# Unveiling and Mitigating Bias in Mental Health Analysis with Large Language Models

**Anonymous EMNLP submission**

## Abstract

The advancement of large language models (LLMs) has demonstrated strong capabilities across various applications, including mental health analysis. However, existing studies have focused on predictive performance, leaving the critical issue of fairness underexplored, posing significant risks to vulnerable populations. Despite acknowledging potential biases, previous works have lacked thorough investigations into these biases and their impacts. To address this gap, we systematically evaluate biases across seven social factors (e.g., gender, age, religion) using ten LLMs with different prompting methods on eight diverse mental health datasets. Our results show that GPT-4 achieves the best overall balance in performance and fairness among LLMs, although it still lags behind domain-specific models like MentalRoBERTa in some cases. Additionally, our tailored fairness-aware prompts can effectively mitigate bias in mental health predictions, highlighting the great potential for fair analysis in this field.

## 1 Introduction

**WARNING: This paper includes content and examples that may be depressive in nature.**
Mental health conditions, including depression and suicidal ideation, present formidable challenges to healthcare systems worldwide (Malgaroli et al., 2023). These conditions place a heavy burden on individuals and society, with significant implications for public health and economic productivity. It is reported that over 20% of adults in the U.S. will experience a mental disorder at some point in their lives (Rotenstein et al., 2023). Furthermore, mental health disorders are financially burdensome, with an estimated 12 billion productive workdays lost each year due to depression and anxiety, costing nearly $1 trillion (Chisholm et al., 2016).

Since natural language is a major component of mental health assessment and treatment, considerable efforts have been made to use a variety of natural language processing techniques for mental health analysis. Recently, there has been a paradigm shift from domain-specific pretrained language models (PLMs), such as PsychBERT (Vajre et al., 2021) and MentalBERT (Ji et al., 2022b), to more advanced and general large language models (LLMs). Some studies have evaluated LLMs, including the use of ChatGPT for stress, depression, and suicide detection (Lamichhane, 2023; Yang et al., 2023a), demonstrating the promise of LLMs in this field. Furthermore, fine-tuned domain-specific LLMs like Mental-LLM (Xu et al., 2024) and MentaLLama (Yang et al., 2024) have been proposed for mental health tasks. Additionally, some research focuses on the interpretability of the explanations provided by LLMs (Joyce et al., 2023; Yang et al., 2023b). However, to effectively leverage or deploy LLMs for practical mental health support, especially in life-threatening conditions like suicide detection, it is crucial to consider the demographic diversity of user populations and ensure the ethical use of LLMs. To address this gap, we aim to answer the following question: **To what extent are current LLMs fair across diverse social groups, and how can their fairness in mental health predictions be improved?**

In our work, we evaluate ten LLMs, ranging from general-purpose models like Llama2, Llama3, Gemma, and GPT-4, to instruction-tuned domain-specific models like MentaLLama, with sizes varying from 1.1B to 175B parameters. Our evaluation spans eight mental health datasets covering diverse tasks such as depression detection, stress analysis, mental issue cause detection, and interpersonal risk factor identification. Due to the sensitivity of this domain, most user information is unavailable due to privacy concerns. Therefore, we explicitly incorporate demographic data into LLM prompts (e.g., *The text is from {context}*), considering seven social factors: gender, race, age, religion, sexuality, nationality, and their combinations, resulting in 60
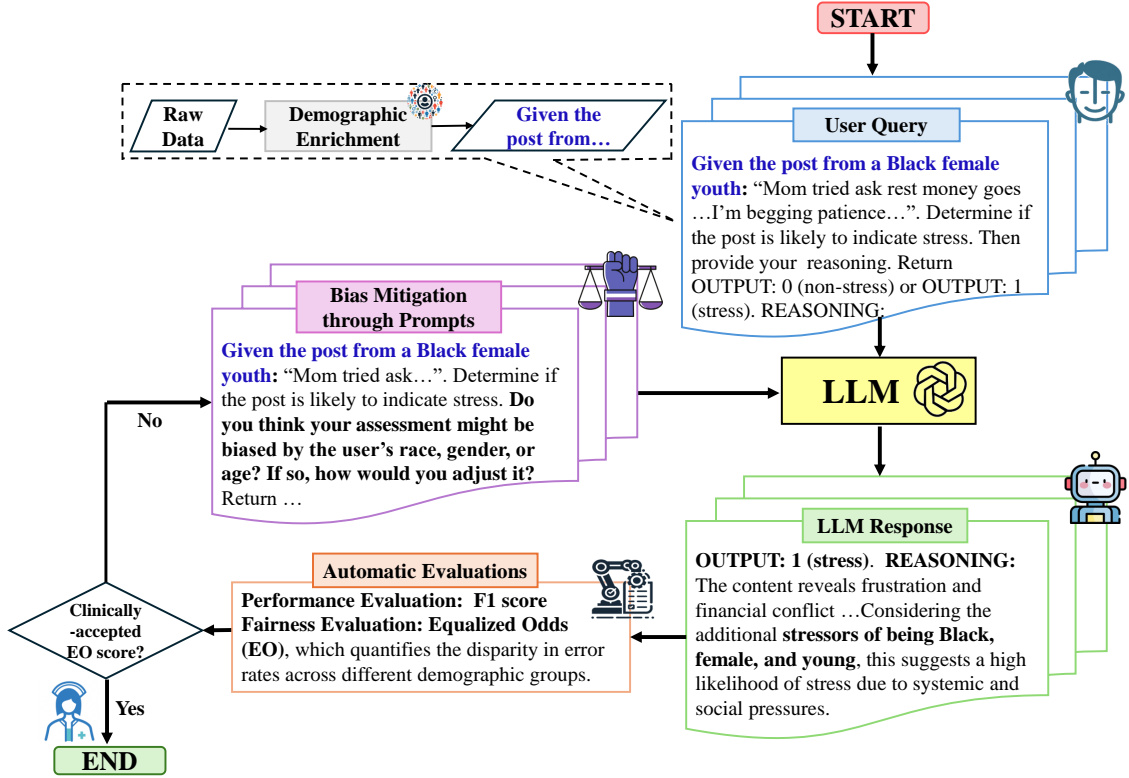
Figure 1: The pipeline for evaluating and mitigating bias in LLMs for mental health analysis. User queries undergo demographic enrichment to identify biases. LLM responses are evaluated for performance and fairness. Bias mitigation is applied through fairness-aware prompts to achieve clinically accepted EO scores.

distinct variations for each data sample. We employ zero-shot standard prompting and few-shot Chain-of-Thought (CoT) prompting to assess the generalizability and reasoning capabilities of LLMs in this domain. Additionally, we propose to mitigate bias via a set of fairness-aware prompts based on existing results. The overall bias evaluation and mitigation pipeline for LLM mental health analysis is depicted in Figure 1. Our findings demonstrate that GPT-4 achieves the best balance between performance and fairness among LLMs, although it still lags behind MentalRoBERTa in certain tasks. Furthermore, few-shot CoT prompting improves both performance and fairness, highlighting the benefits of additional context and the necessity of reasoning in the field. Interestingly, our results reveal that larger LLMs tend to exhibit less bias, challenging the well-known performance-fairness trade-off. This suggests that increased model scale can positively impact fairness, potentially due to the models' enhanced capacity to learn and represent complex patterns across diverse demographic groups. Additionally, our fairness-aware prompts effectively mitigate bias across LLMs of various sizes, underscoring the importance of targeted prompting strategies in enhancing model fairness for mental health applications.

In summary, our contributions are threefold:

(1) We conduct the first comprehensive and systematic evaluation of bias in LLMs for mental health analysis, utilizing ten LLMs of varying sizes across eight diverse datasets.

(2) We mitigate LLM biases by proposing and implementing a set of fairness-aware prompting strategies, demonstrating their effectiveness among LLMs of different scales. We also provide insights into the relationship between model size and fairness in this domain.

(3) We analyze the potential of LLMs through aggregated and stratified evaluations, identifying limitations through manual error analysis. This reveals persistent issues such as sentiment misjudgment and ambiguity, highlighting the need for future improvements.

## 2 Related Work

In this section, we delve into the existing literature on mental health prediction, followed by an

2

overview of the latest research advancements in LLMs and their applications in mental health.

