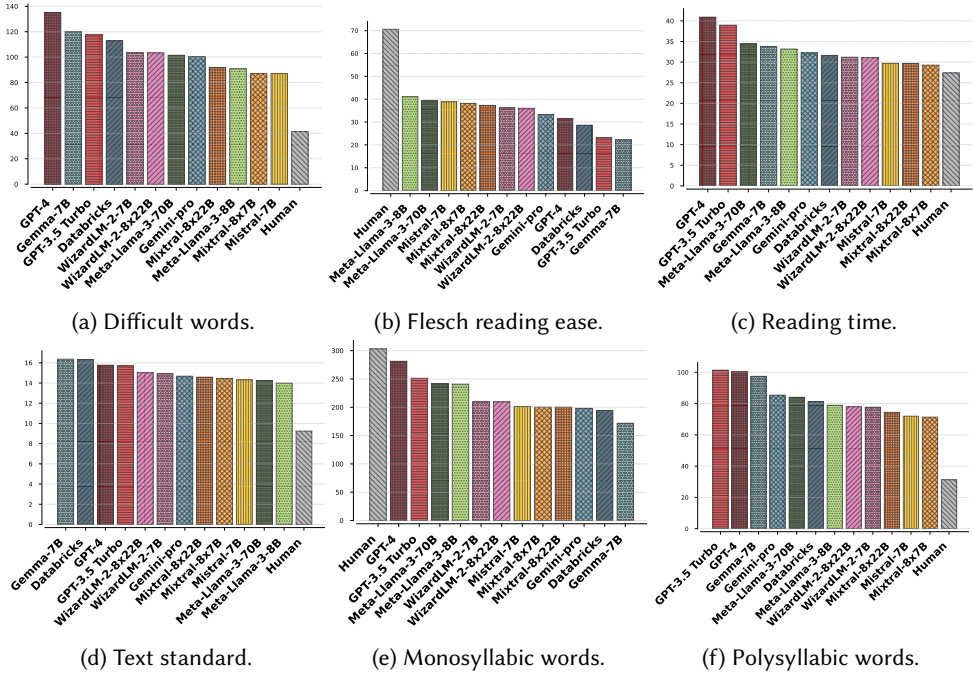
---

(a) Difficult words.

(b) Flesch reading ease.

(c) Reading time.

(d) Text standard.

(e) Monosyllabic words.

(f) Polysyllabic words.

Fig. 4. Readability Statistics.

## 4.2 RQ2: Variance in Text Generation

The inner similarity analysis presented in Figure 2 reveals significant variance in text generation among different LLMs. To illustrate this, Figure 5 employs the UMAP algorithm [75] to visualize high-dimensional text data in a comprehensible two-dimensional space, for six different prompts. This visualization offers a clear representation of the distribution and relationships among the text outputs, facilitating the identification of patterns and distinctions across the diverse LLMs under examination.

At first glance, we note that some LLMs display notable variability in text generation, whereas others exhibit a more uniform and consistent output. For instance, it is observed that Mistral 7B demonstrates low variance in text generation, while Gemini-pro-1.5 exhibits higher variance. This disparity suggests that Mistral 7B tends to produce more consistent and uniform text outputs, while Gemini-pro-1.5 introduces greater variability in its generated text. This insight into the diversity of LLM behavior is crucial for understanding the nuances of these models.

> **Summary of Key Findings for RQ2**
>
> This analysis reveals that some models produce more consistent and uniform outputs, while others show greater variability, highlighting the diversity in LLM behavior.

## 4.3 RQ3: Classification Performance

The ability to differentiate between text written by humans and that generated by LLMs holds significant importance in various contexts. In applications such as content moderation, misinformation detection, and ensuring ethical use of AI, being able to identify the source of text content is critical.

(a) Prompt 1.                              (b) Prompt 2.                              (c) Prompt 3.

