
AUTOMATED MULTI-LABEL ANNOTATION FOR MENTAL HEALTH ILLNESSES USING LARGE LANGUAGE MODELS

Abdelrahaman A. Hassan[†], Radwa J. Hanafy^{†,‡} and Mohammed E. Fouda[†]

[†]Compumacy for Artificial Intelligence Solutions, Cairo, Egypt.

[‡] Department of Behavioural Health- Saint Elizabeths Hospital, Washington DC, 20032.

fouda@compumacy.com

ABSTRACT

The growing prevalence and complexity of mental health disorders present significant challenges for accurate diagnosis and treatment, particularly in understanding the interplay between co-occurring conditions. Mental health disorders, such as depression and Anxiety, often co-occur, yet current datasets derived from social media posts typically focus on single-disorder labels, limiting their utility in comprehensive diagnostic analyses. This paper addresses this critical gap by proposing a novel methodology for cleaning, sampling, labeling, and combining data to create versatile multi-label datasets. Our approach introduces a synthetic labeling technique to transform single-label datasets into multi-label annotations, capturing the complexity of overlapping mental health conditions. To achieve this, two single-label datasets are first merged into a foundational multi-label dataset, enabling realistic analyses of co-occurring diagnoses. We then design and evaluate various prompting strategies for large language models (LLMs), ranging from single-label predictions to unrestricted prompts capable of detecting any present disorders. After rigorously assessing multiple LLMs and prompt configurations, the optimal combinations are identified and applied to label six additional single-disorder datasets from RMHD. The result is SPAADE-DR, a robust, multi-label dataset encompassing diverse mental health conditions. This research demonstrates the transformative potential of LLM-driven synthetic labeling in advancing mental health diagnostics from social media data, paving the way for more nuanced, data-driven insights into mental health care.

Keywords Large Language Models (LLMs), Mental Health Assessment, Stress, Anxiety, Depression, PTSD, ADHD, Eating Disorder, Suicide, Zero-shot Learning, Data Annotation, Multi-label classification, Mental Disorders Comorbidity

1 Introduction

Mental illnesses are a major global health issue, affecting millions of people across all age groups and societies. Mental disorders not only impact the well-being and quality of life of individuals but also contribute to a significant burden on healthcare systems worldwide. Mental health conditions encompass a broad spectrum, including disorders like depression, anxiety, and PTSD, as well as more severe outcomes like suicide [1]. For instance, depression affects approximately 280 million people globally, accounting for 3.8% of the population. Among adults, 5% are affected, with a higher prevalence in women (6%) compared to men (4%)[2]. Anxiety disorders, the most common of all mental health conditions, impact around 301 million people globally, with 4% of the world's population experiencing these disorders[3]. Moreover, the severe consequences of untreated mental illness are starkly evident in the global suicide rate, where an estimated 726,000 people take their own lives each year, making it one of the leading causes of death, especially among young people aged 15 to 29[4]. These statistics illustrate the widespread nature and severity of mental illnesses, underscoring the urgent need for effective methods of early diagnosis and intervention. In recent years, social media platforms have emerged as valuable resources for mental health research, providing a rich source of data that can shed light on the prevalence, nature, and impact of these disorders. Researchers have utilized social media data to assess mental health conditions, identify individuals at risk, and develop interventions to improve well-being. The

availability of massive datasets, coupled with advanced computational methods, has fueled the development of novel approaches for mental health diagnosis using social media.

Mental illnesses, while traditionally diagnosed through clinical interviews and standardized assessments, are increasingly being detected through social media. With the widespread use of platforms such as X (formerly Twitter), Reddit, and Facebook, individuals often express their mental health struggles in online forums. This fact offers an unprecedented opportunity for researchers and clinicians to tap into large-scale, real-time data to identify patterns indicative of mental health conditions. Studies have shown that certain linguistic markers, such as negative sentiment, changes in language use, and self-referential statements, can be correlated with mental health issues like depression and anxiety[5, 6]. For example, research has demonstrated that individuals with depression often use more first-person pronouns, indicating a heightened focus on the self, along with an increased use of negative emotion words such as "sad" or "hopeless". Similarly, anxiety-related posts frequently contain words reflecting uncertainty and worry, alongside expressions of physical symptoms like insomnia or fatigue. By mining these linguistic patterns, social media have become a valuable resource for detecting early signs of mental health deterioration[7]. However, there are also significant challenges in utilizing social media for mental illness diagnosis. Issues of privacy, the need for ethical guidelines, and the potential for misinterpretation of self-expressed symptoms are ongoing concerns[8]. Despite these limitations, the potential for social media to provide real-time, large-scale insights into population-level mental health trends is a growing area of research, offering a complementary tool to traditional clinical diagnosis.

Recent advancements in natural language processing (NLP) and the development of large language models (LLMs) have significantly enhanced the potential for automated mental illness diagnosis. LLMs such as GPT, LLAMA and Phi, and their derivatives are capable of processing large volumes of unstructured text from social media and other platforms, allowing them to identify linguistic patterns associated with mental health conditions. These models can capture nuanced language use, which traditional machine learning methods may overlook. For instance, LLMs have been used to detect signs of depression, anxiety, and even suicidal ideation through the analysis of social media posts. By leveraging pre-trained models, researchers have demonstrated that these models can outperform previous methods, such as lexicon-based approaches like LIWC and feature extraction techniques like TF-IDF, which rely on manual feature engineering and are limited in capturing nuanced linguistic patterns [9, 10, 11]. These advancements highlight the ability of LLMs to capture subtle mental health signals, such as shifts in tone, sentiment, and specific word usage [12], providing a significant improvement over earlier approaches.

Furthermore, LLMs can be fine-tuned to adapt to the specific context of mental health data, improving their accuracy and making them highly versatile in handling diverse datasets. One of the key advantages of LLMs in this domain is their ability to perform multi-label classification, where a single post can be labeled with multiple mental health conditions. This capability is crucial for detecting co-occurring mental illnesses, such as depression and anxiety, which are often diagnosed together in clinical settings [13]. Additionally, LLMs have the unique capability to not only provide classifications but also offer explanations for their decisions, making them particularly valuable in understanding the rationale behind a diagnosis [14, 15]. In the context of mental health, this can help clinicians better interpret why the model labeled a post with certain conditions based on linguistic features or patterns identified in the text. By providing these explanations, LLMs help bridge the gap between automated diagnosis and human understanding, ensuring that the outputs are interpretable and clinically relevant [16]. Moreover, the scalability of LLMs enables large-scale deployment across various social media platforms, potentially providing real-time insights into mental health trends on a population level [17].

However, there are still challenges to overcome when applying LLMs to mental illness diagnosis. Ensuring model accuracy across diverse populations and linguistic variations, as well as addressing ethical concerns around privacy and the responsible use of data, remain important areas for ongoing research [18]. Despite these hurdles, LLMs represent a significant step forward in the automation and scalability of mental illness diagnosis, offering a promising tool for improving early detection and intervention efforts. One of the key challenges in advancing automated mental illness diagnosis using large language models (LLMs) lies in the lack of publicly available multi-label mental illness datasets. Current datasets typically focus on single-label classification, where each instance (e.g., a social media post) is annotated for only one mental health condition, such as depression or anxiety. However, in real-world clinical settings, mental illnesses often co-occur; individuals may experience multiple conditions simultaneously, such as depression with anxiety or other mental disorders. The absence of multi-label datasets limits the ability of models to learn associations between these co-occurring conditions. Multi-label classification, where a post could be tagged with more than one condition, would enable more accurate and realistic diagnostic models. By identifying and addressing these associations, models could better reflect the complexity of mental health conditions and offer more nuanced insights for early intervention.

In the current research landscape, several studies explore LLMs for data annotation [19] and synthetic labeling, predominantly focusing on enhancing data diversity and improving model performance for single-label classification tasks. However, these approaches do not address the need to expand single-label datasets into multi-label formats

through zero-shot synthetic labeling—a technique that would allow LLMs to predict multiple potential labels for each instance without prior multi-label training [20, 21]. This gap is particularly evident in mental health research, where multi-label datasets could drastically improve the diagnostic accuracy and effectiveness of models. Despite the proven benefits of multi-label classification in other domains [22], mental health research has yet to see the same level of dataset development. This gap presents a significant barrier to progress, as the lack of comprehensive, multi-label datasets prevents models from fully capturing the intricacies of mental health symptoms. Without these datasets, LLMs cannot leverage the correlations between various mental illnesses, limiting their diagnostic accuracy and effectiveness.