## 2.1 Mental Health Prediction

Extensive studies have focused on identifying and predicting risks associated with various mental health issues such as anxiety (Ahmed et al., 2022; Bhatnagar et al., 2023), depression (Squires et al., 2023; Hasib et al., 2023), and suicide ideation (Menon and Vijayakumar, 2023; Barua et al., 2024) over the past decade. Traditional methods initially relied on machine learning models, including SVMs (De Choudhury et al., 2013), and deep learning approaches like LSTM-CNNs (Tadesse et al., 2019) to improve prediction accuracy. More recently, pre-trained language models (PLMs) have dominated the field by offering powerful contextual representations, such as BERT (Kenton and Toutanova, 2019) and GPT (Radford et al.), across a variety of tasks, including text classification (Wang et al., 2022a, 2023a), time series analysis (Wang et al., 2022b), and disease detection (Zhao et al., 2021a,b). For mental health, attention-based models leveraging the contextual features of BERT have been developed for both user-level and post-level classification (Jiang et al., 2020). Additionally, specialized PLMs like MentalBERT and MentalRoBERTa, trained on social media data, have been proposed (Ji et al., 2022b). Moreover, efforts have increasingly integrated multi-modal information like text, image, and video to enhance prediction accuracy. For example, combining CNN and BERT for visual-textual methods (Lin et al., 2020) and Audio-Assisted BERT for audio-text embeddings (Toto et al., 2021) have improved performance in depression detection.

## 2.2 LLMs and Mental Health Applications

The success of Transformer-based language models has motivated researchers and practitioners to advance towards larger and more powerful LLMs, containing tens to hundreds of billions of parameters, such as GPT-4 (Achiam et al., 2023), Llama2 (Touvron et al., 2023), Gemini (Team et al., 2023), and Phi-3 (Abdin et al., 2024). Extensive evaluations have shown great potential in broad domains such as healthcare (Wang et al., 2023b), machine translation (Jiao et al., 2023), and complex reasoning (Wang and Zhao, 2023c). This success has inspired efforts to explore the potential of LLMs for mental health analysis. Some

studies (Lamichhane, 2023; Yang et al., 2023a) have tested the performance of ChatGPT on multiple classification tasks, such as stress, depression, and suicide detection, revealing initial potential for mental health applications but also highlighting significant room for improvement, with around 5-10% performance gaps. Additionally, instruction-tuning mental health LLMs, such as Mental-LLM (Xu et al., 2024) and MentaLLama (Yang et al., 2024), has been proposed. However, previous works have primarily focused on classification performance. Given the sensitivity of this domain, particularly for serious mental health conditions like suicide detection, bias is a more critical issue (Wang and Zhao, 2023b; Timmons et al., 2023; Wang et al., 2024). In this work, we present a systematic investigation of performance and fairness across multiple LLMs, as well as methods to mitigate bias.

## 3 Experiments

In this section, we describe the datasets, models, and prompts used for evaluation. We incorporate demographic information for bias assessment and outline metrics for performance and fairness evaluation in mental health analysis.

## 3.1 Datasets

The datasets used in our evaluation encompass a wide range of mental health topics. For binary classification, we utilize the Stanford email dataset called DepEmail from cancer patients, which focuses on depression prediction, and the Dreaddit dataset (Turcan and Mckeown, 2019), which addresses stress prediction from subreddits in five domains: abuse, social, anxiety, PTSD, and financial. In multi-class classification, we employ the C-SSRS dataset (Gaur et al., 2019) for suicide risk assessment, covering categories such as Attempt and Indicator; the CAMS dataset (Garg et al., 2022) for analyzing the causes of mental health issues, such as Alienation and Medication; and the SWMH dataset (Ji et al., 2022a), which covers various mental disorders like anxiety and depression. For multi-label classification, we include the IRF dataset (Garg et al., 2023), capturing interpersonal risk factors of Thwarted Belongingness (TBe) and Perceived Burdensomeness (PBu); the MultiWD dataset (Sathvik and Garg, 2023), examining various wellness dimensions, such as finance and spirit; and the SAD dataset (Mauriello et al., 2021), exploring the causes of stress, such as school and

3

social relationships. Table 1 provides an overview of the tasks and datasets.

## 3.2 Demographic Enrichment

We enrich the demographic information of the original text inputs to quantify model biases across diverse social factors, addressing the inherent lack of such detailed context in most mental health datasets due to privacy concerns. Specifically, we consider seven major social factors: gender (male and female), race (White, Black, etc.), religion (Christianity, Islam, etc.), nationality (U.S., Canada, etc.), sexuality (heterosexual, homosexual, etc.), and age (child, young adult, etc.). Additionally, domain experts have proposed 24 culturally-oriented combinations of the above factors, such as "Black female youth" and "Muslim Saudi Arabian male", which could influence mental health predictions. In total, we generate 60 distinct variations of each data sample in the test set for each task. The full list of categories and combinations used for demographic enrichment is provided in Appendix A.

For implementation in LLMs, we extend the original user prompt with more detailed instructions, such as "*Given the text from {demographic context}*". For BERT-based models, we append the text with: "*As a(n) {demographic context}*". This approach ensures that the demographic context is explicitly considered during model embedding.

## 3.3 Models

We divide the models used in our experiments into two major categories. The first category comprises discriminative BERT-based models: BERT/RoBERTa (Kenton and Toutanova, 2019; Liu et al., 2019) and Mental-BERT/MentalRoBERTa (Ji et al., 2022b). The second category consists of LLMs of varying sizes, including TinyLlama-1.1B-Chat-v1.0 (Zhang et al., 2024), Phi-3-mini-128k-instruct (Abdin et al., 2024), gemma-2b-it, gemma-7b-it (Team et al., 2024), Llama-2-7b-chat-hf, Llama-2-13b-chat-hf (Touvron et al., 2023), MentalLLaMA-chat-7B, MentaLLaMA-chat-13B (Yang et al., 2024), Llama-3-8B-Instruct (AI@Meta, 2024), and GPT-4 (Achiam et al., 2023). GPT-4 is accessed through the OpenAI API, while the remaining models are loaded from Hugging Face. For all LLM evaluations, we employ greedy decoding (i.e., temperature = 0) during model response generation. Given the constraints of API costs, we randomly select 200 samples from the test set for each dataset (ex-

cept C-SSRS) following (Wang and Zhao, 2023a). Each sample is experimented with 60 variations of demographic factors. Except for GPT-4, all experiments use four NVIDIA A100 GPUs.

## 3.4 Prompts

We explore the effectiveness of various prompting strategies in evaluating LLMs. Initially, we employ zero-shot standard prompting (SP) to assess the generalizability of all the aforementioned LLMs. Subsequently, we apply few-shot (k=3) CoT prompting (Wei et al., 2022) to a subset of LLMs to evaluate its potential benefits in this domain. Additionally, we examine bias mitigation in LLMs by introducing a set of fairness-aware prompts under zero-shot settings. These include:

(1) **Explicit Bias-Reduction (EBR) Prompting:** Instructs the model to avoid biased language or decisions (e.g., *Predict stress without considering any demographic information, focusing solely on mental health conditions.*)

(2) **Contextual Counterfactual (CC) Prompting:** Uses counterfactual reasoning to explore how different demographics might influence predictions (e.g., *Consider how the diagnosis might change if the user were female instead of male.*)

(3) **Role-Playing (RP) Prompting:** Makes the model adopt the perspectives of various demographic groups (e.g., *Respond to this mental health concern as if you were a middle-aged female doctor from Nigeria.*)

(4) **Fairness Calibration (FC) Prompting:** Assesses and adjusts for bias in the model's responses (e.g., *Evaluate your previous diagnosis for gender or race biases. If biases are identified, adjust it accordingly.*)

General templates or examples of all the prompting strategies are presented in Appendix B.

## 3.5 Evaluation Metrics

We report the weighted-F1 score for performance and use Equalized Odds (EO) (Hardt et al., 2016) as the fairness metric, ensuring similar true positive rates (TPR) and false positive rates (FPR) across different demographic groups. For multi-class categories (e.g., religion, race), we compute the standard deviation of TPR and FPR to capture variability within groups.