(d) Prompt 4.                              (e) Prompt 5.                              (f) Prompt 6.
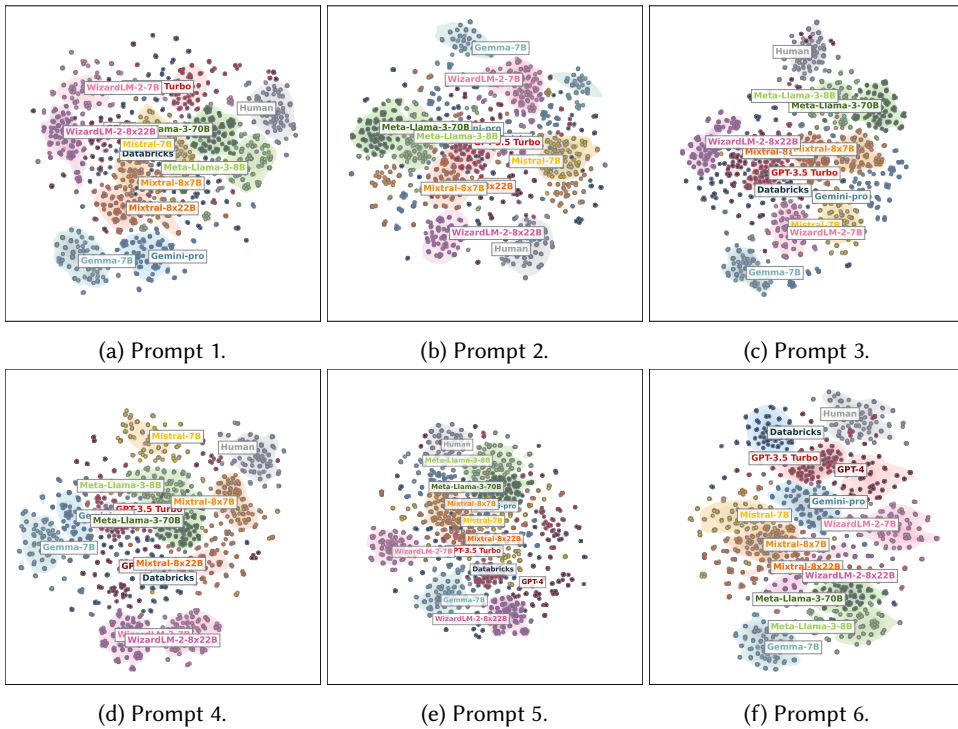
Fig. 5. UMAP projection of high-dimensional text generated into a two-dimensional space for visualization.

Furthermore, understanding which specific LLM has generated a text can provide valuable insights into the model's biases, tendencies, and potential shortcomings. This capability not only enhances transparency and accountability in AI applications but also aids in addressing ethical concerns surrounding the use of language models. Moreover, it equips users, researchers, and policymakers with the necessary tools to assess the reliability and trustworthiness of textual content, thereby encouraging a more responsible and conscious integration of LLM-generated content across various sectors.

To achieve this objective, we make use of BERT [14], several variants of DeBERTa-v3 [76] (DeBERTa-v3-xsmall, DeBERTa-v3-small, DeBERTa-v3-base), and an XGBoost [77] model with Bag of Words features. We partition our dataset into TrainVal-Test sets using a split of 60%, 20%, and 20%, respectively. The performance of these classifiers is detailed in Table 4, which presents the achieved classification metrics and reveals the following insights: **BERT** outperforms all DeBERTa variants as well as XGBoost with Bag-of-Word (BoW). BERT achieves the highest overall accuracy of **0.7095** and demonstrates superior performance across multiple LLMs, with the highest F1-scores in most categories. This highlights its robustness in identifying text generated by different models. **DeBERTa-v3-base** shows strong performance, particularly by achieving the highest F1-score (**0.7231**) for the DBRX class. Among the DeBERTa variants, it consistently outperforms the smaller models, demonstrating the effectiveness of a larger model architecture in learning various stylometric patterns. **DeBERTa-v3-small** and **DeBERTa-v3-xsmall** provide reasonable performance but lag behind BERT and the larger DeBERTa model. Finally, **XGBoost-BoW** exhibits the highest F1-score (**0.9906**) for the Human class, indicating its strength in fitting strongly to human-written text. However, it generally performs worse than the neural network-based models

in identifying texts generated by different LLMs, reflecting its limitations in handling more complex language patterns.

**Summary of Key Findings for RQ3**

This analysis suggests that distinguishing between texts written by humans and those generated by LLMs, as well as identifying which specific LLM produced a given text, appears to be a relatively straightforward task. This ability to reliably identify the source of generated text not only enhances transparency but also provides valuable insights into the unique characteristics and potential biases of each LLM, contributing to more informed and ethical use of these technologies.

## 4.4 RQ4: Language Markers Analysis

In this section, our objective is to explore the language markers, distinctive characteristics, and vocabulary exhibited by LLMs. A general method for measuring the amount of information that a feature (i.e., a word) $x_j$ provides w.r.t. predicting a class label $y$ (i.e., the LLM generating the text or the human author) is to calculate its Point-wise Mutual Information (PMI) [78]. A high PMI value indicates a more informative feature. We leverage this information to rank and select only the most positive features (words), which are then used to generate the word clouds depicted in Figure 6. The language used by different language models varies significantly, even when given the same instruction. These observations show that each language model has distinct vocabulary and linguistic styles. This diversity highlights their unique strengths and potential applications, offering valuable insights into their capabilities.



(a) Gemini-pro.  (b) Gemma-7B.  (c) GPT-4.  (d) GPT-3.5.

(e) Mistral-7B.  (f) Mixtral-8x7B.  (g) Mixtral-8x22B.  (h) Databricks.

(i) Llama-3-8B.  (j) Llama-3-70B.  (k) Wizard-2-7B.  (l) Wizard-2-8x22B.

Fig. 6. Language markers.