This research focuses on leveraging the capabilities of large language models (LLMs) to address the critical gap in multi-label mental illness datasets. LLMs play a pivotal role in automating data annotation, offering a scalable and efficient solution for creating annotated datasets without requiring extensive manual effort. By using LLMs, more complex and nuanced datasets can be generated to capture the co-occurrence of mental health conditions. This automation allows for the rapid generation of labeled data, providing a valuable resource for future research and applications in mental health diagnosis. The primary contribution of this work is the development of a method to synthetically transform any single-label dataset into a multi-label dataset using LLMs in a zero-shot setting. Unlike existing approaches that depend on labeled examples for each condition, this method enables LLMs to infer potential co-occurring conditions within a given text. This transformation converts single-label datasets into representations that capture a more realistic and multi-dimensional view of mental health. Furthermore, this approach is applied to create a new, comprehensive multi-label mental illnesses dataset labeled for six distinct mental disorders. This dataset reflects co-occurring conditions and provides a valuable resource for training and evaluating multi-label classification models in mental health. This innovation not only simplifies the process of dataset annotation but also enhances data quality, making it better suited for real-world clinical applications. Ultimately, this approach improves the accuracy and relevance of predictive models for mental health diagnosis, paving the way for more advanced and comprehensive AI-driven mental health assessments. To the best of our knowledge, this is the first work to address this gap in the dataset and annotation.

The rest of this paper is organized as follows: Section 2 provides a review of relevant datasets, including Dreddit, DepSeverity, and the Reddit Mental Health Dataset (RMHD), as well as the large language models (LLMs) and evaluation metrics used. Section 3 outlines the DepSeverity-Dreddit multi-label dataset, describes the different prompt template types (single-label, multi-label, and unrestricted prompts), and details the evaluation process on multi-label datasets, including the SPAADE-DR dataset’s preparation. Section 4 presents the findings from evaluations on the DepSeverity-Dreddit dataset and the SPAADE-DR dataset, covering aspects such as data cleaning, sampling, and multi-label labeling. Section 5 explores the comorbidity between disorders and evaluates prompts and LLMs based on six key disorders. Finally, the paper concludes with a discussion in Section 6, highlighting key takeaways and suggesting areas for future research.

2 Resources and Methods

In this section, the key resources and methodologies underlying the study are presented, including the datasets, large language models (LLMs), and evaluation metrics used. By integrating multiple datasets and leveraging advanced LLMs, the goal is to develop a framework capable of capturing complex, real-world mental health patterns. The evaluation metrics are designed to support this goal, offering a robust framework to assess model performance in both single- and multi-label classifications.

2.1 Datasets

This paper experiments with seven primary single-label datasets, each focusing on a specific mental health condition, as well as the RMHD dataset, which contains data on multiple mental health conditions. These datasets are individually labeled for conditions like depression, anxiety, ADHD, eating disorders, PTSD, and emergent conditions, such as suicide. Below is an overview of each dataset used in this research.

2.1.1 Dreddit

The Dreddit dataset [23] contains 190,000 Reddit posts across ten subreddits in the five domains of abuse, social, anxiety, PTSD, and financial. These posts span a diverse range of content, including personal narratives, advice-seeking posts, and emotional expressions. Of these, 3,553 post segments are manually labeled with stress indicators by human annotators using Amazon Mechanical Turk, creating a supervised training set for stress detection. With its large scale, rich variety of content, and detailed annotations, Dreddit is an invaluable resource for researchers studying stress in social media, providing insights across multiple real-world domains.

2.1.2 DepSeverity

The DepSeverity dataset [24] is derived from the same 3,553 Reddit posts used in the Dreaddit dataset, focusing on depression severity. These posts are labeled with four clinical levels of depression severity: Minimal, Mild, Moderate, and Severe, based on the Depressive Disorder Annotation (DDA) scheme[25]. The labeling process utilized six clinical resources: the Diagnostic and Statistical Manual of Mental Disorders, 5th Edition (DSM-5) [1], the Behavioral Risk Factors Surveillance System (BRFSS), the Harvard Department of Psychiatry National Depression Screening Day Scale (HANDS), the Patient Health Questionnaire (PHQ-9)[26, 27], the Quick Inventory of Depressive Symptomatology (QIDS-SR), and the Columbia Suicide Severity Rating Scale(C-SSRS) [28].

Originally created as a binary dataset for detecting the presence of depression, it is later expanded to include four severity levels, providing a more nuanced understanding of users’ mental health. This transformation helps better assess users’ condition, supporting potential intervention and treatment strategies. As noted by the dataset authors, at least 10 posts from a user are necessary for reliably predicting their depression severity, emphasizing the importance of longitudinal data for early intervention. With its detailed annotation and substantial volume of posts, DepSeverity offers a rich resource for researchers studying depression in social media. It spans various types of content such as personal narratives, advice-seeking, and emotional expressions, making it ideal for developing models that detect both the presence and severity of depression.

2.1.3 Reddit Mental Health Dataset

The Reddit Mental Health Dataset (RMHD) [29] consists of a large collection of Reddit posts extracted from 28 subreddits, including both mental health-focused communities and general interest groups. The dataset spans posts from 17 mental health subreddits such as r/depression, r/Anxiety, and r/SuicideWatch, as well as 11 non-mental health subreddits such as r/conspiracy, r/legaladvice, r/personalfinance.

The dataset includes posts from unique users across two critical timeframes: prepandemic (November 2018 to November 2019) and mid-pandemic (January to April 2020), allowing for comparative analyses of mental health trends before and during the COVID-19 pandemic. All posts are preprocessed to include only English-language content, and posts from bots, advertisements, and duplicate users are removed to ensure data quality. The subreddit each post belongs to is used as a label, associating the post with specific mental health conditions or general topics discussed in the community.

In addition to the raw text data, a variety of features are extracted from each post to enhance its utility for research. These include sentiment analysis using VADER, lexical counts, and the application of the Linguistic Inquiry and Word Count (LIWC) tool to assess semantic and grammatical categories such as emotions, pronouns, and body references. Moreover, lexicons are manually developed to track topics like suicidality, economic stress, isolation, and substance use, offering insights into the mental health conditions discussed in these posts. With its wide range of posts and comprehensive feature extraction, RMHD serves as a valuable resource for analyzing mental health trends. It is particularly suited for studies on mental health conditions, topic modeling, and trend analysis across large online communities.

2.2 LLMs

In this paper, experiments are done using five models of varying sizes, architectures, and model families to ensure comprehensive coverage across different approaches. These models are selected for their diversity in source and design, allowing for a robust comparison in performance.

GPT-4o-mini: Developed by OpenAI, GPT-4o-mini is a scaled-down variant of GPT-4o, designed to balance performance with computational efficiency. It features fewer parameters than the full GPT-4o model, making it more accessible for smaller-scale tasks without sacrificing core capabilities in natural language understanding and generation. The model is optimized for various NLP applications, such as text classification, summarization, and dialogue systems.

Llama-3 70b: Developed by Meta, Llama-3 70b is a 70-billion-parameter model designed to excel in a range of natural language processing tasks. It features multilingual support, coding capabilities, and advanced reasoning. Llama-3 benefits from an improved data curation pipeline and optimized scaling laws, allowing it to match the performance of leading models like GPT-4 across various benchmarks. Despite its size, Llama-3 has been fine-tuned to balance computational efficiency with high performance in both general and domain-specific tasks.

Mistral NeMo 12b: Mistral NeMo is a state-of-the-art model with 12 billion parameters, developed in collaboration with NVIDIA. It features an expansive 128k context window, offering exceptional performance in reasoning, world knowledge, and code generation within its size category. Mistral NeMo is particularly strong in multilingual applications, supporting a wide range of languages, including English, French, German, Chinese, and Arabic. The model uses the highly efficient Tekken tokenizer, which significantly improves compression rates across various languages and source code, outperforming previous Mistral models.