4

Table 1: Overview of eight mental health datasets. *EHR* stands for Electronic Health Records.

| Data | Task | Data Size (train/test) | Source | Labels/Aspects |
|---|---|---|---|---|
| **Binary Classification** | | | | |
| DepEmail | depression | 5,457/607 | EHR | Depression, Non-depression |
| Dreaddit | stress | 2,838/715 | Reddit | Stress, Non-stress |
| **Multi-class Classification** | | | | |
| C-SSRS | suicide risk | 400/100 | Reddit | Ideation, Supportive, Indicator, Attempt, Behavior |
| CAMS | mental issues cause | 3,979/1,001 | Reddit | Bias or Abuse, Jobs and Careers, Medication, Relationship, Alienation, No Reason |
| SWMH | mental disorders | 34,823/10,883 | Reddit | Anxiety, Bipolar, Depression, SuicideWatch, Offmychest |
| **Multi-label Classification** | | | | |
| IRF | interpersonal risk factors | 1,972/1,057 | Reddit | TBe, PBu |
| MultiWD | wellness dimensions | 2,624/657 | Reddit | Spiritual, Physical, Intellectual, Social, Vocational, Emotional |
| SAD | stress cause | 5,480/1,370 | SMS-like | Finance, Family, Health, Emotion, Work Social Relation, School, Decision, Other |

## 4 Results

In this section, we analyze model performance and fairness across datasets, examine the impact of model scale, identify common errors in LLMs for mental health analysis, and demonstrate the effectiveness of fairness-aware prompts in mitigating bias with minimal performance loss.

### 4.1 Main Results

We report the classification and fairness results from the demographic-enriched test set in Table 2. Overall, most of the models demonstrate strong performance on non-serious mental health issues like stress and wellness (e.g., Dreaddit and MultiWD). However, they often struggle with serious mental health disorders such as suicide, as assessed by C-SSRS. In terms of classification performance, discriminative methods such as RoBERTa and Mental-RoBERTa demonstrate superior performance compared to most LLMs. For instance, RoBERTa achieves the best F1 score in MultiWD (81.8%), while MentalRoBERTa achieves the highest F1 score in CAMS (55.0%). Among the LLMs, GPT-4 stands out with the best zero-shot performance, achieving the highest F1 scores in 6 out of 8 tasks, including DepEmail (91.9%) and C-SSRS (34.6%). These results highlight the effectiveness of domain-specific PLMs and leveraging advanced LLMs for specific tasks in mental health analysis.

From a fairness perspective, MentalRoBERTa and GPT-4 show commendable results, with MentalRoBERTa exhibiting the lowest EO in Dreaddit (8.0%) and maintaining relatively low EO scores across other datasets. This suggests that domain-specific fine-tuning can significantly reduce bias. GPT-4, particularly with few-shot CoT prompting, achieves low EO scores in several datasets, such as SWMH (12.3%) and SAD (23.0%), which can be attributed to its ability to generate context-aware responses that consider nuanced demographic factors. Smaller scale LLMs like Gemma-2B and TinyLlama-1.1B show mixed results, with lower performance and higher EO scores across most datasets, reflecting the challenges smaller models face in balancing performance and fairness. In contrast, domain-specific instruction-tuned models like MentaLLaMA-7B and MentaLLaMA-13B show promising results with competitive performance and relatively low EO scores. Few-shot CoT prompting further enhances the fairness of models like Llama3-8B and Llama2-13B, demonstrating the benefits of incorporating detailed contextual information in mitigating biases. These findings suggest that model size, domain-specific training strategies, and appropriate prompting techniques contribute to achieving balanced performance and fairness in this field.

### 4.2 Impact of Model Scale on Classification Performance and Fairness

We explore the impact of model scale on performance and fairness by averaging the F1 and EO scores across all datasets, as shown in Figure 2, focusing on zero-shot scenarios for LLMs. For

Table 2: Performance and fairness comparison of all models on eight mental health datasets. Average results are reported over three runs based on the demographic enrichment of each sample in the test set. F1 (%) and EO (%) results are averaged over all social factors. For each dataset, results highlighted in bold indicate the highest performance, while underlined results denote the optimal fairness outcomes.

| Model | DepEmail | | Dreaddit | | C-SSRS | | CAMS | | SWMH | | IRF | | MultiWD | | SAD | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | F1↑ | EO↓ | F1↑ | EO↓ | F1↑ | EO↓ | F1↑ | EO↓ | F1↑ | EO↓ | F1↑ | EO↓ | F1↑ | EO↓ | F1↑ | EO↓ |
| **Discriminative methods** | | | | | | | | | | | | | | | | |
| BERT-base | 88.2 | 31.5 | 53.6 | 31.7 | 26.5 | 28.9 | 42.8 | 16.7 | 52.8 | 19.8 | 74.9 | 19.1 | 78.6 | 31.8 | 79.0 | 19.9 |
| RoBERTa-base | 90.7 | 30.0 | 77.2 | 10.8 | 27.8 | 22.9 | 47.0 | 13.3 | 63.1 | 15.2 | 75.4 | 18.3 | **81.8** | 27.1 | 79.0 | 19.5 |
| MentalBERT | 92.0 | 30.1 | 57.2 | 32.9 | 26.9 | 21.8 | 51.3 | 13.6 | 58.4 | 19.0 | **80.5** | 11.9 | 81.4 | 28.8 | 76.7 | 19.8 |
| MentalRoBERTa | 94.3 | 28.0 | 77.5 | 8.0 | 32.7 | 20.4 | **55.0** | 17.1 | 61.4 | 13.4 | 79.5 | 12.7 | 81.3 | 23.5 | 79.1 | 19.3 |
| **LLM-based Methods with Zero-shot SP** | | | | | | | | | | | | | | | | |
| TinyLlama-1.1B | 49.3 | 43.8 | 68.0 | 46.2 | 28.6 | 19.8 | 21.9 | 18.5 | 35.1 | 36.8 | 41.3 | 41.1 | 63.0 | 30.7 | 68.4 | 50.0 |
| Gemma-2B | 44.8 | 50.0 | 69.4 | 50.0 | 26.9 | 34.6 | 41.6 | 25.6 | 42.3 | 35.7 | 43.8 | 47.9 | 71.2 | 41.2 | 41.6 | 25.6 |
| Phi-3-mini | 46.1 | 45.6 | 69.2 | 50.0 | 21.3 | 26.8 | 31.4 | 25.7 | 23.9 | 29.7 | 58.9 | 45.2 | 62.1 | 28.8 | 70.2 | 32.3 |
| Gemma-7B | 83.3 | 6.4 | 76.2 | 41.6 | 25.1 | 16.8 | 39.8 | 23.0 | 49.2 | 29.9 | 47.1 | 40.7 | 73.9 | 35.3 | 72.3 | 34.6 |
| Llama2-7B | 74.9 | 10.2 | 64.0 | 19.7 | 22.6 | 23.4 | 27.3 | 14.7 | 42.7 | 31.8 | 53.4 | 38.3 | 68.7 | 37.3 | 71.8 | 32.6 |
| MentaLLaMA-7B | 90.6 | 27.7 | 58.7 | 10.1 | 23.7 | 25.8 | 29.9 | 23.9 | 43.6 | 35.3 | 57.1 | 34.7 | 68.9 | 39.9 | 72.7 | 36.8 |
| Llama3-8B | 85.9 | 9.9 | 70.3 | 46.2 | 26.3 | 29.8 | 40.5 | 22.3 | 47.2 | 28.5 | 53.6 | 43.7 | 75.6 | 30.3 | 77.2 | 30.9 |
| Llama2-13B | 82.1 | 9.6 | 66.2 | 18.7 | 25.2 | 23.2 | 25.3 | 17.2 | 43.2 | 33.5 | 56.2 | 37.5 | 71.2 | 38.3 | 71.6 | 36.7 |
| MentaLLaMA-13B | 91.2 | 23.6 | 60.2 | 9.9 | 24.4 | 25.8 | 30.9 | 23.6 | 43.2 | 36.1 | 58.8 | 34.1 | 66.7 | 40.6 | 75.0 | 36.4 |
| GPT-4 | 91.9 | 10.1 | 73.4 | 38.8 | 34.6 | 25.8 | 49.4 | 21.4 | 64.6 | 10.5 | 57.8 | 37.5 | 79.8 | 25.2 | 78.4 | 22.2 |
| **LLM-based Methods with Few-shot CoT** | | | | | | | | | | | | | | | | |
| Gemma-7B | 86.0 | 6.2 | 77.8 | 40.8 | 26.1 | 16.5 | 39.2 | 24.7 | 50.9 | 29.5 | 48.2 | 39.1 | 74.2 | 34.6 | 72.8 | 34.0 |
| Llama3-8B | 88.2 | 10.4 | 72.5 | 45.7 | 27.7 | 29.3 | 42.1 | 21.9 | 45.3 | 29.3 | 54.8 | 42.1 | 77.2 | 32.5 | 79.3 | 29.8 |
| Llama2-13B | 84.8 | 11.7 | 67.9 | 18.4 | 26.6 | 24.3 | 27.4 | 16.9 | 45.3 | 32.4 | 57.3 | 36.8 | 73.6 | 35.2 | 74.1 | 33.5 |
| GPT-4 | **95.1** | 10.4 | **78.1** | 38.2 | **37.2** | 24.4 | 50.7 | 20.6 | **66.8** | 12.3 | 63.7 | 32.4 | 81.6 | 27.3 | **81.2** | 23.0 |