Table 4. Performance Comparison of AI Models in Predicting Authorship of LLM's Texts.

| Model | Accuracy | F1-Score | | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | Human | GPT-3.5 | GPT-4 | Gemini-pro | Mixtral-8x7B | Mistral-7B | Meta-Llama-3-8B | Gemma-7B | Meta-Llama-3-70B | DBRX | WizardLM-2-8x22B | WizardLM-2-7B | Mixtral-8x22B |
| bert-base-cased | **0.7095** | 0.9146 | **0.7457** | **0.7128** | **0.7970** | 0.6982 | 0.6826 | 0.6822 | **0.8762** | **0.6803** | 0.6803 | **0.7502** | **0.6393** | **0.6555** |
| DeBERTa-v3-xsmall | 0.5790 | 0.9201 | 0.6471 | 0.5022 | 0.7204 | 0.5441 | 0.4338 | 0.6038 | 0.8062 | 0.6066 | 0.6577 | 0.4702 | 0.4651 | 0.3952 |
| DeBERTa-v3-small | 0.6112 | 0.9199 | 0.6578 | 0.5594 | 0.7413 | 0.5825 | 0.4993 | 0.6459 | 0.8277 | 0.6122 | 0.6719 | 0.5509 | 0.4554 | 0.4495 |
| DeBERTa-v3-base | 0.6513 | 0.9771 | 0.7012 | 0.5937 | 0.7724 | 0.6483 | 0.6130 | 0.6673 | 0.8357 | 0.5597 | **0.7231** | 0.5524 | 0.5779 | 0.5002 |
| XGBoost-BoW | 0.5359 | **0.9906** | 0.5644 | 0.4696 | 0.5243 | 0.4786 | 0.4680 | 0.5218 | 0.6168 | 0.5539 | 0.6173 | 0.5166 | 0.5297 | 0.3868 |

Summary of Key Findings for RQ4

This analysis reveals that different LLMs exhibit distinct vocabulary and linguistic styles, as evidenced by the significant variation in language markers, making them easily recognizable and distinguishable.

## 4.5 RQ5: Bias and Ethics in LLMs

In this section, we explore whether certain LLMs adhere more closely to ethical standards by effectively reducing the propagation of biased stereotypes, thereby aligning more closely with ethical guidelines in AI development and usage.

*4.5.1 Methodology.* Analyzing bias and discrimination in LLM outputs may not be a straightforward task, as these models do not explicitly exhibit racial or gender biases and stereotypes. Therefore, one effective approach to uncovering embedded bias in these LLMs is through an *embedding-based* method.
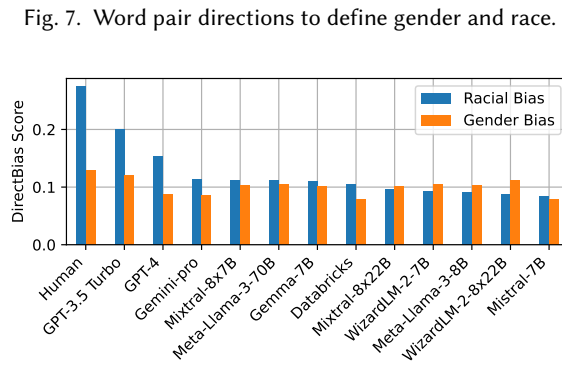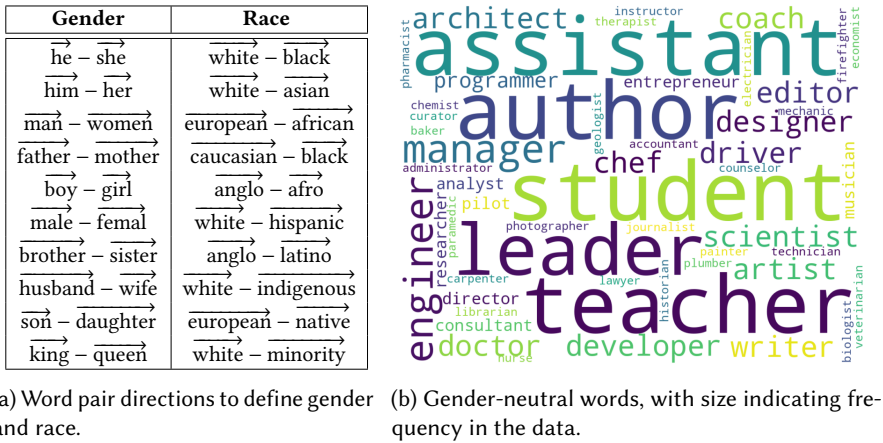
A word embedding is a representation that encodes each word $w$ as a d-dimensional vector (i.e., $w \in \mathbb{R}^d$) [11, 79]. These embeddings are trained using word co-occurrence within text corpora, leveraging paradigmatic similarity, where words with similar meanings frequently occur in similar contexts and are thus interchangeable. Consequently, words with related semantic meanings tend to have vectors that are close together in the embedding space. Moreover, the vector differences between words in these embeddings can capture the relationships between them. For example, in the analogy "man is to king as woman is to $x$", simple arithmetic on the embedding vectors reveals that the best match for $x$ is "queen", as the vector difference between "man" and "king" mirrors that between "woman" and "queen". Building on the analysis by Bolukbasi et al. [59], we aim to analyze word embeddings to identify and quantify the embedded stereotypes and biases in LLM outputs.