Phi-3.5-MoE: The Phi-3.5-MoE model, developed by Microsoft, uses a Mixture of Experts (MoE) architecture and features a total of 42 billion parameters, with 6.6 billion active parameters at any given time during inference. This model is designed to activate only a subset of its parameters (two out of sixteen experts) depending on the complexity of the task, which allows it to optimize resource usage while maintaining high performance across various natural language processing (NLP) tasks. This dynamic routing of experts enables the model to excel in diverse tasks such as reasoning, multilingual support, and code generation, while still being computationally efficient. By testing Phi-3 Medium, it is obvious that Phi-3.5-MoE performed better across our benchmarks, leading us to choose it over the Phi-3 Medium model.

Gemma-2 9b: Developed by Google DeepMind, Gemma-2 9b is a 9-billion-parameter model designed to deliver state-of-the-art performance in natural language understanding tasks. This model benefits from several enhancements, such as knowledge distillation and advanced attention mechanisms like Grouped-Query Attention (GQA). Despite its relatively smaller size compared to larger models, Gemma-2 9b demonstrates competitive performance across a range of benchmarks, offering a practical balance between size and accuracy for real-world applications.

Table 1: Summary of large language models used in the experiments

Model	Size (Parameters)	Source	API Service Provider
Gemma 2 9b [30]	9B	Google	Groq
GPT-4o mini [31]	Proprietary	OpenAI	OpenAI
Llama 3 70b [32]	70B	Meta	Groq
Mistral NeMo [33]	7B	Mistral AI	Mistral AI
Phi-3.5-MoE [34]	42B	Phi-1 Labs	Azure AI

2.3 Evaluation Metrics

The evaluation of LLM performance is conducted using a variety of metrics, grouped into three categories: per-class, overall (multi-label), and multi-class. Metrics applied for per-class and multi-class evaluations include:

- **Balanced Accuracy (BA):**

$$\mathbf{BA} = \frac{1}{N} \sum_{i=0}^{N-1} \text{Recall}_i \quad (1)$$

where:

- N : The total number of classes in the classification problem.
- Recall_i : The recall for the i^{th} label or class.

- **F1-Score (F1):** (Binary F1-score)

$$F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

- **Precision (CP):**

$$\text{CP} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

- **Recall (CR):**

$$\text{CR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

where TP donates True Positives, FP False Positives, TN True Negatives, and FN False Negatives.

Balanced accuracy is chosen over regular accuracy to address the imbalance present in most datasets. The F1-Score is incorporated to evaluate the model’s performance, as it balances precision and recall, making it especially valuable when both false positives and false negatives have significant implications. For multi-label evaluation, the Micro F1-Score is employed to ensure equal consideration of all classes in the assessment. Overall metrics, including Balanced Accuracy, Precision, and Recall, are calculated using the combined values of True Positives, False Positives, True Negatives, and False Negatives across all classes. These overall values are computed as follows:

- **Overall TP, FP, TN, and FN:**

$$\text{Overall TP} = \sum_{i=1}^n \text{TP}_i$$

$$\text{Overall FP} = \sum_{i=1}^n \text{FP}_i$$

$$\text{Overall TN} = \sum_{i=1}^n \text{TN}_i$$

$$\text{Overall FN} = \sum_{i=1}^n \text{FN}_i$$

where n is the number of classes or instances.

Hamming Loss is also utilized to evaluate performance in the multi-label classification setting. This metric quantifies the fraction of labels incorrectly predicted across all instances and labels, making it especially useful for multi-label tasks where each instance can have multiple associated labels. By focusing on individual label errors, Hamming Loss provides a granular view of the model’s accuracy across all labels, ensuring that each label’s misclassification is accounted for. It is calculated as follows:

$$\text{Hamming Loss} = \frac{1}{N \cdot L} \sum_{i=1}^N \sum_{j=1}^L \mathbb{I}(y_{ij} \neq \hat{y}_{ij})$$

where N is the total number of instances, L is the total number of labels, y_{ij} is the true label for instance i and label j , \hat{y}_{ij} is the predicted label for instance i and label j , $\mathbb{I}(\cdot)$ is the indicator function, which is 1 if $y_{ij} \neq \hat{y}_{ij}$ and 0 otherwise.

A lower Hamming Loss indicates better model performance, with 0 representing perfect label prediction for every instance.

The Odds Ratio (OR), a measure of association between two events, is employed to analyze associations and comorbidities among mental disorders. Commonly used in medical and statistical analyses, the OR helps assess how strongly the presence or absence of one condition is associated with another. For two binary variables. This allows us to assess how likely one disorder is to occur alongside another, offering insights into potential comorbidities, the OR is calculated as follows:

$$\text{OR} = \frac{\frac{a}{b}}{\frac{c}{d}} = \frac{a \times d}{b \times c}$$

Where **a** = Number of cases where both the exposure and outcome are present (both event A and B occur). **b** = Number of cases where the exposure is present but the outcome is absent (event A occurs, but not event B). **c** = Number of cases where the exposure is absent but the outcome is present (event A does not occur, but event B occurs). **d** = Number of cases where neither the exposure nor the outcome is present (neither event A nor B occur).

3 Methodology

The methodology aims at developing a robust, multi-label dataset for mental health diagnostics by synthetically labeling social media posts, leveraging various prompt strategies with LLMs. This section outlines the structured workflow for dataset creation and labeling, as shown in Figure 1. The Depseverity-Dreaddit merged dataset, labeled for depression and stress, is used to test three distinct prompt strategies—single-label, multi-label, and unrestricted—to guide LLMs in annotating multiple disorders. After evaluating the prompt effectiveness across LLMs, the optimal prompt-LLM combination is used to synthetically label RMHD dataset, resulting in a multi-label dataset that reflects complex mental health profiles.

3.1 Depseverity-Dreaddit as Multi-Label Dataset

The Depseverity and Dreaddit datasets consist of identical posts from the same user base, each independently labeled for either depression or stress. To utilize this overlap, these datasets are merged into a unified multi-label dataset, providing a resource labeled for both depression and stress. The resulting dataset distribution is shown in Table [2]. This new dataset, annotated by psychiatrists, enables a more nuanced analysis of mental health conditions, where users may exhibit symptoms of both disorders simultaneously. Merging these datasets enriches the available data for training and evaluating models capable of multi-label classification, thereby improving the accuracy and applicability of mental health detection tools.

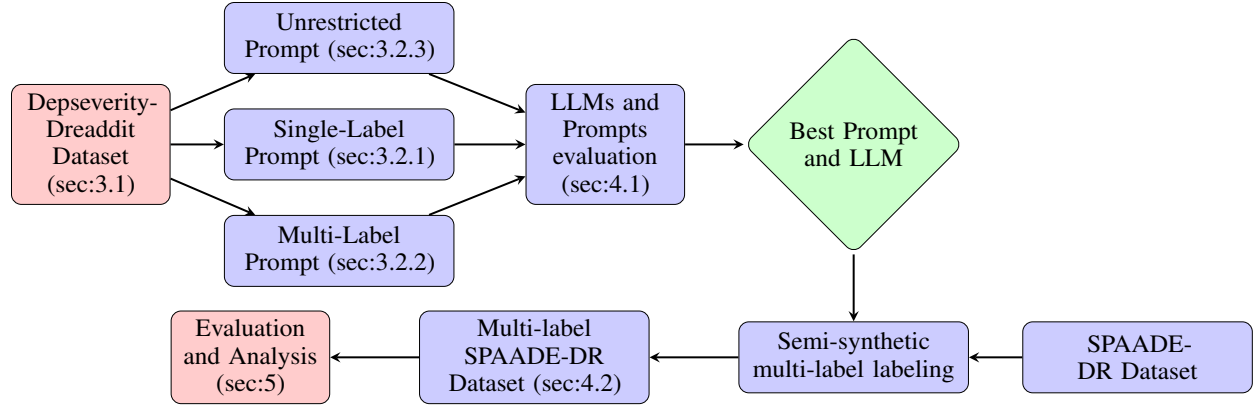


Figure 1: Process Workflow

Table 2: The distribution of the depseverity-dreaddit dataset

Disorder	Stress		
	Label	Negative	Positive
Depression	Negative	1532	1041
	Positive	154	814

3.2 Prompt Template

The experiments utilize various types of prompts to diagnose mental illness in social media posts, with a particular focus on multi-label classification. Single-label, multi-label, and unrestricted prompts are tested across various datasets to evaluate their performance and efficiency. Each prompt template is paired with a specific parser to convert the model’s responses (e.g., "yes," "no," or the name of a disorder) into binary labels.