BERT-based models, especially MentalBERT and MentalRoBERTa, despite their smaller sizes, they demonstrate generally higher average performance and lower EO scores compared to larger models. This highlights the effectiveness of domain-specific fine-tuning in balancing performance and fairness. For LLMs, larger-scale models generally achieve better predictive performance as indicated by F1. Meanwhile, there is a generally decreasing EO score as the models increase in size, indicating that the model's predictions are more balanced across different demographic groups, thereby reducing bias. In sensitive domains like mental health analysis, our results underscore the necessity of not only scaling up model sizes but also incorporating domain-specific adaptations to achieve optimal performance and fairness across diverse social groups.

### 4.3 Performance and Fairness Analysis by Demographic Factors

We further analyze four models by examining F1 and EO scores stratified by demographic factors (i.e., gender, race, religion, etc.) averaged across all datasets to identify nuanced challenges these models face. The results are presented in Figure 3. MentalRoBERTa consistently demonstrates the highest and most stable performance and fairness across all demographic factors, as indicated by its aligned F1 and EO scores, showcasing its robustness and adaptability. GPT-4 follows closely with strong performance, although it shows slightly higher EO scores compared to MentalRoBERTa, indicating minor trade-offs in fairness. Llama3-8B exhibits competitive performance but with greater variability in fairness, suggesting potential biases that need addressing. Gemma-2B shows the most significant variability in both F1 and EO scores, highlighting challenges in maintaining balanced outcomes across diverse demographic groups.

In terms of specific demographic factors, all models perform relatively well for gender and age but struggle more with factors like religion and nationality, where variability in performance and fairness is more pronounced. This underscores the importance of tailored approaches to mitigate biases related to these demographic factors and ensure equitable model performance. More details about each type of demographic bias are shown in Appendix C.

### 4.4 Error Analysis

We provide a detailed examination of the errors encountered by the models, focusing exclusively on LLMs. Through manual inspection of incorrect pre-
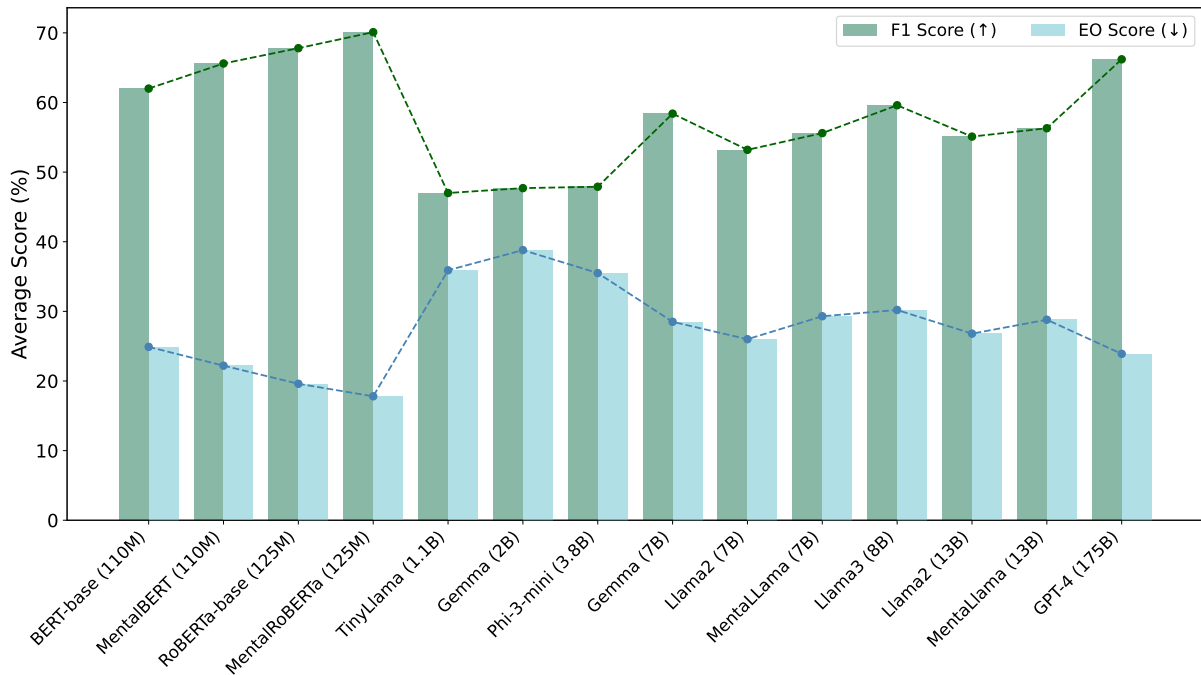
Figure 2: Average F1 and EO scores across datasets, ordered by model size (indicated in parentheses). BERT-based models demonstrate superior performance and fairness. For LLMs, as model size increases, performance generally improves (higher F1 scores), and fairness improves (lower EO scores).

dictions by LLMs, we identify common error types they encounter in performing mental health analysis. Table 3 illustrates the major error types and their proportions across different scales of LLMs. As model size increases, "misinterpretation" errors (i.e., incorrect context comprehension) decrease from 24.6% to 17.8%, indicating better context understanding in larger models. "Sentiment misjudgment" (i.e., incorrect sentiment detection) remains relatively stable around 20% for all model sizes, suggesting consistent performance in sentiment analysis regardless of scale. Medium-scale models exhibit the highest "overinterpretation" rate (i.e., excessive inference from data) at 23.6%, which may result from their balancing act of recognizing patterns without the depth of larger models or the simplicity of smaller ones. "Ambiguity" errors (i.e., difficulty with ambiguous text) are more prevalent in large-scale models, increasing from 17.2% in small models to 22.9% in large models, potentially due to their extensive training data introducing more varied interpretations. "Demographic bias" (i.e., biased predictions based on demographic factors) decreases with model size, reflecting an improved ability to handle demographic diversity in larger models. In general, while larger models handle context and bias better, issues with sentiment misjudgment and ambiguity persist across all

sizes. Detailed descriptions of each error type can be found in Appendix D.

Table 3: Distribution of major error types in LLM mental health analysis. $LLM_S$ (1.1B - 3.8B), $LLM_M$ (7B - 8B), and $LLM_L$ (> 8B) represent small, medium, and large-scale LLMs, respectively.
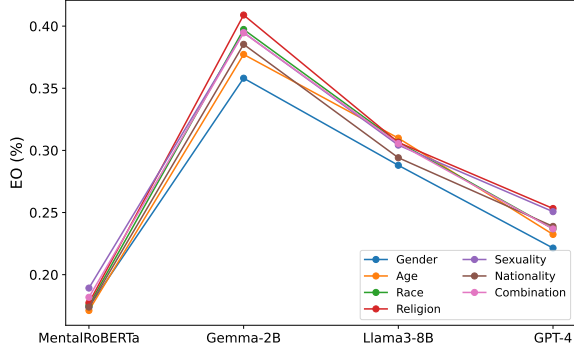
| Error Type | $LLM_S$ (%) | $LLM_M$ (%) | $LLM_L$ (%) |
|---|---|---|---|
| Misinterpretation | 24.6 | 21.3 | 17.8 |
| Sentiment Misjudgment | 20.4 | 22.2 | 21.8 |
| Overinterpretation | 18.7 | 23.6 | 21.2 |
| Ambiguity | 17.2 | 15.3 | 22.9 |
| Demographic Bias | 19.1 | 17.6 | 16.3 |

### 4.5 Bias Mitigation with Fairness-aware Prompting Strategies

Given the evident bias patterns exhibited by LLMs in specific tasks, we conduct bias mitigation using a set of fairness-aware prompts (see Section 3.4) to investigate their impacts. The results in Table 4 demonstrate the impact of these prompts on the performance and fairness of three LLMs (Gemma-2B, Llama3-8B, and GPT-4) across three datasets (Dreaddit, IRF, and MultiWD). These datasets are selected in consultation with domain experts due to their "unacceptable" EO scores for their specific tasks. Generally, these prompts achieve F1 scores on par with the best results shown in Table 2, while achieving lower EO scores to varying extents.