In our experiments, we trained 50-dimensional word embeddings for each LLM's text data using the CBOW Word2Vec architecture, with a context window of 5 words and a minimum count of 1, ensuring that all words with a total frequency lower than 1 were ignored. Once the embeddings were trained, we first evaluated them by comparing the results of various analogies against those produced by the original Word2Vec embeddings [11]. This allowed us to assess the quality and consistency of our trained embeddings in capturing semantic relationships between words. We utilized the Gensim Python library to train and evaluate our word embeddings[3].

*4.5.2 Identifying the Bias Subspace.* In this work, we focus on two primary types of bias: **gender bias** and **racial bias**. As noted by Bolukbasi et al. [59], individual word pairs do not always behave as expected because a word can have multiple meanings depending on the context. To better estimate bias, Bolukbasi et al. proposed aggregating multiple paired comparisons to more accurately identify the bias direction subspace. By combining several word pair directions, such as $\overrightarrow{she} - \overrightarrow{he}$, $\overrightarrow{woman} - \overrightarrow{man}$, $\overrightarrow{white} - \overrightarrow{black}$, and $\overrightarrow{european} - \overrightarrow{african}$, we are able to identify a significant gender or racial direction $g \in \mathbb{R}^d$ that effectively captures the underlying bias in the embedding. Formally, the bias direction subspace $g \in \mathbb{R}^d$ is estimated as follows:

$$\overrightarrow{g} = \frac{1}{|P|} \sum_{(w_1, w_2) \in P} (\overrightarrow{w_1} - \overrightarrow{w_2}) \tag{1}$$

---

[3]https://radimrehurek.com/gensim/

| Gender | Race |
|---|---|
| $\overrightarrow{he} - \overrightarrow{she}$ | $\overrightarrow{white} - \overrightarrow{black}$ |
| $\overrightarrow{him} - \overrightarrow{her}$ | $\overrightarrow{white} - \overrightarrow{asian}$ |
| $\overrightarrow{man} - \overrightarrow{women}$ | $\overrightarrow{european} - \overrightarrow{african}$ |
| $\overrightarrow{father} - \overrightarrow{mother}$ | $\overrightarrow{caucasian} - \overrightarrow{black}$ |
| $\overrightarrow{boy} - \overrightarrow{girl}$ | $\overrightarrow{anglo} - \overrightarrow{afro}$ |
| $\overrightarrow{male} - \overrightarrow{femal}$ | $\overrightarrow{white} - \overrightarrow{hispanic}$ |
| $\overrightarrow{brother} - \overrightarrow{sister}$ | $\overrightarrow{anglo} - \overrightarrow{latino}$ |
| $\overrightarrow{husband} - \overrightarrow{wife}$ | $\overrightarrow{white} - \overrightarrow{indigenous}$ |
| $\overrightarrow{son} - \overrightarrow{daughter}$ | $\overrightarrow{european} - \overrightarrow{native}$ |
| $\overrightarrow{king} - \overrightarrow{queen}$ | $\overrightarrow{white} - \overrightarrow{minority}$ |

(a) Word pair directions to define gender and race.

(b) Gender-neutral words, with size indicating frequency in the data.

Fig. 7. Word pair directions to define gender and race.



Fig. 8. DirectBias scores.

where $P$ is the list of word pair directions shown in Figure 7a, and $(w_1, w_2)$ is a pair of words in $P$. The direction represented by $g$ allows us to quantify direct biases in word associations, offering a more comprehensive understanding of how bias manifests in the embeddings.

*4.5.3 Estimating Direct Bias .* To measure direct bias, we first defined a list $N$ of 50 words that are expected to be gender-neutral, as illustrated in Figure 7b. Then, given a gender-neutral word from $N$, and the gender direction $g$ learned earlier, we estimate the direct bias of an embedding using cosine similarity, as suggested in [59]:

$$b_w = cos(\overrightarrow{w}, \overrightarrow{g}) \tag{2}$$

where a positive value of $b_w$ indicates that $w$ is more strongly associated with male, white, European, or Caucasian, while a negative value of $b_w$ suggests a stronger association with female, Black, African, or Asian. Finally, we estimate the overall direct bias of the embeddings as follows:

$$DirectBias = \frac{1}{|N|} \sum_{w \in N} |b_w| \tag{3}$$

where a lower value of *DirectBias* indicates a lower level of bias.