3.2.1 Single-Label Prompts

Single-label prompts are designed to diagnose one mental illness at a time, focusing the model’s attention on identifying whether the user exhibits symptoms of a specific disorder. This approach is simple yet effective for situations where only one condition is being evaluated. The model is prompted to make a binary decision, determining if the user shows clear signs of the targeted mental illness. A prompt inspired by those in [35] is utilized.

The following prompt template is used to guide the model in identifying a specific mental health condition:

```

[Task]
Analyze the following social media post to determine if the writer exhibits clear
symptoms of {The target disorder} according to provided guidelines.

[Guidelines]
- Concise Response: Respond only with 'Yes' (exhibits clear symptoms of {The target
disorder}) or 'No' (Does not exhibit clear symptoms of {The target disorder}).
- No Explanations: Don't provide explanations for your assessment.
- Ambiguity: If the post is unclear, choose the most probable label.

[Post]
{The Post}
  
```

Figure 2: Single-label Binary Prompt Template for Identifying Mental Health Conditions

This prompt template in Fig:2 allows for focused diagnosis of a single condition, minimizing confusion and maintaining clarity. The model’s binary response (Yes/No) delivers a straightforward mechanism for assessing the user’s symptoms

according to specific clinical guidelines. When adapted for multi-label problems, this approach becomes resource-intensive, as it requires applying the same prompt multiple times—once for each mental illness or label being evaluated. While this method allows the model to focus exclusively on diagnosing a single illness per prompt, thereby reducing confusion, it incurs significantly higher computational costs.

3.2.2 Multi-Label Prompts

Various multi-label prompt templates are developed and tested to evaluate whether users exhibit multiple mental illnesses simultaneously. For this purpose, two primary approaches are identified:

Template 1: Multi-Class Classification Approach

This template frames the problem as a multi-class classification task. In cases with two mental illnesses (e.g., depression and stress), the model categorizes the post into one of four classes: "Depressed," "Stressed," "Depressed and Stressed," or "Normal". This approach accounts for scenarios where users display symptoms of more than one condition. However, as the number of mental illnesses (n) increases, the number of possible classes grows exponentially (2^n), which complicates classification and can confuse the model. As a result, this approach becomes less practical for higher numbers of labels due to the explosion in the number of classes.

Template 2: Multi-Label Classification Approach

This template treats the problem as a true multi-label classification task, asking the model to determine whether any of the $n+1$ conditions (including "Normal") is present. In the case of depression and stress, the model is prompted to assess whether the poster is "Depressed", "Stressed", allowing for any combination of these labels or else he is "Normal". This method simplifies the classification task by avoiding the exponential class growth seen in Template 1. Instead of handling a growing number of classes, the model only deals with a manageable set of label combinations, making it more efficient as the number of conditions increases.

The multi-label approach ensures greater flexibility in identifying overlapping symptoms and prevents the complexity associated with multi-class classification. Both templates are evaluated for their effectiveness in addressing posts that exhibit co-occurring mental health symptoms. The prompt templates in Fig: 3 guides the model in identifying a number of specific mental health conditions, in this case depression and stress. These templates are designed to provide clear and concise assessments while maintaining the flexibility to handle multiple conditions, making them effective tools for multi-label classification tasks.

3.2.3 Unrestricted Prompt

Unrestricted prompts are also explored, enabling the model to diagnose multiple mental illnesses without being restricted to predefined categories. The aim of this approach is to give the model more freedom in detecting any mental health condition based solely on the content of the post, rather than being confined to a limited set of labels. The unrestricted prompt is designed to assess the model's ability to identify a broader spectrum of mental health issues and are compared to the more structured single-label and multi-label prompts. The flexibility of this approach enables the model to diagnose multiple co-occurring conditions or determine if the poster exhibits no symptoms at all.

As a psychiatrist, analyze the provided social media post to determine if the writer exhibits clear symptoms of any mental illness according to provided guidelines.

[Guidelines]

- The poster could have multiple illnesses at the same time; otherwise, he is normal.
- Without explaining your reasoning, respond with any combination of mental illnesses names separated by comas only and if the poster doesn't have any mental illness just answer with "Normal"
- Answer with basic mental illnesses names only like "Depression" without further details or complex names.

[Post]

{The post}

Figure 4: Unrestricted Binary Prompt Template for Identifying Mental Health Conditions

[Task]
Analyze the following social media post to determine if the writer exhibits clear symptoms of {Depression or Stress} according to provided guidelines.

[Guidelines]
- The poster could have multiple illnesses at the same time; otherwise, he is normal.
- Concise Response: Respond with one of these 4 words only ["Depressed", "Stressed", "Depressed and Stressed", "Normal"].
- No Explanations: Don't provide explanations for your assessment.
- Ambiguity: If the post is unclear, choose the most probable label.

[Post]
{The Post}

Multi Label Binary Template 1

[Task]
As a psychiatrist, analyze the provided social media post to determine if the writer exhibits clear symptoms of {Depression or Stress} according to provided guidelines.

[Guidelines]
- The poster could have multiple illnesses at the same time; otherwise, he is normal.
- Concise Response: Respond with any combination of these 2 mental illness names or "Normal" only ["Depressed", "Stressed"]. If the poster doesn't have any mental illness, just answer with "Normal".
- No Explanations: Don't provide explanations for your assessment.
- Ambiguity: If the post is unclear, choose the most probable label.

[Post]
{The post}

Multi Label Binary Template 2

Figure 3: Multi-Label Binary Prompt Templates for Identifying Mental Health Conditions

This template in Fig:4 provides the model with the flexibility to identify a wide range of mental health conditions without the constraints of predefined labels, making it a valuable tool for exploring diverse mental health issues in user-generated content.

3.3 Evaluation on a Multi-label dataset

The performance of various prompt templates and models is evaluated using the Depseverity-Dreaddit dataset. The goal is to identify the most effective prompt-template and model combinations for accurate multi-label classification. Additionally, the impact of majority voting across model predictions is examined to determine its potential for improving accuracy and reducing biases in the results. This analysis provides insights into how different prompts and models handle overlapping mental health conditions in a multi-label context.

3.4 SPAADE-DR dataset

An extensive search reveals a lack of existing, accurate multi-label datasets for mental health classification. To fill this gap, a new semi-synthetic multi-label dataset, SPAADE-DR (Suicidal Ideation, PTSD, Anxiety, ADHD, Depression, Eating Disorder Diagnosis from Reddit posts), is created by combining and labeling multiple single-label datasets, each dedicated to a specific mental health condition. For this task, the RMHD dataset is utilized, encompassing data for various mental disorders. Specifically, datasets for conditions such as ADHD, anxiety, depression, eating disorders, PTSD, and suicide are incorporated. By merging these single-label datasets, a comprehensive multi-class dataset is constructed to represent multiple mental health conditions. This allows for more sophisticated multi-label labeling,

analysis, and model evaluation, providing a valuable resource for further research into mental health classification. See section:4.2 for more details.

Following the evaluation of the prompt templates and model performance, the optimal combinations are applied to label the SPAADE-DR dataset. the best-performing prompt-template and model combinations identified from the Depseverity-Dreaddit evaluation are utilized to predict and label the remaining five conditions for each post, ensuring consistency and accuracy in the multi-label classification. After labeling the SPAADE-DR dataset, it is used to evaluate the performance of both multi-label and unrestricted prompt templates on the newly created six-label dataset. This evaluation aims to assess how well the models perform when the number of labels increases in a multi-class setting. Additionally, it allows us to explore the models’ ability to identify correlations between various mental illnesses and their associated symptoms. Testing both multi-label and unrestricted prompts aims to uncover the strengths and limitations of the models in capturing the complexity of overlapping mental health conditions.