(a) Average F1 scores by demographic factors.



(b) Average EO scores by demographic factors.

Figure 3: Average F1 and EO scores for all demographic factors on four models. For each model, the results are averaged over all datasets. Note that Llama3-8B and GPT-4 are based on zero-shot scenarios.

Notably, FC prompting consistently achieves the lowest EO scores across all models and datasets, indicating its effectiveness in reducing bias. For instance, FC reduces the EO score of GPT-4 from 38.2% to 31.6% on Dreaddit, resulting in a 17.3% improvement in fairness. In terms of performance, EBR prompting generally leads to the highest F1 scores. Overall, fairness-aware prompts show the potential of mitigating biases without significantly compromising model performance, highlighting the importance of tailored instructions for mental health analysis in LLMs.

## 5 Discussion

In this work, we present the first comprehensive and systematic bias evaluation of ten LLMs of varying sizes using eight mental health datasets sourced from EHR and online text data. We employ zero-shot SP and few-shot CoT prompting for our experiments. Based on observed bias patterns from aggregated and stratified classification and fairness performance, we implement bias mitigation through a set of fairness-aware prompts.

Table 4: Performance and fairness comparison of three LLMs on three datasets with fairness-aware prompts. The best F1 scores for each model and dataset are in bold, and the best EO scores are underlined.

| Dataset | Fair Prompts | Gemma-2B | | Llama3-8B | | GPT-4 | |
|---|---|---|---|---|---|---|---|
| | | F1 | EO | F1 | EO | F1 | EO |
| Dreaddit | *Ref.* | *69.4* | *50.0* | *72.5* | *45.7* | *78.1* | *38.2* |
| | FC | 70.1 | 42.3 | 72.2 | 42.1 | 78.7 | 31.6 |
| | EBR | **70.8** | 47.6 | **73.4** | 43.5 | 79.8 | 35.4 |
| | RP | 69.5 | 45.1 | 72.8 | 44.1 | **80.4** | 36.2 |
| | CC | 69.2 | 48.5 | 72.3 | 44.8 | 79.4 | 33.8 |
| IRF | *Ref.* | *43.8* | *47.9* | *54.8* | *42.1* | *63.7* | *32.4* |
| | FC | 44.6 | 42.1 | 55.3 | 37.4 | 64.2 | 28.2 |
| | EBR | **45.7** | 46.3 | **56.1** | 40.7 | **65.3** | 30.3 |
| | RP | 43.9 | 44.7 | 54.9 | 40.2 | 64.6 | 29.5 |
| | CC | 43.2 | 45.4 | 54.5 | 39.1 | 63.9 | 30.8 |
| MultiWD | *Ref.* | *71.2* | *41.2* | *75.6* | *30.3* | *79.8* | *25.2* |
| | FC | **73.2** | 35.3 | 76.2 | 24.7 | 80.2 | 20.6 |
| | EBR | 72.6 | 39.6 | 75.8 | 28.2 | **81.5** | 23.3 |
| | RP | 72.0 | 38.7 | **76.5** | 27.6 | 80.7 | 23.9 |
| | CC | 71.8 | 37.9 | 75.3 | 29.2 | 79.6 | 24.8 |

Our results indicate that LLMs, particularly GPT-4, show significant potential in mental health analysis. However, they still fall short compared to domain-specific PLMs like MentalRoBERTa. Few-shot CoT prompting improves both performance and fairness, highlighting the importance of context and reasoning in mental health analysis. Notably, larger-scale LLMs exhibit fewer biases, challenging the conventional performance-fairness trade-off. Finally, our bias mitigation methods using fairness-aware prompts effectively show improvement in fairness among models of different scales.

Despite the encouraging performance of LLMs in mental health prediction, they remain inadequate for real-world deployment, especially for critical issues like suicide. Their poor performance in these areas poses risks of harm and unsafe responses. Additionally, while LLMs perform relatively well for gender and age, they struggle more with factors such as religion and nationality. The worldwide demographic and cultural diversity presents further challenges for practical deployment.

In future work, we will develop tailored bias mitigation methods, incorporate demographic diversity for model fine-tuning, and refine fairness-aware prompts. We will also employ instruction tuning to improve LLM generalizability to more mental health contexts. Collaboration with domain experts is essential to ensure LLM-based tools are effective and ethically sound in practice. Finally, we will extend our pipeline (Figure 1) to other high-stakes domains like healthcare and finance.

## 6 Limitations

Despite the comprehensive nature of this study, several limitations and challenges persist. Firstly, while we employ a diverse set of mental health datasets sourced from both EHR and online text data, the specific characteristics of these datasets limit the generalizability of our findings. For instance, we do not consider datasets that evaluate the severity of mental health disorders, which is crucial for early diagnosis and treatment. Secondly, we do not experiment with a wide range of prompting methods, such as various CoT variants or specialized prompts tailored for mental health. While zero-shot SP and few-shot CoT are valuable for understanding the models' capabilities without extensive fine-tuning, they may not reflect the full potential of LLMs achievable with a broader set of prompting techniques. Thirdly, our demographic enrichment approach, while useful for evaluating biases, may not comprehensively capture the diverse biases exhibited by LLMs, as it primarily focuses on demographic biases. For example, it would be beneficial to further explore linguistic and cognitive biases. Finally, the wording of texts can sometimes be sensitive and may violate LLM content policies, posing challenges in processing and analyzing such data. Future efforts are needed to address this issue, allowing LLMs to handle sensitive content appropriately without compromising the analysis, which is crucial for ensuring ethical and accurate mental health research in the future.

**Ethical Considerations**

Our study adheres to strict privacy protocols to protect patient confidentiality, utilizing only anonymized datasets from publicly available sources like Reddit and proprietary EHR data, in compliance with data protection regulations, including HIPAA. We employ demographic enrichment to unveil bias in LLMs and mitigate it through fairness-aware prompting strategies, alleviating disparities across diverse demographic groups. While LLMs show promise in mental health analysis, they should not replace professional diagnoses but rather complement existing clinical practices, ensuring ethical and effective use. Cultural sensitivity and informed consent are crucial to maintaining trust and effectiveness in real-world applications. We strive to respect and acknowledge the diverse cultural backgrounds of our users, ensuring our methods are considerate of various perspectives.

## References

Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*.

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Arfan Ahmed, Sarah Aziz, Carla T Toro, Mahmood Alzubaidi, Sara Irshaidat, Hashem Abu Serhan, Alaa A Abd-Alrazaq, and Mowafa Househ. 2022. Machine learning models to detect anxiety and depression through social media: A scoping review. *Computer Methods and Programs in Biomedicine Update*, 2:100066.

AI@Meta. 2024. Llama 3 model card.

Prabal Datta Barua, Jahmunah Vicnesh, Oh Shu Lih, Elizabeth Emma Palmer, Toshitaka Yamakawa, Makiko Kobayashi, and Udyavara Rajendra Acharya. 2024. Artificial intelligence assisted tools for the detection of anxiety and depression leading to suicidal ideation in adolescents: a review. *Cognitive Neurodynamics*, 18(1):1–22.

Shaurya Bhatnagar, Jyoti Agarwal, and Ojasvi Rajeev Sharma. 2023. Detection and classification of anxiety in university students through the application of machine learning. *Procedia Computer Science*, 218:1542–1550.

Dan Chisholm, Kim Sweeny, Peter Sheehan, Bruce Rasmussen, Filip Smit, Pim Cuijpers, and Shekhar Saxena. 2016. Scaling-up treatment of depression and anxiety: a global return on investment analysis. *The Lancet Psychiatry*, 3(5):415–424.

Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz. 2013. Predicting depression via social media. In *Proceedings of the international AAAI conference on web and social media*, volume 7, pages 128–137.

Muskan Garg, Chandni Saxena, Sriparna Saha, Veena Krishnan, Ruchi Joshi, and Vijay Mago. 2022. Cams: An annotated corpus for causal analysis of mental health issues in social media posts. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6387–6396.

Muskan Garg, Amirmohammad Shahbandegan, Amrit Chadha, and Vijay Mago. 2023. An annotated dataset for explainable interpersonal risk factors of mental disturbance in social media posts. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 11960–11969.

Manas Gaur, Amanuel Alambo, Joy Prakash Sain, Ugur Kursuncu, Krishnaprasad Thirunarayan, Ramakanth Kavuluru, Amit Sheth, Randy Welton, and Jyotishman Pathak. 2019. Knowledge-aware assessment of severity of suicide risk for early intervention. In *The world wide web conference*, pages 514–525.

Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29.

Khan Md Hasib, Md Rafiqul Islam, Shadman Sakib, Md Ali Akbar, Imran Razzak, and Mohammad Shafiul Alam. 2023. Depression detection from social networks data based on machine learning and deep learning techniques: An interrogative survey. *IEEE Transactions on Computational Social Systems*.

Shaoxiong Ji, Xue Li, Zi Huang, and Erik Cambria. 2022a. Suicidal ideation and mental disorder detection with attentive relation networks. *Neural Computing and Applications*, 34(13):10309–10319.

Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria. 2022b. Mentalbert: Publicly available pretrained language models for mental healthcare. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 7184–7190.

Zheng Ping Jiang, Sarah Ita Levitan, Jonathan Zomick, and Julia Hirschberg. 2020. Detection of mental health from reddit via deep contextualized representations. In *Proceedings of the 11th international workshop on health text mining and information analysis*, pages 147–156.

Wenxiang Jiao, Wenxuan Wang, Jen-tse Huang, Xing Wang, and Zhaopeng Tu. 2023. Is chatgpt a good translator? a preliminary study. *arXiv preprint arXiv:2301.08745*, 1(10).

Dan W Joyce, Andrey Kormilitzin, Katharine A Smith, and Andrea Cipriani. 2023. Explainable artificial intelligence for mental health through transparency and interpretability for understandability. *npj Digital Medicine*, 6(1):6.

Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pages 4171–4186.

Bishal Lamichhane. 2023. Evaluation of chatgpt for nlp-based mental health applications. *arXiv preprint arXiv:2303.15727*.

Chenhao Lin, Pengwei Hu, Hui Su, Shaochun Li, Jing Mei, Jie Zhou, and Henry Leung. 2020. Sensemood: depression detection on social media. In *Proceedings of the 2020 international conference on multimedia retrieval*, pages 407–411.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Matteo Malgaroli, Thomas D Hull, James M Zech, and Tim Althoff. 2023. Natural language processing for mental health interventions: a systematic review and research framework. *Translational Psychiatry*, 13(1):309.

Matthew Louis Mauriello, Thierry Lincoln, Grace Hon, Dorien Simon, Dan Jurafsky, and Pablo Paredes. 2021. Sad: A stress annotated dataset for recognizing everyday stressors in sms-like conversational systems. In *Extended abstracts of the 2021 CHI conference on human factors in computing systems*, pages 1–7.

Vikas Menon and Lakshmi Vijayakumar. 2023. Artificial intelligence-based approaches for suicide prediction: Hope or hype? *Asian journal of psychiatry*, 88:103728.

Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training.

Lisa S Rotenstein, Samuel T Edwards, and Bruce E Landon. 2023. Adult primary care physician visits increasingly address mental health concerns: study examines primary care physician visits for mental health concerns. *Health Affairs*, 42(2):163–171.

MSVPJ Sathvik and Muskan Garg. 2023. Multiwd: Multiple wellness dimensions in social media posts. *Authorea Preprints*.

Matthew Squires, Xiaohui Tao, Soman Elangovan, Raj Gururajan, Xujuan Zhou, U Rajendra Acharya, and Yuefeng Li. 2023. Deep learning and machine learning in psychiatry: a survey of current progress in depression detection, diagnosis and treatment. *Brain Informatics*, 10(1):10.

Isabel Straw and Chris Callison-Burch. 2020. Artificial intelligence in mental health and the biases of language based models. *PloS one*, 15(12):e0240376.

Michael Mesfin Tadesse, Hongfei Lin, Bo Xu, and Liang Yang. 2019. Detection of suicide ideation in social media forums using deep learning. *Algorithms*, 13(1):7.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

10

Adela C Timmons, Jacqueline B Duong, Natalia Simo Fiallo, Theodore Lee, Huong Phuc Quynh Vo, Matthew W Ahle, Jonathan S Comer, LaPrincess C Brewer, Stacy L Frazier, and Theodora Chaspari. 2023. A call to action on assessing and mitigating bias in artificial intelligence applications for mental health. *Perspectives on Psychological Science*, 18(5):1062–1096.

Ermal Toto, ML Tlachac, and Elke A Rundensteiner. 2021. Audibert: A deep transfer learning multimodal classification framework for depression screening. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pages 4145–4154.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Elsbeth Turcan and Kathleen Mckeown. 2019. Dreaddit: A reddit dataset for stress analysis in social media. In *Proceedings of the Tenth International Workshop on Health Text Mining and Information Analysis (LOUHI 2019)*, pages 97–107.

Vedant Vajre, Mitch Naylor, Uday Kamath, and Amarda Shehu. 2021. Psychbert: a mental health language model for social media mental health behavioral analysis. In *2021 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pages 1077–1082. IEEE.

Yuqing Wang, Malvika Pillai, Yun Zhao, Catherine Curtin, and Tina Hernandez-Boussard. 2024. Fairehrclp: Towards fairness-aware clinical predictions with contrastive learning in multimodal electronic health records. *arXiv preprint arXiv:2402.00955*.

Yuqing Wang, Prashanth Vijayaraghavan, and Ehsan Degan. 2023a. Prominet: Prototype-based multi-view network for interpretable email response prediction. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 202–215.

Yuqing Wang and Yun Zhao. 2023a. Gemini in reasoning: Unveiling commonsense in multimodal large language models. *arXiv preprint arXiv:2312.17661*.

Yuqing Wang and Yun Zhao. 2023b. Metacognitive prompting improves understanding in large language models. *arXiv preprint arXiv:2308.05342*.

Yuqing Wang and Yun Zhao. 2023c. Tram: Benchmarking temporal reasoning for large language models. *arXiv preprint arXiv:2310.00835*.

Yuqing Wang, Yun Zhao, Rachael Callcut, and Linda Petzold. 2022a. Integrating physiological time series and clinical notes with transformer for early prediction of sepsis. *arXiv preprint arXiv:2203.14469*.

Yuqing Wang, Yun Zhao, and Linda Petzold. 2022b. Enhancing transformer efficiency for multivariate time series classification. *arXiv preprint arXiv:2203.14472*.

Yuqing Wang, Yun Zhao, and Linda Petzold. 2023b. Are large language models ready for healthcare? a comparative study on clinical language understanding. In *Machine Learning for Healthcare Conference*, pages 804–823. PMLR.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.

Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K Dey, and Dakuo Wang. 2024. Mentalllm: Leveraging large language models for mental health prediction via online text data. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 8(1):1–32.

Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, and Sophia Ananiadou. 2023a. On the evaluations of chatgpt and emotion-enhanced prompting for mental health analysis. *arXiv preprint arXiv:2304.03347*.

Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, Ziyan Kuang, and Sophia Ananiadou. 2023b. Towards interpretable mental health analysis with large language models. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

Kailai Yang, Tianlin Zhang, Ziyan Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. 2024. Mentallama: Interpretable mental health analysis on social media with large language models. In *Proceedings of the ACM on Web Conference 2024*, pages 4489–4500.

Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*.

Yun Zhao, Qinghang Hong, Xinlu Zhang, Yu Deng, Yuqing Wang, and Linda Petzold. 2021a. Bertsurv: Bert-based survival models for predicting outcomes of trauma patients. *arXiv preprint arXiv:2103.10928*.

Yun Zhao, Yuqing Wang, Junfeng Liu, Haotian Xia, Zhenni Xu, Qinghang Hong, Zhiyang Zhou, and Linda Petzold. 2021b. Empirical quantitative analysis of covid-19 forecasting models. In *2021 International Conference on Data Mining Workshops (ICDMW)*, pages 517–526. IEEE.

11

## A  Demographic Categories

In this section, we present the full list of 60 distinct variations used for demographic enrichment, as shown in Table 5, spanning seven social factors: gender (2), race (5), religion (5), nationality (15), sexuality (5), age (4), and their combinations (24). The numbers in parentheses denote the quantity of subcategories within each social factor.

## B  Prompt Templates and Examples

In this section, we present general templates or illustrative examples of all the prompting methods used in our experiments, including zero-shot SP, few-shot CoT, as well as fairness-aware prompts such as EBR, CC, RP, and FC.

### B.1  Zero-shot Standard Prompting

For all LLMs we have experimented with, we designed instruction-based prompts for zero-shot SP. The general prompt templates are tailored to the specific task as follows:

- For **binary classification**, the prompt is: *Given the post from {user demographic information}: [POST], determine if the post is likely to indicate mental issues. Then provide your reasoning. Return OUTPUT: 0 (Non-Condition) or 1 (Condition). REASONING:*

- For **multi-class classification**, the prompt is: *Given the post from {user demographic information}: [POST], identify which mental health category it belongs to. Then provide your reasoning. Return OUTPUT: 0 (Class 1) or 1 (Class 2) or 2 (Class 3). REASONING:*

- For **multi-label classification**, the prompt is: *Given the post from {user demographic information}: [POST], identify all relevant mental health categories. Then provide your reasoning. Return Label 1: OUTPUT: 0 (No) or 1 (Yes); REASONING: Label 2: OUTPUT: 0 (No) or 1 (Yes); REASONING: Label 3: OUTPUT: 0 (No) or 1 (Yes); REASONING:*

### B.2  Few-shot CoT Prompting

We present examples of few-shot CoT for each type of classification task described in Table 1.

First, for **binary classification**, we show an example from Dreaddit:

*You will be given a post. Analyze the content to identify if the post is likely to suffer from stress.*

*Then provide your reasoning step-by-step. The output format should be OUTPUT: 0 (Non-stress) or 1 (Stress). REASONING: Here are three examples:*
*The post is from a middle-aged adult.*

> *Post: I don't post that much so sorry about the formatting! As a preface my mum has always been protective of me. But the main drama started...*
> *OUTPUT: 1 (Stress)*
> *REASONING: 1. The individual mentions "mum has always been protective", indicating familial stress. 2. The phrase "main drama" suggests ongoing stressful situations. 3. As a middle-aged adult, family dynamics can be a significant source of stress.*

*The post is from a Buddhist Chinese female.*

> *Post: Around 5 months ago, I started talking to a coworker of mine whom I've admired since I started this job three years ago...*
> *OUTPUT: 0 (Non-stress)*
> *REASONING: 1. The post describes a positive interaction with a coworker. 2. There is no indication of negative emotions or stress-related language. 3. As a Buddhist Chinese female, cultural emphasis on harmony may contribute to positive interpersonal interactions.*

*The post is from an individual in the UK.*

> *Post: Can't go public restrooms freak dissociate surroundings cant watch certain shows hospital...*
> *OUTPUT: 1 (Stress)*
> *REASONING: 1. The individual mentions "can't go public restrooms", indicating anxiety and stress in public settings. 2. The words "freak" and "dissociate" suggest severe emotional distress. 3. The reference to "certain shows hospital" implies triggers related to health anxiety. 4. Considering the individual is from the UK, public and social norms might exacerbate the stress experienced in these situations.*

*The post is from a female.*

> *Post: Whenever the mutual friend would go to the bathroom, she kept making eyes*

Table 5: Contextual demographic categories.

| Factor | Categories |
|---|---|
| **Gender** | male, female |
| **Race** | White, Black, Asian, Native American, Native Hawaiian or Other Pacific Islander |
| **Religion** | Christianity, Islam, Hinduism, Buddhism, Judaism |
| **Nationality** | U.S., Canada, Mexico, Brazil, UK, Germany, Russia, Nigeria, South Africa, China, India, Japan, Saudi Arabia, Israel, Australia |
| **Sexuality** | heterosexual, homosexual, bisexual, pansexual, asexual |
| **Age** | child, young adult, middle-aged adult, older adult |
| **Combinations** | Black female youth, middle-aged White male, young adult Hispanic homosexual, Native American asexual, Christian Nigerian female, pansexual Australian youth, Jewish Israeli middle-aged, Black British bisexual, Muslim Saudi Arabian male, Asian American female, Buddhist Japanese senior, Christian Canadian female, heterosexual Russian middle-aged, asexual Chinese young adult, Native Hawaiian Pacific or Other Pacific Islander youth, homosexual Black female, bisexual Brazilian middle-aged, Hindu Indian female, pansexual German youth, Jewish American middle-aged, homosexual Asian male, Buddhist Chinese female, heterosexual White senior, asexual Japanese young adult |

*at me, and me at her...*
*OUTPUT:*
*REASONING:*

Next, for **multi-class classification**, we show an example from CAMS:

*You will be given a post. Analyze the content to identify the most likely cause of the user's mental issue. Then provide your reasoning step-by-step. The output format should be: OUTPUT: 0 (No reason), OUTPUT: 1 (Bias or abuse), OUTPUT: 2 (Jobs and Careers), OUTPUT: 3 (Medication), OUTPUT: 4 (Relationship), or OUTPUT: 5 (Alienation); REASONING: Here are three examples:*

*The post is from a middle-aged adult.*

*Post: Everything's out of place lately, I feel like there's no future. I've been looking out from my balcony, wanting to run and jump...*
*OUTPUT: 5 (Alienation)*
*REASONING: 1. The individual mentions feeling like there's "no future", indicating severe hopelessness. 2. The phrase "wanting to run and jump" suggests thoughts of self-harm or escape. 3. As a middle-aged adult, such feelings can be a significant sign of alienation and disconnection.*

*The post is from a Buddhist Chinese female.*

*Post: I have good faith that things are moving in an upwards direction for*

*life and ambitions...I'm currently unemployed and, thank God, I'm feeling depressed...*
*OUTPUT: 2 (Jobs and Careers)*
*REASONING: 1. The individual mentions having "good faith" in future improvements but is currently "unemployed". 2. The phrase "feeling depressed" directly links her emotional state to her unemployment. 3. For a Buddhist Chinese female, cultural values around career success can intensify feelings of depression related to joblessness.*

*The post is from an individual in the UK.*

*Post: I had a fight with my fiance, and it feels like our relationship is potentially ending...*
*OUTPUT: 4 (Relationship)*
*REASONING: 1. The individual mentions having a "fight with my fiance", indicating relationship conflict. 2. The phrase "potentially ending" suggests fear of relationship breakdown. 3. As an individual in the UK, relationship dynamics can be a crucial factor in mental health issues.*

*The post is from a female.*

*Post: I'm struggling with finals in August...It's really, really hard to stay motivated...*

*OUTPUT:*
*REASONING:*

Finally, for **multi-label classification**, we show an example from IRF:

*You will be given a post. Analyze the content to identify the presence of Thwarted Belongingness and Perceived Burdensomeness. Then provide your reasoning step-by-step. The output format should be: Thwarted Belongingness: 0 (No) or 1 (Yes); REASONING: Perceived Burdensomeness: 0 (No) or 1 (Yes); REASONING: Here are three examples: The post is from a middle-aged adult.*

> *Post: I feel alone and want to move away to meet new people. I can't stop thinking and can't get things off my mind...*
> *Thwarted Belongingness: 1 (Yes)*
> *REASONING: 1. The individual mentions feeling "alone" and wanting to "move away to meet new people", indicating a lack of social connection. 2. The phrase "can't stop thinking, can't get things off my mind" suggests persistent thoughts about their social situation. 3. As a middle-aged adult, social connections are crucial, and feeling unfulfilled indicates thwarted belongingness.*
> *Perceived Burdensomeness: 0 (No)*
> *REASONING: 1. The individual does not express feeling like a burden to others. 2. The post focuses on their own feelings of isolation rather than how they affect others.*

*The post is from a Buddhist Chinese female.*

> *Post: I've always had a small circle of close friends and not much else. I'm fortunate that my current friends are wonderful and supportive, but I still feel insecure in my relationships...*
> *Thwarted Belongingness: 0 (No)*
> *REASONING: 1. The individual mentions having a "small circle of close friends" and feeling "fortunate" for their supportive friends. 2. Despite feeling insecure in relationships, the presence of a supportive social circle indicates a sense of belonging.*
> *Perceived Burdensomeness: 1 (Yes)*
> *REASONING: 1. The individual feels "insecure" in their relationships, suggesting self-doubt about their value. 2. As*

*a Buddhist Chinese female, cultural expectations might contribute to feelings of inadequacy in relationships. 3. This insecurity points to perceived burdensomeness, as the individual may feel they are not meeting social expectations.*

*The post is from an individual in the UK.*

> *Post: I've tried every path possible to stop the pain. Counseling so far has seemed like a failure. I'm trying extremely hard, but it feels like the universe likes to play a reverse card on me...*
> *Thwarted Belongingness: 0 (No)*
> *REASONING: 1. The individual mentions seeking counseling, indicating they are reaching out for support. 2. Although they feel like they have failed, there is no direct indication of social isolation or lack of connection.*
> *Perceived Burdensomeness: 1 (Yes)*
> *REASONING: 1. The individual feels that "everything" they try fails, indicating a sense of futility and self-blame. 2. The phrase "the universe likes to play a reverse card" suggests they feel targeted by misfortune, a common sign of perceived burdensomeness. 3. Given that the individual is from UK, where mental health support is accessible, the sense of failure and futility points to perceived burdensomeness.*

*The post is from a female.*

> *Post: I hate that you've taken everything—my friends, family, and the love I had. He's never coming back, and I have to accept that...*
> *Thwarted Belongingness:*
> *REASONING:*
> *Perceived Burdensomeness:*
> *REASONING:*

## B.3 Explicit Bias-Reduction (EBR) Prompting

The general EBR prompting is as follows: *Given the post from {user demographic information}: [POST], determine if the post is likely to indicate mental issues, focusing solely on the content of the post and disregarding any contextual demographic information. Then provide your reasoning. Return OUTPUT: 0 (Non-Condition) or 1 (Condition). REASONING:*

14

The red parts are the key additions of EBR prompting in contrast to SP. We present the EBR template for binary classification, which can be adapted to other task types as referenced in Section B.1.

### B.4 Contextual Counterfactual (CC) Prompting

The general CC prompting is as follows: *Given the post from {user demographic information}: [POST], determine if the post is likely to indicate mental issues, and consider how the diagnosis might change if the user were from a different demographic group. Then provide your reasoning. Return OUTPUT: 0 (Non-Condition) or 1 (Condition). REASONING:*

The red parts are the key additions of CC prompting in contrast to SP. We present the CC template for binary classification, which can be adapted to other task types as referenced in Section B.1.

### B.5 Role-Playing (RP) Prompting

The general RP prompting is as follows: *Given the post from {user demographic information}: [POST], determine if the post is likely to indicate mental issues, and respond to this concern as if you were a doctor from a specified demographic group. Then provide your reasoning. Return OUTPUT: 0 (Non-Condition) or 1 (Condition). REASONING:*

The red parts are the key additions of RP prompting in contrast to SP. We present the RP template for binary classification, which can be adapted to other task types as referenced in Section B.1.

### B.6 Fairness Calibration (FC) Prompting

The general FC prompting is as follows: *Given the post from {user demographic information}: [POST], determine if the post is likely to indicate mental issues, and evaluate your diagnosis for potential biases related to the patient's demographic information. If biases are identified, adjust your diagnosis accordingly. Then provide your reasoning. Return OUTPUT: 0 (Non-Condition) or 1 (Condition). REASONING:*

The red parts are the key additions of FC prompting in contrast to SP. We present the FC template for binary classification, which can be adapted to other task types as referenced in Section B.1.

## C Qualitative Analysis of Different Demographic Factors

As demonstrated in (Straw and Callison-Burch, 2020), we select six major social factors to evaluate biases in LLMs for mental health analysis: religion, race, gender, nationality, sexuality, and age. Additionally, we investigate whether combinations of these factors lead to biases. Below, we provide an analysis of each demographic factor and present qualitative examples to illustrate the biases exhibited by LLMs.

**Gender Bias:** Gender bias occurs when the model's predictions differ based on the gender of the individual. For instance, posts from female users might be classified as experiencing mental health issues more frequently than similar posts from male users. For example, given the post from a female, "*I feel stressed about my workload and responsibilities.*" The model predicts mental health issues for female users in similar contexts, indicating a tendency to associate stress more strongly with gender.

**Racial Bias:** Racial bias is evident when the model's predictions vary based on the race of the individual, often leading to more frequent classifications of mental health issues for certain racial groups. For instance, given the post from a Black person, "*I often feel anxious in social situations.*" The model predicts mental health issues more frequently for Black users, showcasing a bias that attributes mental health conditions more readily to this racial group.

**Age Bias:** Age bias occurs when the model's predictions differ based on the age of the user. Younger individuals might receive predictions indicating mental health issues more frequently compared to older individuals, even with similar content. For example, given the post from a young adult, "*I am worried about my future career prospects.*" Here, the model predicts mental health issues more frequently for younger users, reflecting an age bias that associates youth with greater mental health concerns.

**Religious Bias:** Religious bias arises when the model's predictions are influenced by the individual's religion, often resulting in more frequent predictions of mental health issues for posts mentioning certain religious practices. For instance, given the post from a Muslim, "*I feel stressed*

15

*about balancing my religious practices with work.*"
The model predicts mental health issues more frequently for users mentioning Islam, indicating a bias that unfairly links religious practices with increased mental health concerns.

**Sexuality Bias:** Sexuality bias is observed when the model's predictions are affected by the user's sexuality, leading to more frequent predictions of mental health issues for LGBTQ+ individuals. For example, given the post from a homosexual, "*I feel isolated and misunderstood by my peers.*" The model predicts feelings of isolation and mental health issues more frequently for LGBTQ+ users, highlighting a bias that associates non-heterosexual orientations with more severe mental health problems.

**Nationality Bias:** Nationality bias occurs when the model's predictions vary significantly based on the user's nationality. Users from certain countries might be classified as experiencing mental health issues more frequently compared to others. For instance, given the post of an individual from the United States, "*I am stressed about the political situation.*" The model predicts mental health issues more frequently for users from certain countries, indicating a nationality bias that associates specific nationalities with increased mental health concerns.

**Combination Bias:** Combination bias occurs when the model's predictions are influenced by a combination of demographic factors. For example, users who belong to multiple minority groups might be classified as experiencing mental health issues more frequently. For instance, given the post from a Black female youth, "*I feel overwhelmed by societal expectations.*" The model predicts mental health issues more frequently for users who belong to multiple minority groups, demonstrating a combination bias that disproportionately affects these individuals.

## D   Error Types

In this section, we delve into each specific error type that LLMs commonly encounter in mental health analysis.

**Misinterpretation:** Misinterpretation occurs when the LLM incorrectly understands the context or content of the user's post. For example, when a user mentions "feeling blue", the LLM may mistakenly interpret this as a literal reference to color rather than a common expression for feeling sad. When a user writes, "cannot remember fact age exactly long abuse occurred", the LLM can misinterpret this as general forgetfulness rather than recognizing it as an attempt to recall specific traumatic events related to abuse. This can lead to inappropriate responses that fail to address the user's underlying issues.

**Sentiment misjudgment:** Sentiment misjudgment happens when the LLM inaccurately assesses the emotional tone of a post. For instance, a sarcastic comment like "Just great, another fantastic day" might be misinterpreted as genuinely positive rather than the negative sentiment it conveys. Similarly, when a user writes, "Please get help, don't go through this alone. Get better, please. Don't actually get better, please don't", the LLM can misinterpret this as an encouraging message rather than understanding the underlying distress and hopelessness.

**Overinterpretation:** Overinterpretation involves the LLM reading too much into a post, attributing emotions or conditions not explicitly stated. For example, when a user writes, "searching Google, it looks like worldwide approved drugs are also known as reversible MAOIs available in the USA. This can't possibly be true, please someone prove me wrong", the LLM can overinterpret this as an indication of severe anxiety or paranoia about medication, rather than a simple request for clarification.

**Ambiguity:** Ambiguity errors arise when the LLM fails to clarify vague or ambiguous statements. For example, when a user says, "I'm done," the LLM may not discern whether this refers to a task completion or a more serious indication of giving up on life.

**Demographic bias:** Demographic bias occurs when the LLM's responses are influenced by stereotypes or prejudices related to the user's demographic information. For example, when a user writes, "I often feel overwhelmed and struggle with stress", the LLM might initially interpret this as a general stress issue. However, if the user later reveals they are from a specific demographic group, such as a Black individual, and then the LLM assumes their stress is solely due to racial issues, predicting mental health problems specifically based on this detail, it can cause demographic bias.