*4.5.4 Bias assessment.* Figure 8 shows the DirectBias score, calculated using Equation 3, for each model across different bias dimensions. At first glance, it is apparent that, overall, all models exhibit
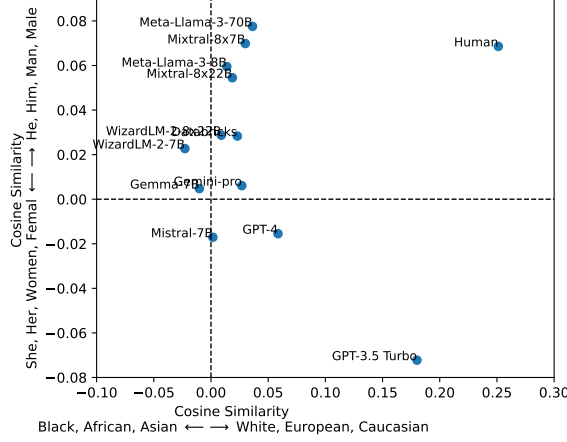
Fig. 9. Two-dimensional scatter plot of the association score between each LLM and each bias dimension.



(a) Gemini-pro.　　(b) Gemma-7B.　　(c) GPT-4.　　(d) GPT-3.5.

(e) Mistral-7B.　　(f) Mixtral-8x7B.　　(g) Mixtral-8x22B.　　(h) Databricks.

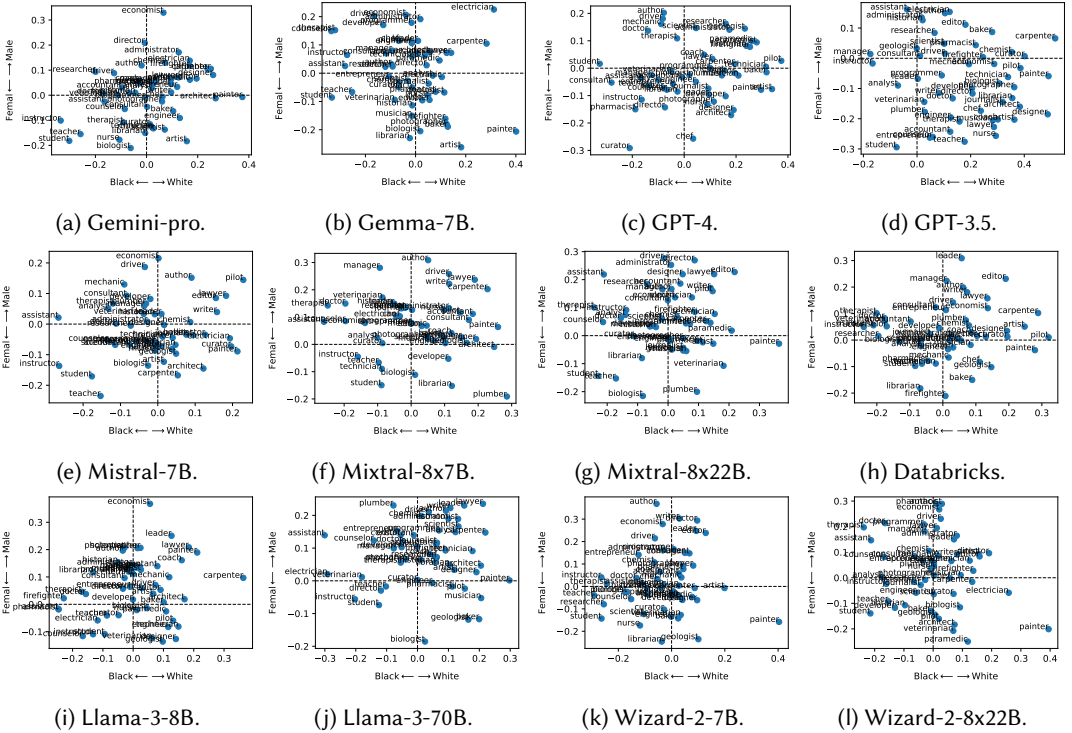(i) Llama-3-8B.　　(j) Llama-3-70B.　　(k) Wizard-2-7B.　　(l) Wizard-2-8x22B.

Fig. 10. Two-dimensional scatter plot of the association score between each occupation of Figure 7b and each bias dimension.

a relatively similar level of latent gender bias. However, human-generated texts, along with GPT-3.5 and GPT-4, appear to contain the most latent racial bias.

To further analyze which models are more strongly associated with specific bias dimensions, we refer to Figure 9, which presents the average bias scores calculated using Equation 2 for each model

across various bias dimensions. Here, we observe that some models show stronger associations with particular biases. For instance, human-generated texts tend to have a strong association with white males, whereas models like GPT-3.5 and GPT-4 exhibit a stronger latent association with white females, suggesting that these models lean more towards feminist viewpoints compared to others. Additionally, we note that Gemma-7B and Gemini-pro are positioned closer to the origin (0,0), indicating that they are the most balanced models in terms of bias.