4 Results & Discussion

This section presents the results of the synthetic labeling methodology, highlighting the effectiveness of various prompt strategies and their influence on multi-label mental health diagnostics. Each prompt-LLM combination is evaluated on the Depseverity-Dreaddit dataset to analyze how single-label, multi-label, and unrestricted prompts affect diagnostic accuracy across various disorders. Insights gained from this analysis inform the selection of the optimal prompt-LLM configuration for creating a nuanced multi-label dataset. Finally, the implications of the results are discussed, emphasizing the dataset’s potential to advance multi-disorder mental health assessment in social media contexts.

4.1 Evaluation on Depseverity-Dreaddit Dataset

In the initial experiment, various LLMs and prompt templates are evaluated using a balanced subset of the Depseverity-Dreaddit dataset. The goal of this evaluation is to identify the most effective LLM-prompt combinations for accurate multi-label classification of mental health conditions. Additionally, majority voting across models is tested to assess its impact on improving overall accuracy.

The metrics chosen are used to assess the models’ ability to predict each label (depression and stress) separately, their overall multi-label classification accuracy, and their performance in multi-class classification (combining all possible classes). From the results in Table 3, it is clear that the best-performing LLM is Llama-3 70b, closely followed by GPT-4o-mini and Phi-3.5 MoE. Llama-3 70b demonstrates the highest scores using all prompts across most metrics, particularly in multi-label classification, where it achieves the highest overall balanced accuracy (0.78) with a lower hamming loss (0.24). GPT-4o-mini and Phi-3.5 MoE also perform strongly, with competitive results in the same metrics, making them viable alternatives.

In terms of prompt templates, the single-label prompt template consistently outperforms the multi-label templates. This is evident from the higher precision, recall, and F1-scores when evaluating the LLMs on each label (depression and stress). The single-label prompts provide a more focused diagnosis for each condition, contributing to the models’ higher accuracy in both multi-label and multi-class evaluations. Although multi-label templates capture co-occurring conditions, they generally result in lower precision and higher Hamming loss, which indicates some confusion in classifying posts with overlapping symptoms.

Additionally, after testing majority voting between models to see if combining predictions from multiple LLMs improves accuracy, it is obvious that while majority voting slightly improves recall, it does not significantly impact precision or F1-scores, and in some cases, it leads to higher Hamming loss. As a result, majority voting does not offer a clear advantage over individual model predictions. In conclusion, for multi-label classification on the Depseverity-Dreaddit dataset, the combination of Llama-3 70b and the single-label prompt template delivers the most accurate and reliable results. This combination effectively handles the complexities of diagnosing both depression and stress simultaneously.

Table 3: Comparisons between prompt templates and LLMs on the Depseverity-Dreaddit dataset are conducted using several metrics: per-category balanced accuracy (CBA), F1-measure (CF1), precision (CP), recall (CR), overall balanced accuracy (OBA), overall precision (OP), and overall recall (OR)[36, 37], hamming loss (HL)[38], and Multi-class balanced accuracy (BA).

Prompt	LLM	Depression				Stress				Multi-Label					Multi-Class BA
		CBA	CF1	CP	CR	CBA	CF1	CP	CR	GBA	OF1	OP	OR	HL	
Single-Label	Llama-3 70b	0.74	0.61	0.54	0.71	0.75	0.81	0.70	0.96	0.78	0.74	0.64	0.88	0.24	0.46
	Gemma-2 9b	0.72	0.58	0.44	0.84	0.64	0.75	0.60	1.00	0.71	0.69	0.54	0.94	0.34	0.39
	Phi-3.5-MoE	0.75	0.62	0.52	0.75	0.74	0.79	0.70	0.92	0.77	0.73	0.63	0.87	0.25	0.46
	GPT-4o-mini	0.74	0.61	0.47	0.84	0.65	0.76	0.61	1.00	0.73	0.71	0.56	0.94	0.31	0.40
	Mistral-Nemo	0.59	0.48	0.32	0.96	0.52	0.70	0.54	1.00	0.57	0.60	0.44	0.99	0.51	0.27
	Majority Vote	0.72	0.58	0.43	0.88	0.63	0.75	0.60	1.00	0.70	0.69	0.53	0.95	0.35	0.38
Multi-Label-1	Llama-3 70b	0.71	0.58	0.47	0.73	0.78	0.82	0.73	0.93	0.76	0.73	0.63	0.87	0.26	0.48
	Gemma-2 9b	0.69	0.55	0.46	0.69	0.64	0.75	0.61	0.98	0.71	0.68	0.56	0.88	0.32	0.39
	Phi-3.5-MoE	0.73	0.59	0.47	0.79	0.66	0.74	0.63	0.89	0.71	0.68	0.57	0.86	0.32	0.42
	GPT-4o-mini	0.71	0.57	0.45	0.77	0.72	0.79	0.67	0.96	0.74	0.71	0.59	0.89	0.29	0.44
	Mistral-Nemo	0.64	0.45	0.65	0.34	0.59	0.71	0.58	0.93	0.70	0.65	0.59	0.73	0.31	0.36
	Majority Vote	0.71	0.58	0.49	0.71	0.66	0.76	0.62	0.99	0.73	0.70	0.58	0.89	0.30	0.42
Multi-Label-2	Llama-3 70b	0.74	0.61	0.54	0.70	0.71	0.76	0.69	0.85	0.75	0.71	0.64	0.80	0.26	0.48
	Gemma-2 9b	0.70	0.55	0.41	0.84	0.63	0.75	0.60	0.98	0.69	0.67	0.53	0.93	0.36	0.37
	Phi-3.5-MoE	0.72	0.58	0.45	0.82	0.67	0.72	0.65	0.79	0.70	0.66	0.57	0.80	0.32	0.42
	GPT-4o-mini	0.69	0.56	0.60	0.52	0.66	0.72	0.64	0.83	0.72	0.67	0.63	0.72	0.28	0.44
	Mistral-Nemo	0.70	0.56	0.41	0.86	0.53	0.70	0.54	1.00	0.65	0.65	0.49	0.95	0.41	0.32
	Majority Vote	0.74	0.60	0.50	0.75	0.64	0.75	0.61	0.96	0.73	0.70	0.57	0.89	0.31	0.42
Unrestricted	Llama-3 70b	0.71	0.57	0.48	0.69	0.71	0.76	0.69	0.85	0.73	0.69	0.61	0.80	0.28	0.44
	Gemma-2 9b	0.70	0.56	0.42	0.84	0.72	0.74	0.72	0.76	0.70	0.66	0.57	0.79	0.32	0.43
	Phi-3.5-MoE	0.70	0.55	0.41	0.85	0.73	0.74	0.74	0.74	0.69	0.66	0.57	0.78	0.32	0.42
	GPT-4o-mini	0.68	0.54	0.40	0.80	0.73	0.77	0.72	0.82	0.70	0.67	0.57	0.82	0.32	0.41
	Mistral-Nemo	0.51	0.43	0.28	0.91	0.52	0.63	0.54	0.78	0.50	0.54	0.40	0.82	0.57	0.27
	Majority Vote	0.68	0.54	0.41	0.79	0.72	0.76	0.71	0.81	0.70	0.67	0.57	0.81	0.32	0.42

4.2 SPAADE-DR Dataset

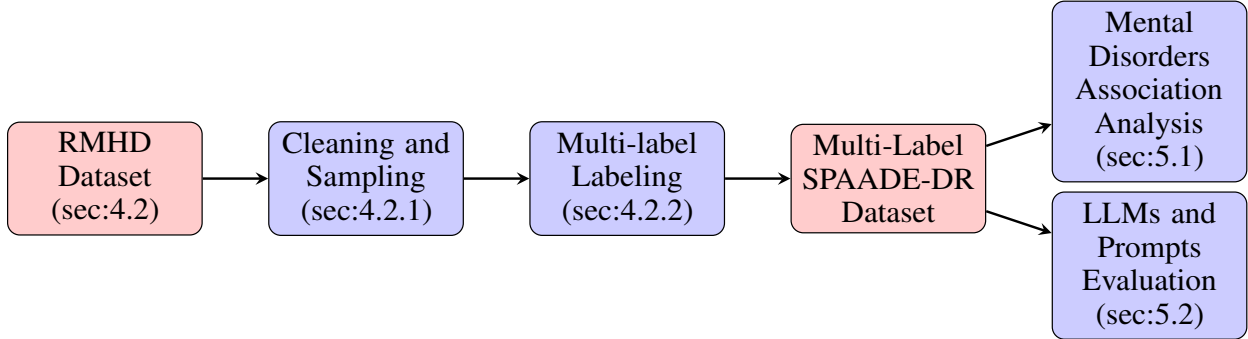


Figure 5: Workflow for the SPAADE-DR Dataset Process

As mentioned earlier, the RMHD dataset originally consists of posts from 17 mental health subreddits and 11 non-mental health subreddits. For this study, posts from r/conspiracy, r/jokes, r/teaching, r/personalfinance, and r/legaladvice are selected as control samples, and posts from r/adhd, r/anxiety, r/depression, r/EDAnonymous (Eating Disorder), r/ptsd, and r/suicidewatch as samples representing various mental health disorders. Since the dataset uses the subreddit each post originates from as its label, it is not entirely clean. Some posts are incorrectly labeled as positive for a disorder simply because they are posted in a particular mental health subreddit, even if they do not reflect the condition. To address this, a cleaning and sampling process is applied, followed by multi-label labeling, as outlined in the workflow shown in Figure 5.