Finally, Figure 10 presents a two-dimensional scatter plot illustrating the association score between each occupation from Figure 7b and ou two bias dimensions. Several interesting stereotypes emerge from the data. For example, in Databricks, a leader and a manager are more strongly associated with being white males, whereas in Gemini-pro, a nurse is more closely associated with being a Black female.

> **Summary of Key Findings for** `RQ5`
>
> We note that all models exhibit relatively similar levels of latent gender bias in general. However, certain models demonstrate stronger associations with specific bias dimensions. For instance, GPT-3.5 and GPT-4 show a stronger association with females, suggesting a tendency towards feminist viewpoints. Also, in terms of racial bias, GPT-3.5 and GPT-4 display the highest levels of latent racial bias, particularly associating leadership roles with white males. Finally, models like Gemma-7B and Gemini-pro are identified as the most balanced models.

## 5 CONCLUSION AND FUTURE WORK

In conclusion, our analysis sheds light on critical aspects of Large Language Models (LLMs). The observed low similarity within LLMs, distinctive inter-LLM writing styles, varying degrees of variance in text generation, successful classification outcomes, and discernible language markers underscore the nuanced and complex nature of LLM behavior. Moreover, we demonstrate that LLMs differ in their associations with gender and racial biases, with models like GPT-3.5 and GPT-4 showing a stronger tendency towards feminist viewpoints and latent racial bias. Balanced models such as Gemma-7B and Gemini-pro are identified as exhibiting the least bias overall. These findings contribute to a deeper understanding of LLM capabilities, providing valuable insights for future advancements in natural language processing and model interpretability.

Future work involves exploring explainability techniques, where the focus extends beyond detecting whether a text is authored by a human or generated by an LLM. The aim is to explore and articulate the reasons behind the model's predictions, providing a more comprehensive understanding of the decision-making process.

## REFERENCES

[1] Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Harkirat Singh Behl, Xin Wang, Sébastien Bubeck, Ronen Eldan, Adam Tauman Kalai, Yin Tat Lee, and Yuanzhi Li. Textbooks are all you need, 2023.

[2] Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Théophile Gervet, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mixtral of experts, 2024.

[3] Gemini Team. Gemini: A family of highly capable multimodal models, 2023.

[4] 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. Mistral 7b, 2023.

[5] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, Daniele Mazzotta, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo. The falcon series of open language models, 2023.

[6] Meta GenAI. Llama 2: Open foundation and fine-tuned chat models, 2023.

[7] Google. Palm 2 technical report, 2023.

[8] OpenAI. Gpt-4 technical report, 2024.

[9] Gemma Team. Gemma: Open models based on gemini research and technology, 2024.

[10] Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. Wizardlm: Empowering large language models to follow complex instructions, 2023.

[11] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013.

[12] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[13] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.

[14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *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 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

[15] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.

[16] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc., 2020.

[17] Shiyuan Huang, Siddarth Mamidanna, Shreedhar Jangam, Yilun Zhou, and Leilani H. Gilpin. Can large language models explain themselves? a study of llm-generated self-explanations, 2023.

[18] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. A survey of large language models, 2023.

[19] Yuxiang Zhou, Jiazheng Li, Yanzheng Xiang, Hanqi Yan, Lin Gui, and Yulan He. The mystery and fascination of llms: A comprehensive survey on the interpretation and analysis of emergent abilities, 2023.

[20] Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for large language models: A survey. *ACM Trans. Intell. Syst. Technol.*, jan 2024. Just Accepted.

[21] Sun et al. Trustllm: Trustworthiness in large language models, 2024.

[22] Lingbo Mo, Boshi Wang, Muhao Chen, and Huan Sun. How trustworthy are open-source llms? an assessment under malicious demonstrations shows their vulnerabilities, 2023.

[23] Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei Li, Jinfeng Bai, and Minlie Huang. Unveiling the implicit toxicity in large language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1322–1338, Singapore, December 2023. Association for Computational Linguistics.

[24] Ameet Deshpande, Vishvak Murahari, Tanmay Rajpurohit, Ashwin Kalyan, and Karthik Narasimhan. Toxicity in chatgpt: Analyzing persona-assigned language models. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1236–1270, Singapore, December 2023. Association for Computational Linguistics.

[25] Shrimai Prabhumoye, Mostofa Patwary, Mohammad Shoeybi, and Bryan Catanzaro. Adding instructions during pretraining: Effective way of controlling toxicity in language models. In Andreas Vlachos and Isabelle Augenstein, editors, *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*,

pages 2636–2651, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics.

[26] Adrian de Wynter, Xun Wang, Alex Sokolov, Qilong Gu, and Si-Qing Chen. An evaluation on large language model outputs: Discourse and memorization. *Natural Language Processing Journal*, 4:100024, September 2023.

[27] Seth Neel and Peter Chang. Privacy issues in large language models: A survey, 2024.

[28] Robin Staab, Mark Vero, Mislav Balunović, and Martin Vechev. Beyond memorization: Violating privacy via inference with large language models, 2023.