4.2.1 Cleaning and Sampling

Initially, the dataset is cleaned using unsupervised and semi-supervised learning techniques, including KMeans, DBSCAN, and One-Class SVM. These methods leverage both the features provided with the dataset and word embeddings to identify and remove outliers or irrelevant posts. However, these techniques fail to produce satisfactory results. Furthermore, due to the large size of the RMHD dataset, a smaller, cleaner subset must be sampled for further analysis.

To address both cleaning and sampling efficiently, 600 posts are initially sampled from each mental health disorder subreddit. The LLaMA-3 70b model, identified as the best-performing LLM, is then employed to evaluate and classify these posts. Based on the model’s predictions, posts predicted as negative for the corresponding disorder are manually reviewed, and true negatives are removed. Following the cleaning process, 500 posts are selected from the remaining subset of 600 for each disorder, ensuring a clean and representative sample for further analysis.

4.2.2 Multi-label labeling

After sampling, the SPAADE-DR dataset is labeled using the most effective prompt and LLMs identified from the evaluation conducted on the Depseverity-Dreaddit dataset. The single-label prompt is applied with LLaMA-3 70b, GPT-4o-mini, and Phi-3.5 MoE to ensure accurate labeling. Since each sample originates from a specific mental disorder subreddit, the original label is retained as the true label for that condition. The remaining five disorders are then annotated using the single-label prompt, ensuring comprehensive multi-label classification. This process allows us to transform the dataset from single-label to multi-label, enabling more comprehensive analysis of co-occurring mental health conditions across the dataset. After labeling, the distribution of each label is calculated to represent the new multi-label characteristics of the dataset, as presented in Table:4.

Table 4: The Dataset’s Label distribution after conversion to multi-label using Llama-3 70b, GPT-4o mini, and Phi-3.5-MoE

LLM	Label\Disorder	ADHD	Anxiety	Depression	Eating Disorder	PTSD	Suicide
Llama-3 70b	Positive	1546	2547	1902	564	1651	971
	Negative	1954	953	1598	2936	1849	2529
Phi-3.5-MoE	Positive	791	3003	2190	584	1496	1153
	Negative	2709	497	1310	2916	2004	2347
GPT-4o-mini	Positive	1783	3125	2419	735	1796	1477
	Negative	1717	375	1081	2765	1704	2023

From Table:4, it is evident that the label distributions produced by Phi-3.5-MoE are consistently skewed towards either the positive or negative side for all disorders. In contrast, LLaMA-3 70b and GPT-4o-mini exhibit mostly balanced distributions across the disorders.

5 Additional Analysis

After constructing the multi-label SPAADE-DR dataset, further analyses are conducted to explore associations and comorbidities among various mental disorders, providing deeper insights into their association and comorbidity. Additionally, the performance of multi-label and unrestricted prompts is re-examined to evaluate how they adapt to an increase in the number of labels, from 2 to 6, compared to single-label prompts.

5.1 Comorbidities between Mental Disorders

To explore the comorbidities between mental disorders, a contingency heatmap (as shown in Fig. 6) is created to visually represent the relationships between 2 disorders, A and B. In this visualization, the y-axis lists the first disorder A, while the x-axis shows the second disorder B. Each cell in the heatmap displays the proportion of samples that exhibit disorder B status (positive or negative) within groups defined by disorder A status (positive or negative). For example, the cell indicating the percentage of samples negative for PTSD among those positive for suicide offers insight into the co-occurrence trends between these two conditions.

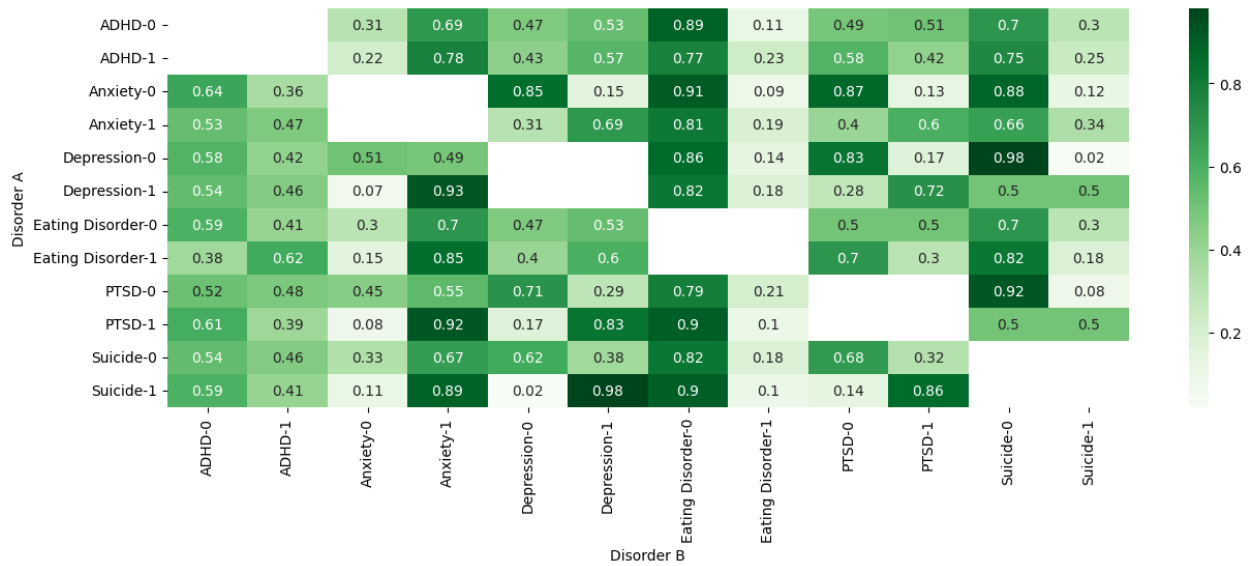


Figure 6: Contingency matrix showing associations between mental disorders (comorbidity)

This analysis reveals several key insights into the comorbidities between mental disorders:

- From figure:6 it is clear that the likelihood of suicide without depression is low, with only 1.5% of suicide-positive samples being depression-negative. This suggests a strong association between depression and suicidal tendencies.
- Among samples positive for PTSD, only 7.6% do not also have an anxiety disorder, indicating a high comorbidity rate between PTSD and anxiety.
- Depression and anxiety show a strong association, with 93% of individuals diagnosed with depression also exhibiting symptoms of anxiety. This high comorbidity rate highlights the close relationship between these disorders, which often occur together in clinical settings.

The odds ratio (OR) between disorders is then calculated, as depicted in Fig. 7. The OR serves as a statistical measure of association between two events, commonly used in case-control studies, medical research, and analyses of relationships within datasets. It helps determine how strongly the presence or absence of one event is associated with the presence or absence of another event.