[29] Li Du, Yequan Wang, Xingrun Xing, Yiqun Ya, Xiang Li, Xin Jiang, and Xuezhi Fang. Quantifying and attributing the hallucination of large language models via association analysis, 2023.

[30] Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. Hallucination is inevitable: An innate limitation of large language models, 2024.

[31] Alessandro Bruno, Pier Luigi Mazzeo, Aladine Chetouani, Marouane Tliba, and Mohamed Amine Kerkouri. Insights into classifying and mitigating llms' hallucinations, 2023.

[32] Cheolhyoung Lee, Kyunghyun Cho, and Wanmo Kang. Mixout: Effective regularization to finetune large-scale pretrained language models. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net, 2020.

[33] Yuxiang Wu, Yu Zhao, Baotian Hu, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. An efficient memory-augmented transformer for knowledge-intensive NLP tasks. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5184–5196, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.

[34] Susmit Jha, Sumit Kumar Jha, Patrick Lincoln, Nathaniel D. Bastian, Alvaro Velasquez, and Sandeep Neema. Dehalluci-nating large language models using formal methods guided iterative prompting. In *2023 IEEE International Conference on Assured Autonomy (ICAA)*, pages 149–152, 2023.

[35] Xinlin Peng, Ying Zhou, Ben He, Le Sun, and Yingfei Sun. Hidding the ghostwriters: An adversarial evaluation of AI-generated student essay detection. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 10406–10419, Singapore, December 2023. Association for Computational Linguistics.

[36] Shangbin Feng, Herun Wan, Ningnan Wang, Zhaoxuan Tan, Minnan Luo, and Yulia Tsvetkov. What does the bot say? opportunities and risks of large language models in social media bot detection, 2024.

[37] Chujie Gao, Dongping Chen, Qihui Zhang, Yue Huang, Yao Wan, and Lichao Sun. Llm-as-a-coauthor: The challenges of detecting llm-human mixcase, 2024.

[38] R. Corizzo and S. Leal-Arenas. One-class learning for ai-generated essay detection. *Applied Sciences*, 13(13):7901, 2023.

[39] Pengyu Wang, Linyang Li, Ke Ren, Botian Jiang, Dong Zhang, and Xipeng Qiu. SeqXGPT: Sentence-level AI-generated text detection. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1144–1156, Singapore, December 2023. Association for Computational Linguistics.

[40] Harika Abburi, Michael Suesserman, Nirmala Pudota, Balaji Veeramani, Edward Bowen, and Sanmitra Bhattacharya. Generative ai text classification using ensemble llm approaches, 2023.

[41] Fatemehsadat Mireshghallah, Justus Mattern, Sicun Gao, Reza Shokri, and Taylor Berg-Kirkpatrick. Smaller language models are better black-box machine-generated text detectors, 2023.

[42] Tharindu Kumarage and Huan Liu. Neural authorship attribution: Stylometric analysis on large language models, 2023.

[43] Aparna Garimella, Akhash Amarnath, Kiran Kumar, Akash Pramod Yalla, Anandhavelu N, Niyati Chhaya, and Balaji Vasan Srinivasan. He is very intelligent, she is very beautiful? On Mitigating Social Biases in Language Modelling and Generation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4534–4545. Association for Computational Linguistics, August 2021.

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

[45] Patrick Schramowski, Cigdem Turan, Nico Andersen, Constantin A. Rothkopf, and Kristian Kersting. Large pre-trained language models contain human-like biases of what is right and wrong to do. *Nature Machine Intelligence*, 4(3):258–268, March 2022. Publisher: Nature Publishing Group.

[46] Abubakar Abid, Maheen Farooqi, and James Zou. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, page 298–306, New York, NY, USA, 2021. Association for Computing Machinery.

[47] Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. "kelly is a warm person, joseph is a role model": Gender biases in llm-generated reference letters, 2023.

[48] Sourojit Ghosh and Aylin Caliskan. Chatgpt perpetuates gender bias in machine translation and ignores non-gendered pronouns: Findings across bengali and five other low-resource languages. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '23, page 901–912, New York, NY, USA, 2023. Association for Computing Machinery.

[49] Tarek Naous, Michael J. Ryan, Alan Ritter, and Wei Xu. Having beer after prayer? measuring cultural bias in large language models, 2024.

[50] Christoph Treude and Hideaki Hata. She elicits requirements and he tests: Software engineering gender bias in large language models. In *2023 IEEE/ACM 20th International Conference on Mining Software Repositories (MSR)*, pages 624–629, 2023.

[51] Roma Patel and Ellie Pavlick. "was it "stated" or was it "claimed"?: How linguistic bias affects generative language models. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10080–10095, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.

[52] Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. "I'm sorry to hear that": Finding new biases in language models with a holistic descriptor dataset. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9180–9211, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics.

[53] Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. In-context impersonation reveals large language models'strengths and biases. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 72044–72057. Curran Associates, Inc., 2023.