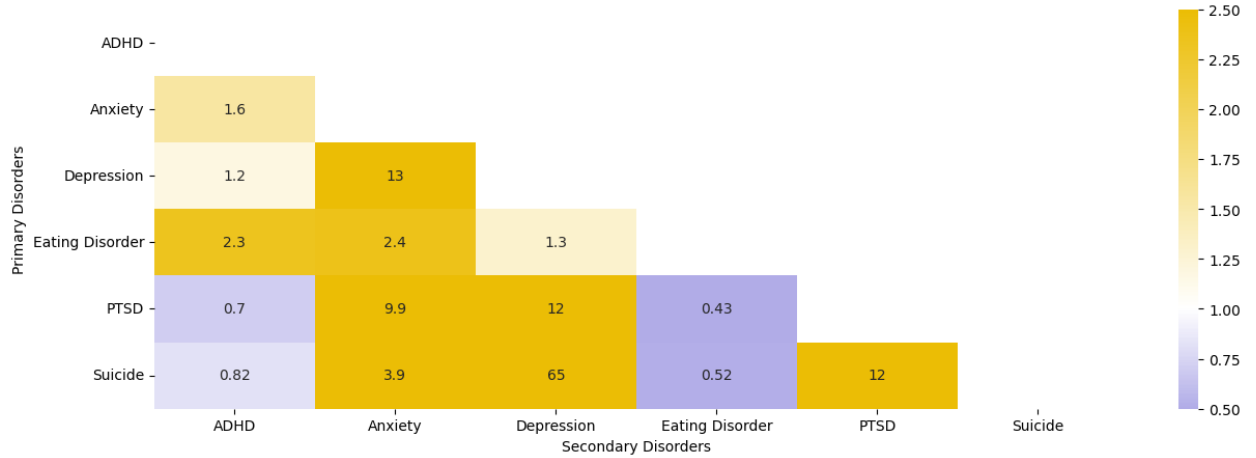


Figure 7: Odds Ratio between mental disorders

Figure 7 shows that suicide and depression are the two most positively associated disorders, while PTSD and eating disorders are the two most negatively associated disorders. Additionally, the combinations of ADHD & Suicide and ADHD & Depression are the least associated disorders. These comorbidity patterns provide important insights into the co-occurrence of mental disorders and underline the importance of multi-label classification approaches, which can capture these complex relationships within mental health data. This analysis of the association and comorbidity between disorders is aligned with real-world patterns documented in other studies [39, 40].

5.2 Prompts and LLMs evaluation on 6 disorders

The multi-label and unrestricted prompts are re-evaluated on the labeled SPAADE-DR dataset to examine how the performance of each model is affected when the number of disorders increases from 2 (as in the Depseverity-Dreaddit dataset) to 6. The evaluation utilizes metrics such as CBA, CF1, OBA, OF1, HL and multi-class BA.

Table 5: Multi-label and Unrestricted prompts evaluation ana 6 disorders multi-labeled RMHD Dataset

Prompt	Disorder LLM Metric	ADHD		Anxiety		Depression		Eating Disorder		PTSD		Suicide		Multi-label			Multi-Class
		BA	F1	BA	F1	BA	F1	BA	F1	BA	F1	BA	F1	GBA	OF1	HL	BA
Multi-Label	Llama-3 70b	0.51	0.06	0.63	0.63	0.67	0.63	0.51	0.04	0.61	0.45	0.71	0.58	0.64	0.50	0.33	0.09
	GPT-4o-mini	0.66	0.48	0.74	0.77	0.78	0.82	0.96	0.86	0.69	0.57	0.88	0.84	0.78	0.73	0.20	0.16
	Phi-3.5-MoE	0.51	0.07	0.62	0.68	0.59	0.65	0.89	0.85	0.55	0.45	0.69	0.55	0.65	0.56	0.33	0.08
Unrestricted	Llama-3 70b	0.50	0.02	0.68	0.81	0.69	0.71	0.47	0.08	0.57	0.43	0.74	0.61	0.66	0.58	0.32	0.06
	GPT-4o-mini	0.58	0.28	0.76	0.88	0.81	0.82	0.80	0.68	0.62	0.40	0.84	0.80	0.76	0.71	0.21	0.13
	Phi-3.5-MoE	0.50	0.02	0.64	0.77	0.61	0.57	0.64	0.38	0.53	0.35	0.68	0.54	0.62	0.52	0.35	0.08
Single_Label	Llama-3 70b	0.97	0.97	1.00	1.00	0.99	0.99	0.98	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.01	0.90
	GPT-4o-mini	0.76	0.75	0.69	0.89	0.83	0.87	0.95	0.84	0.89	0.88	0.89	0.78	0.86	0.85	0.15	0.18
	Phi-3.5-MoE	0.68	0.54	0.74	0.91	0.87	0.89	0.91	0.87	0.85	0.83	0.92	0.85	0.86	0.84	0.14	0.17

From Table 5, it is evident that in contrast to the first experiment, where Llama-3 70b emerged as the top-performing model closely followed by GPT-4o-mini and Phi-3.5-MoE, GPT-4o-mini outperforms the other models when the number of disorders in multi-label and unrestricted prompts increases. Meanwhile, the performance of Llama-3 70b and Phi-3.5-MoE significantly declines under these conditions. This indicates that GPT-4o-mini is the most robust and consistent LLM.

Additionally, when using the single-label prompt, Llama-3 70b achieves the highest scores. This is likely due to the alignment of the Llama-3 70b model with the prompt structure used during the data labeling process. Notably, both GPT-4o-mini and Phi-3.5-MoE attain BA and F1 scores above 0.80 across all mental disorders except for ADHD and Anxiety when using the single-label prompt. They both also achieve an overall balanced accuracy (OBA) of 0.86, which is notably high. This suggests that single-label prompts are the most robust, primarily because, despite being computationally more expensive, their performance remains unaffected by the number of labels.

6 Conclusion & Future Directions

This paper explores the potential of large language models (LLMs) to assist in data annotation for mental health research, particularly in the context of social media platforms. Mental health disorders such as depression, anxiety, and PTSD pose significant challenges on a global scale. Social media offer an extensive and valuable source of data that can be harnessed for research into these conditions. Leveraging large language models (LLMs) enables us to improve dataset accuracy, by making them more reflective of real-world data and increasing the efficiency of data annotation processes, thus expanding the volume of high-quality annotated data. This, in turn, can facilitate more effective diagnosis and intervention strategies, ultimately contributing to improved mental health outcomes.

6.1 Conclusion

The integration of LLMs into mental health research presents a promising approach for tackling the challenges of large-scale data annotation within this field. The capability of LLMs to process and interpret complex language patterns allows them not only to aid in the identification and diagnosis of mental health conditions through social media analysis but also to enhance the quality and depth of dataset annotations. Notably, LLMs can identify subtle indicators and comorbidities among mental health conditions that traditional methods might miss.

This study leverages LLMs to transition from single-label to multi-label annotation, facilitating more nuanced classification of mental health conditions. Various prompting techniques were tested to identify the most effective strategies, revealing that despite the added complexity of increasing the number of labels or classes, the performance of LLMs remained consistent and robust. These findings highlight the potential of LLMs to revolutionize data annotation processes in mental health research. Additionally, a novel multi-label dataset encompassing six distinct mental disorders was developed, expanding the existing pool of resources and enabling more sophisticated analyses in mental health research.

6.2 Future Directions

While LLMs show significant promise, several avenues remain for future exploration:

- **Improving model precision for specific mental health conditions:** LLMs need to be further fine-tuned to detect nuances between different mental health disorders. Customizing models for specific conditions such as depression, anxiety, and severe disorders like psychosis will enhance their diagnostic effectiveness.
- **Extending multi-label annotation to other domains:** The multi-label annotation method facilitated by LLMs in this study can be expanded beyond mental health to other domains. Future work should explore how this approach can be applied to various types of data, including text, images, and across different languages.
- **Real-time monitoring and intervention:** LLMs could eventually be integrated into real-time monitoring systems, which analyze social media data to detect emerging mental health issues. Such systems could provide timely alerts to health professionals or directly to users when concerning behavioral patterns are identified, potentially preventing crises.

In summary, the use of LLMs for mental health data annotation is still in its early stages, but the potential benefits are considerable. Continued refinement of these models, along with addressing current challenges, could position LLMs as a key component in the future of mental health research and interventions.

References

- [1] American Psychiatric Association. *Diagnostic and Statistical Manual of Mental Disorders: DSM-5-TR*. American Psychiatric Association Publishing, Washington, DC, 5th, text revision edition, 2022.
- [2] World Health Organization. Depression, 2023. Accessed: September 14, 2024.
- [3] World Health Organization. Anxiety disorders, 2022. Accessed: 2024-09-26.
- [4] World Health Organization. *World Mental Health Report: Transforming Mental Health for All*. World Health Organization, Geneva, 2022. Accessed: 2024-09-26.
- [5] Munmun De Choudhury and Sushovan De. Mental health discourse on reddit: Self-disclosure, social support, and anonymity. *Proceedings of the International AAAI Conference on Web and Social Media*, 8(1):71–80, May 2014.
- [6] Glen Coppersmith, Mark Dredze, and Craig Harman. Quantifying mental health signals in Twitter. In Philip Resnik, Rebecca Resnik, and Margaret Mitchell, editors, *Proceedings of the Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 51–60, Baltimore, Maryland, USA, June 2014. Association for Computational Linguistics.
- [7] Sharath Chandra Guntuku, David B Yaden, Margaret L Kern, Lyle H Ungar, and Johannes C Eichstaedt. Detecting depression and mental illness on social media: an integrative review. *Current Opinion in Behavioral Sciences*, 18:43–49, 2017. Big data in the behavioural sciences.
- [8] Philip Resnik, William Armstrong, Leonardo Claudino, Thang Nguyen, Viet-An Nguyen, and Jordan Boyd-Graber. Beyond LDA: Exploring supervised topic modeling for depression-related language in Twitter. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 99–107, Denver, Colorado, June 5 2015. Association for Computational Linguistics.
- [9] Glen Coppersmith, Mark Dredze, Craig Harman, and Kristy Hollingshead. From ADHD to SAD: Analyzing the language of mental health on Twitter through self-reported diagnoses. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 1–10, Denver, Colorado, June 5 2015. Association for Computational Linguistics.
- [10] Stevie Chancellor, Michael L. Birnbaum, Eric D. Caine, Vincent M. B. Silenzio, and Munmun De Choudhury. A taxonomy of ethical tensions in inferring mental health states from social media. In *Proceedings of the Conference on Fairness, Accountability, and Transparency, FAT* ’19*, page 79–88, New York, NY, USA, 2019. Association for Computing Machinery.
- [11] Matthew Matero, Akash Idnani, Youngseo Son, Salvatore Giorgi, Huy Vu, Mohammad Zamani, Parth Limbachiya, Sharath Chandra Guntuku, and H. Andrew Schwartz. Suicide risk assessment with multi-level dual-context language and BERT. In Kate Niederhoffer, Kristy Hollingshead, Philip Resnik, Rebecca Resnik, and Kate Loveys, editors, *Proceedings of the Sixth Workshop on Computational Linguistics and Clinical Psychology*, pages 39–44, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [12] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Nee-lakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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. *CoRR*, abs/2005.14165, 2020.
- [13] Ruiyang Qin, Ryan Cook, Kai Yang, Ahmed Abbasi, David Dobolyi, Salman Seyedi, Emily Griner, Hyeokhyen Kwon, Robert Cotes, Zifan Jiang, and Gari Clifford. Language models for online depression detection: A review and benchmark analysis on remote interviews. *ACM Trans. Manage. Inf. Syst.*, August 2024. Just Accepted.
- [14] 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, 2023.
- [15] Haoyan Luo and Lucia Specia. From understanding to utilization: A survey on explainability for large language models, 2024.
- [16] Yang Liu, Yuanshun Yao, Jean-Francois Ton, Xiaoying Zhang, Ruocheng Guo, Hao Cheng, Yegor Klochkov, Muhammad Faaiz Taufiq, and Hang Li. Trustworthy llms: a survey and guideline for evaluating large language models’ alignment, 2024.
- [17] Biodoumoye George Bokolo and Qingzhong Liu. Deep learning-based depression detection from social media: Comparative evaluation of ml and transformer techniques. *Electronics*, 12(21), 2023.

- [18] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.
- [19] Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, Amrita Bhattacharjee, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. Large language models for data annotation: A survey, 2024.
- [20] Shuo Yang, Zirui Shang, Yongqi Wang, Derong Deng, Hongwei Chen, Qiyuan Cheng, and Xinxiao Wu. Data-free multi-label image recognition via llm-powered prompt tuning, 2024.
- [21] Youri Peskine, Damir Korenčić, Ivan Grubisic, Paolo Papotti, Raphael Troncy, and Paolo Rosso. Definitions matter: Guiding GPT for multi-label classification. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 4054–4063, Singapore, December 2023. Association for Computational Linguistics.
- [22] Jasmin Bogatinovski, Ljupco Todorovski, Saso Dzeroski, and Dragi Kocev. Comprehensive comparative study of multi-label classification methods. *CoRR*, abs/2102.07113, 2021.
- [23] Elsbeth Turcan and Kathleen McKeown. Dreaddit: A reddit dataset for stress analysis in social media. *arXiv preprint arXiv:1911.00133*, 2019.
- [24] Usman Naseem, Adam G Dunn, Jinman Kim, and Matloob Khushi. Early identification of depression severity levels on reddit using ordinal classification. In *Proceedings of the ACM Web Conference 2022*, pages 2563–2572, 2022.
- [25] Danielle Mowery, Craig Bryan, and Mike Conway. Towards developing an annotation scheme for depressive disorder symptoms: A preliminary study using Twitter data. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 89–98, Denver, Colorado, June 5 2015. Association for Computational Linguistics.
- [26] Kurt Kroenke and Robert L Spitzer. The phq-9: A new depression diagnostic and severity measure. *Psychiatric Annals*, 32(9):509–515, 2002.
- [27] Kurt Kroenke, Robert L Spitzer, and Janet B Williams. The phq-9: Validity of a brief depression severity measure. *Journal of General Internal Medicine*, 16(9):606–613, 2001.
- [28] Manas Gaur, Amanuel Alambo, Joy Prakash Sain, Ugur Kursuncu, Krishnaprasad Thirunarayan, Ramakanth Kavuluru, Amit Sheth, Randy Welton, and Jyotishman Pathak. Knowledge-aware assessment of severity of suicide risk for early intervention. In *The World Wide Web Conference, WWW '19*, page 514–525, New York, NY, USA, 2019. Association for Computing Machinery.
- [29] Daniel Low, Laurie Rumker, Tanya Talker, John Torous, Guillermo Cecchi, and Satrajit Ghosh. Natural language processing reveals vulnerable mental health support groups and heightened health anxiety on reddit during covid-19: An observational study. *Journal of medical Internet research*, 22, 09 2020.
- [30] Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- [31] OpenAI. Gpt-4o mini. <https://openai.com/index/gpt-4o-mini-advancing-cost-efficient-intelligence/>, 2024.
- [32] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [33] Mistral AI team. Mistral-nemo. <https://mistral.ai/news/mistral-nemo//>, 2024.
- [34] Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- [35] Abdelrahman Hanafi, Mohammed Saad, Noureldin Zahran, Radwa J. Hanafy, and Mohammed E. Fouda. A comprehensive evaluation of large language models on mental illnesses, 2024.
- [36] Shilong Liu, Lei Zhang, Xiao Yang, Hang Su, and Jun Zhu. Query2label: A simple transformer way to multi-label classification. *CoRR*, abs/2107.10834, 2021.
- [37] Ruijie Yao, Sheng Jin, Lumin Xu, Wang Zeng, Wentao Liu, Chen Qian, Ping Luo, and Ji Wu. Gkgnet: Group k-nearest neighbor based graph convolutional network for multi-label image recognition, 2024.
- [38] Guoqiang Wu and Jun Zhu. Multi-label classification: do hamming loss and subset accuracy really conflict with each other? *CoRR*, abs/2011.07805, 2020.

- [39] Ali M AL-Asadi, Britt Klein, and Denny Meyer. Multiple comorbidities of 21 psychological disorders and relationships with psychosocial variables: A study of the online assessment and diagnostic system within a web-based population. *J Med Internet Res*, 17(3):e55, Feb 2015.
- [40] J. J. McGrath, C. C. W. Lim, O. Plana-Ripoll, Y. Holtz, E. Agerbo, N. C. Momen, P. B. Mortensen, C. B. Pedersen, J. Abdulmalik, S. Aguilar-Gaxiola, and et al. Comorbidity within mental disorders: a comprehensive analysis based on 145 990 survey respondents from 27 countries. *Epidemiology and Psychiatric Sciences*, 29:e153, 2020.