[54] Marion Bartl, Malvina Nissim, and Albert Gatt. Unmasking Contextual Stereotypes: Measuring and Mitigating BERT's Gender Bias. In Marta R. Costa-jussà, Christian Hardmeier, Will Radford, and Kellie Webster, editors, *Proceedings of the Second Workshop on Gender Bias in Natural Language Processing*, pages 1–16, Barcelona, Spain (Online), December 2020. Association for Computational Linguistics.

[55] Hadas Kotek, Rikker Dockum, and David Sun. Gender bias and stereotypes in Large Language Models. In *Proceedings of The ACM Collective Intelligence Conference*, CI '23, pages 12–24, New York, NY, USA, November 2023. Association for Computing Machinery.

[56] Li Lucy and David Bamman. Gender and representation bias in GPT-3 generated stories. In Nader Akoury, Faeze Brahman, Snigdha Chaturvedi, Elizabeth Clark, Mohit Iyyer, and Lara J. Martin, editors, *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Virtual, June 2021. Association for Computational Linguistics.

[57] Nirmalendu Prakash and Roy Ka-Wei Lee. Layered Bias: Interpreting Bias in Pretrained Large Language Models. In Yonatan Belinkov, Sophie Hao, Jaap Jumelet, Najoung Kim, Arya McCarthy, and Hosein Mohebbi, editors, *Proceedings of the 6th BlackboxNLP Workshop: Analyzing and Interpreting Neural Networks for NLP*, pages 284–295, Singapore, December 2023. Association for Computational Linguistics.

[58] Moin Nadeem, Anna Bethke, and Siva Reddy. Stereoset: Measuring stereotypical bias in pretrained language models, 2020.

[59] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS'16, page 4356–4364, Red Hook, NY, USA, 2016. Curran Associates Inc.

[60] Thomas Manzini, Lim Yao Chong, Alan W Black, and Yulia Tsvetkov. Black is to criminal as Caucasian is to police: Detecting and removing multiclass bias in word embeddings. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *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 615–621, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

[61] Orestis Papakyriakopoulos, Simon Hegelich, Juan Carlos Medina Serrano, and Fabienne Marco. Bias in word embeddings. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pages 446–457, 2020.

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

[63] Marc-Etienne Brunet, Colleen Alkalay-Houlihan, Ashton Anderson, and Richard Zemel. Understanding the origins of bias in word embeddings. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 803–811. PMLR, 09–15 Jun 2019.

[64] Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. Measuring bias in contextualized word representations. In Marta R. Costa-jussà, Christian Hardmeier, Will Radford, and Kellie Webster, editors, *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 166–172, Florence, Italy, August 2019. Association for Computational Linguistics.

[65] Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, and Slav Petrov. Measuring and reducing gendered correlations in pre-trained models, 2021.

[66] Jaimeen Ahn and Alice Oh. Mitigating language-dependent ethnic bias in BERT. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 533–549, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.

[67] Haoran Zhang, Amy X. Lu, Mohamed Abdalla, Matthew McDermott, and Marzyeh Ghassemi. Hurtful words: quantifying biases in clinical contextual word embeddings. In *Proceedings of the ACM Conference on Health, Inference, and Learning*, CHIL '20, page 110–120, New York, NY, USA, 2020. Association for Computing Machinery.

[68] Shikha Bordia and Samuel R. Bowman. Identifying and Reducing Gender Bias in Word-Level Language Models. In Sudipta Kar, Farah Nadeem, Laura Burdick, Greg Durrett, and Na-Rae Han, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 7–15, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

[69] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue Wang, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekgonul, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models, 2023.

[70] Myra Cheng, Esin Durmus, and Dan Jurafsky. Marked Personas: Using Natural Language Prompts to Measure Stereotypes in Language Models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1504–1532, Toronto, Canada, July 2023. Association for Computational Linguistics.

[71] Scott Crossley, Perpetual Baffour, Tian Yu, Alex Franklin, Meg Benner, and Ulrich Boser. A large-scale corpus for assessing written argumentation: PERSUADE 2.0, August 2023.

[72] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford_alpaca, 2023.

[73] Rudolf Flesch. A new readability yardstick. *Journal of Applied Psychology*, 32(3):221–233, 1948.

[74] Vera Demberg and Frank Keller. Data from eye-tracking corpora as evidence for theories of syntactic processing complexity. *Cognition*, 109(2):193–210, Nov 2008.

[75] Leland McInnes, John Healy, Nathaniel Saul, and Lukas Großberger. Umap: Uniform manifold approximation and projection. *Journal of Open Source Software*, 3(29):861, 2018.

[76] Pengcheng He, Jianfeng Gao, and Weizhu Chen. Debertav3: Improving deberta using electra-style pre-training with gradient-disentangled embedding sharing, 2023.

[77] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16. ACM, August 2016.

[78] Kenneth Ward Church and Patrick Hanks. Word association norms, mutual information, and lexicography. *Computational Linguistics*, 16(1):22–29, 1990.

[79] Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics.