

MENTA: A SMALL LANGUAGE MODEL FOR ON-DEVICE MENTAL HEALTH PREDICTION

Tianyi Zhang^{1*}, Xiangyuan Xue^{2*}, Lingyan Ruan¹, Shiya Fu¹, Feng Xia³, Simon D’Alfonso¹, Vassilis Kostakos¹, Ting Dang¹, Hong Jia²

¹The University of Melbourne ²The University of Auckland ³RMIT University

ABSTRACT

Mental health conditions affect hundreds of millions globally, yet early detection remains limited. While large language models (LLMs) have shown promise in mental health applications, their size and computational demands hinder practical deployment. Small language models (SLMs) offer a lightweight alternative, but their use for social media-based mental health prediction remains largely underexplored. In this study, we introduce Menta, the first optimized SLM fine-tuned specifically for multi-task mental health prediction from social media data. Menta is jointly trained across six classification tasks using a LoRA-based framework, a cross-dataset strategy, and a balanced accuracy-oriented loss. Evaluated against nine state-of-the-art SLM baselines, Menta achieves an average improvement of 15.2% across tasks covering depression, stress, and suicidality compared with the best-performing non-fine-tuned SLMs. It also achieves higher accuracy on depression and stress classification tasks compared to 13B-parameter LLMs, while being approximately $3.25\times$ smaller. Moreover, we demonstrate real-time, on-device deployment of Menta on an iPhone 15 Pro Max, requiring only approximately 3GB RAM. Supported by a comprehensive benchmark against existing SLMs and LLMs, Menta highlights the potential for scalable, privacy-preserving mental health monitoring. Code is available at: <https://xxue752-nz.github.io/menta-project/>

1 INTRODUCTION

Mental health disorders such as depression, anxiety, and suicidality affect hundreds of millions of people worldwide and constitute one of the leading contributors to the global burden of disease (Organization et al., 2017). Between 30% and 50% of people globally experience stress (Piao et al., 2024); an estimated 5.7% of adults suffer from depression (World Health Organization, 2023); and more than 720,000 people die by suicide each year (World Health Organization, 2025a). As mental health issues have continued to rise globally over the past few decades (Goodwin et al., 2022; Piao et al., 2024; Weaver et al., 2025), there is an urgent need to enhance our understanding, diagnosis, and monitoring of these conditions through advanced technologies.

Despite growing demand for mental health support, significant diagnosis gaps persist globally. Barriers such as limited clinic hours, geographic inaccessibility, and workforce shortages hinder timely assessment and early self-awareness (World Health Organization, 2025b). Structural, patient-level, and systemic obstacles further delay diagnosis for those with serious mental illness (Wiesepape et al., 2025). This mismatch between high need and limited resources reflects the constraints of traditional diagnostic models (Stein et al., 2022). Moreover, these traditional delivery of in-person sessions and self-reported questionnaires struggles to scale, and faces systemic bottlenecks, including workforce shortages (Organization, 2025; Brahmabhatt & Schpero, 2024), high costs (Patel et al., 2018), geographic disparities (Yang & Zhang, 2025), and stigma (Clement et al., 2015). Mental health wait times can range from 3 to 18 months, varying by region and provider availability (McMahan et al., 2022; Yang & Zhang, 2025; Patel et al., 2018). With the growing need for early detection and real-time monitoring, developing more seamless methods to identify early signs is critical for delivering

*These authors contributed equally to this work.

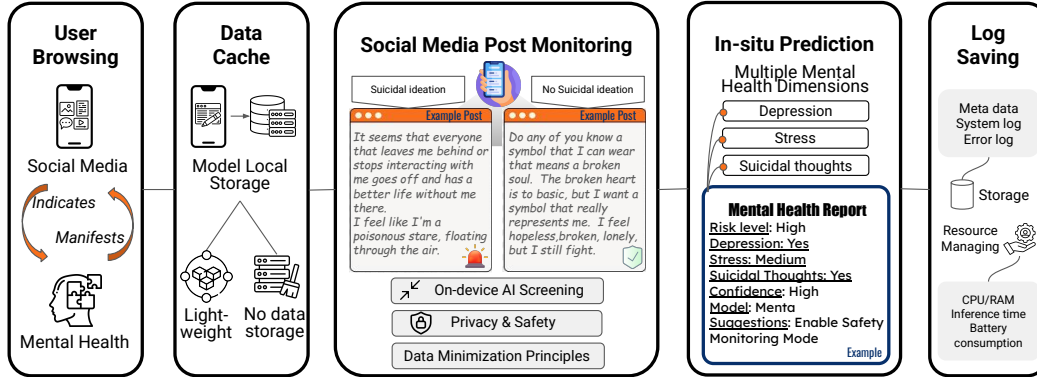


Figure 1: On-device mental health monitoring workflow using SLMs, from user browsing of social media, to data collection, post monitoring, in-situ prediction and meta log saving.

timely support (Patel et al., 2018). This underscores the importance of portable, resource-efficient solutions such as small language models for scalable mental health detection and intervention.

The rise of portable, on-device technologies offers a promising path for the early detection of mental health issues, particularly through passive, real-time analysis of social media posts and smartphone applications. Smartphone-based linguistic data from social media enables scalable, passive monitoring of psychological states, facilitating early identification of mood disturbances, stress, or suicide risk. For example, such systems can support suicide prevention efforts (Areán et al., 2016) and capture patterns in social media use that are relevant to mental well-being (Hamilton et al., 2025). Running locally on personal devices, these models preserve user privacy while analyzing language inputs for indicators of depression or stress (Shin et al., 2023). As a result, on-device deployment enables scalable, privacy-preserving mental health monitoring, extending detection beyond clinical settings and making access easier for a broader population.

In this context, language models offer a powerful means of addressing the limitations of traditional approaches to mental health monitoring (Xu et al., 2024). Language models outperform supervised learning classifiers in mental health prediction tasks for capturing deeper context, information acquisition, and generalization capabilities (Jin et al., 2025), making them more practical for web-scale and real-world mental health monitoring.

However, the application of large language models (LLMs) faces significant challenges: they are computationally expensive to train and deploy, and they raise privacy concerns in terms of processing sensitive data on remote servers (Jin et al., 2025). Their large size also makes them impractical for resource-constrained on-device settings, underscoring the need for more efficient alternatives. By contrast, Small Language Models (SLMs), offer comparable or even superior potential for digital health applications due to their lightweight architectures, lower inference costs, and feasibility of local deployment (Abdin et al., 2024; Jia et al., 2025). Deploying on digital devices locally, a privacy-preserving, on-device workflow is expected for mental health prediction from social media activity using lightweight SLMs which emphasizes in-situ AI screening across multiple mental health dimensions without storing personal data, enabling real-time, low-resource, and ethically-aligned monitoring (Figure 1). Yet, the potential of SLMs for mental health prediction remains largely underexplored.

Despite work on the development of SLMs, there is no agreement yet on a standardized definition of SLMs (Wang et al., 2024). Given the lightweight requirement for such real-time mental health status detection tasks on social media posts, we define SLMs to be models ranging from a minimum of 1B to strictly under 7B parameters are commonly considered ‘small’ in contrast to full-scale LLMs (10B+) while those more than 7B to be large-scale small language models (LSSLMs) (Abdin et al., 2024; Team, 2024).

In this work, we explore the potential of SLMs in predicting digital mental health outcomes, and prove that SLMs, when carefully designed and optimized, serve as practical alternatives for multi-task mental health assessment on-device. To investigate this, we examine SLM performance across six mental health classification tasks spanning stress, depression, and suicidality. We introduce Menta, a compact model fine-tuned using a LoRA-based multi-task framework with cross-dataset training and a balanced accuracy-oriented loss. Menta performs strongly across diverse mental health tasks, particularly in stress and depression detection. We further demonstrate its practical value by deploying it on an iPhone 15 Pro Max, achieving real-time inference with minimal memory usage, enabling private and scalable on-device mental health support. Our key contributions are as follows:

- We present Menta, a fine-tuned SLM tailored for multi-task mental health prediction from social media text, optimized for both accuracy and deployability.
- Menta leverages a weighted-loss, LoRA-based multi-task training framework designed to address data imbalance and maximize generalizability across multiple clinical classification tasks.
- We show that Menta offers a compelling balance between predictive accuracy and computational efficiency, outperforming several SLMs and rivaling LLMs on six mental health task.
- We provide an open-source deployment pipeline and model implementation for on-device inference, supporting privacy-preserving and resource-efficient mental health applications.

2 RELATED WORK

Social Media Posts: Indication of Mental Health. Social media posts offer a valuable lens into individuals’ emotions, thoughts, and mental health status. Globally, approximately 4.8 billion people, nearly 60% of the population, use social media, spending a combined 11.5 billion hours on these platforms each day University of Maine (2023). With the increasing use of social media platforms for emotional expression and peer support, user-generated text has become an important resource for understanding psychological well-being (Hussain et al., 2025; Zhunis et al., 2022) and identifying early warning signs of mental distress (Chancellor & De Choudhury, 2020; De Choudhury et al., 2013). For instance, a study (Kim et al., 2020) built a classifier on Reddit posts to detect multiple mental disorders, while the EmoMent corpus (Atapattu et al., 2022) shows that emotion annotations in user text correlate with mental illness indicators even in non-Western populations. These findings support the potential for social media-based models to serve as early warning systems for mental health support.

Early approaches for textual information analysis among social media posts relied on traditional machine learning and feature-engineered models to classify or detect mental health status. For example, Jiang et al. used contextualized word embeddings such as ELMo and BERT to detect mental health conditions from Reddit posts (Jiang et al., 2020), showing that deep contextualized models significantly outperform traditional bag-of-words or static embeddings for mental health classification tasks. Similarly, Sarkar et al. developed a multi-task learning model that jointly predicts depression and anxiety from Reddit posts, which outperforms single-task baselines by leveraging shared features across related conditions (Sarkar et al., 2022). Another studies developed a scalable deep-learning screening tool for suicide risk with potential clinical implications for early intervention such as CNN (Coppersmith et al., 2018), LSTM (Coppersmith et al., 2018) and SVM (Ji et al., 2018).

More recently, LLMs have been adapted to this domain. A prior study used ChatGPT in zero-shot classification tasks across three mental health domains (stress, depression, suicidality) on social media datasets, validating the ability of LLMs with no prior knowledge in mental health prediction tasks (Lamichhane, 2023). Another study focuses only on the depression detection task from online text while collaborates the finding that fine-tuned LSSLMs achieve great improvements over prior state-of-art models (Shah et al., 2025). Moreover, the interpretability of LLMs in prediction tasks such as depressive symptoms (Bolegave & Bhattacharya, 2025; Chen & Lin, 2025; Belcastro et al., 2025) and emotional states (Yang et al., 2023) have also been explored, with the finding that prompt engineering with emotional cues and few-shot examples improves performance. Importantly, com-

pared to traditional algorithmic models, language models offer not only strong predictive capabilities but also human-readable interpretations and the flexibility to support multiple downstream tasks within a unified framework.

A more comprehensive study introduced Mental-LLM (Xu et al., 2024), demonstrated that instruction-tuned LLMs show large gains in balanced accuracy and can outperform task-specific baselines, with the best finetuned models (Mental-Alpaca and Mental-FLAN-T5) outperform GPT-3.5 by approximately 10.9% and beat GPT-4 by about 4.8%. They also improved model generalizability fine-tuned on diverse datasets. Follow-up work has extended this line of research with domain-specialized models. Similarly, Shi et al. proposed a LSSLM called MentalQLM for mental health tasks by using a base backbone (7B) and then applying a dual LoRA strategy (Hu et al., 2022), achieving comparable or even superior performance to much larger LLMs on several mental health diagnostic tasks (Shi et al., 2025). While these studies demonstrate the feasibility of using LLMs for mental health outcome prediction, significant challenges persist, including concerns over privacy and security (Sarwar, 2025), deployment limitations (Maurya et al., 2025), and high computational costs (Wang et al., 2023).

Small Language Models and Mental Health. To address this gap of application bottlenecks, SLMs emerged as a pragmatic alternative of LLMs (Van Nguyen et al., 2024), optimized for efficiency through techniques such as quantization (Frantar et al., 2022), knowledge distillation (Kang et al., 2023) and light-weight architectures (Wang et al., 2025a). SLMs were designed to deliver core Natural Language Processing capabilities such as text classification (Pecher et al., 2024), summarization (Wang et al., 2025b), and question-answering (Lee et al., 2024), in resource-constrained settings including mobile devices, embedded systems, and edge computing.

Offering a compelling alternative to LLMs, SLMs reduce computational overhead while retaining strong performance in many specialized tasks. Moreover, recent work shows that in classification settings with limited data, well-tuned smaller models can match or even surpass larger models when the task is focused and domain-specific including text classification (Lepagnol et al., 2024; Luo et al., 2023) and healthcare domains (Gondara et al., 2025). These findings highlight that SLMs offer an efficient, scalable, and domain-optimized solution, making them a promising direction for mental-health detection from social media. However, previous work has typically focused on a single task or domain and has not explored multi-domain settings that combine multiple mental health conditions such as depression, stress, and suicidality.

Unlike prior LLM-based (Sarkar et al., 2022; Xu et al., 2024; Yang et al., 2024b; Kim et al., 2024; Shi et al., 2025) which range from 7B to 70B parameters which are impractical for deployment in resource-constrained or privacy-sensitive settings, lightweight SLMs emphasize deployability, reporting inference latency, memory footprint, and approximate cost efficiency. Prior work evaluated several SLMs against LLM baselines across multiple individual mental health understanding tasks with zero-shot and few-shot prompting trained on social media datasets, and concluded that few-shot prompting helps SLMs more (Jia et al., 2025). Also, Kim et al. developed a SLM called mhGPT, which outperforms larger models such as MentaLLaMA and Gemma in mental health tasks despite having far fewer parameters and using less data (Kim et al., 2024). A prior study also suggests SLMs suitable to be integrated in clinical workflows for structured and clinically meaningful tasks beyond generic detection (Aich et al., 2024). However, these studies lack systematic task balancing, rarely optimized with parameter-efficient fine-tuning (LoRA), meaning adaptation is either minimal (prompting) or heavy (full fine-tuning).

Deployment and On-Device Language Models. On-device AI offers a compelling opportunity for mental health monitoring by enabling personalized assessment in real time, while preserving user privacy, reducing latency, and expanding accessibility. Running inference directly on a user’s smartphone avoids sending sensitive linguistic and behavioural data to cloud servers, thereby lowering risks of data leakage and minimizing network dependencies (Mandal et al., 2025). Furthermore, mobile and edge deployment largely improves accessibility by bringing mental health support tools directly to users, addressing traditional barriers such as geographic inaccessibility and long wait times for in-clinic screening (Bunyi et al., 2021). Additionally, because on-device models eliminate round-trip communication with remote servers, they deliver near-instant responses and support real-time detection of changes in language or behaviour indicative of stress or mood deterioration (Ni & Jia, 2025).

Mobile and edge-based mental health monitoring using smartphones and wearable sensors has gained momentum, offering significant advantages in latency, energy efficiency, and data privacy over cloud-dependent systems. Recent work has demonstrated effective deployment of SLMs on consumer devices, such as the Samsung Galaxy S24, for tasks like document assistance (Pham et al., 2024) and health prediction (Wang et al., 2025c), with substantial reductions in memory usage and inference latency. Conceptual analyses further argue that SLMs are better suited than LLMs for real-world, interactive applications due to their lightweight nature and ease of deployment under resource constraints (Belcak et al., 2025). On-device or near-device processing eliminates the need for continuous cloud communication, thereby enabling real-time detection while safeguarding user data—an essential consideration for privacy (Lu et al., 2024).

While recent studies have demonstrated the ability of SLMs on individual mental health tasks in zero- or few-shot settings, existing work remains limited in several key areas. Most notably, prior efforts have primarily focused on single-task classification, often failing to generalize well to more complex, multi-class, and imbalanced datasets. Moreover, little attention has been given to the real-world applicability of these models in terms of on-device deployment under resource constraints. To address these gaps, we propose Menta, a LoRA-based multi-task fine-tuning framework built on Qwen-3, incorporating weighted training to explicitly handle label imbalance. Our method is designed to achieve balanced and consistent performance across six diverse mental health prediction tasks, while maintaining the lightweight efficiency required for practical, privacy-preserving deployment on consumer devices.

3 METHOD

3.1 LANGUAGE MODEL PROMPTING

3.1.1 ZERO-SHOT PROMPTING

Zero-shot prompts were designed to reflect a domain-specific, psychologically grounded context without providing explicit task instructions or label options. Each prompt instructed the model to act as a psychologist evaluating a social media post for indicators of mental health conditions. The prompt construction followed a consistent structure composed of four key elements (Figure 2), designed to encourage the model to rely on contextual understanding of the user’s language and inferred psychological cues:

1. A context statement that framed the post as originating from social media and positioned the model in the role of a psychologist;
2. The user’s text simulating a social media post;
3. A task-specific question corresponding to one of the mental health classification tasks.
4. A response constraint explicitly requiring the model to respond only with a numeric label and avoid hallucination;

3.1.2 PROMPT FORMULATION FOR FEW-SHOT AND FINE-TUNED MODELS

Both few-shot prompting and fine-tuning involve presenting SLMs with illustrative examples that demonstrate how social media posts are categorized into relevant mental health categories. In both setups, the models are instructed to assume the role of a psychologist, evaluating each post based on the examples provided. The prompt is carefully designed to guide the model’s reasoning and classification behavior. Templates of the constructed prompt are shown in Figure 2, including the following elements:

1. A contextual statement framing the post as originating from social media and positioning the model in the role of a psychologist;
2. Example social media posts paired with their corresponding mental health categories;
3. A user text simulating a new social media post;
4. A task-specific question aligned with one of the mental health classification objectives.

Few-shot Prompt Template	Few-shot Example with Full Input-Output
<p>Instruction preview: You are a mental health expert. Your task is to classify the depression level of a person based on their social media post.</p> <p>Here are [number of shot] example(s): <u>Example:</u> Text: "[example]" Question: Which level is the person depressed? Answer: [example] ...</p> <p>Now, please classify the following text: [text] Question: Which level is the person depressed?</p> <p>IMPORTANT: - Respond with ONLY [option1] or [option2]... - [option1] means [meaning1], [option2] means [meaning2], ... - Do not reply anything else</p>	<p>Instruction preview: You are a mental health expert. Your task is to classify the depression level of a person based on their social media post.</p> <p>Here are two examples: <u>Example 1:</u> Text: "Seated on a throne of wire, Mermaids tending to the pyre..." Question: Which level is the person depressed? Answer: minimum <u>Example 2:</u> Text: "Idk Do I tell someone? Do I just quit?..." Question: Which level is the person depressed? Answer: severe</p> <p>Now, please classify the following text: "I'm too unstable. Telling parents is not an option..." Question: Which level is the person depressed?</p> <p>IMPORTANT: - Respond with ONLY `0` or `1` - `1` means depressed, `0` means not depressed - Do not reply anything else</p>
Zero-shot Prompt Template	
[Context Statement] + [User's Social Media Post] + [Task-Specific Question] + [Response Constraint] Your answer is: <eos>	

Figure 2: Prompt instruction templates for mental health classification tasks, specifically Zero-shot prompt, few-shot prompt template, and detailed few-shot prompt with examples.

5. A response constraint explicitly requiring the model to respond only with a numeric label and avoid hallucination;

3.2 MODEL FINE-TUNING

In this section, we focus on Qwen-3, aiming to develop a unified model that balances and enhances performance across distinct mental health tasks. Our tasks span six mental health prediction benchmarks, each formulated as a classification problem. We adopt Low-Rank Adaptation (LoRA) to fine-tune Qwen-3 (4B) for mental health prediction tasks. Instead of updating all model parameters, LoRA injects trainable low-rank matrices into specific weight projections, drastically reducing the number of trainable parameters while preserving reasoning capabilities. In our setting, LoRA adapters are applied to the query and value projections of the transformer attention layers, with rank $r = 16$, scaling factor $\alpha_{\text{LoRA}} = 32$, and dropout 0.05. This configuration achieves efficient adaptation while keeping only $\sim 0.1\%$ of the parameters trainable.

We construct a unified multi-task dataset and apply weighted sampling across tasks to mitigate task imbalance. Specifically, the probability of sampling from task t is proportional to a task weight λ_t , ensuring that underrepresented tasks are adequately seen during training. Within each task, class imbalance is addressed using inverse-frequency weights w_c when computing the loss.

Traditional fine-tuning strategies typically rely on cross-entropy (CE) loss, which optimizes the log-likelihood of the true class. While effective, CE loss can bias the model toward majority classes, leading to degraded performance on minority classes. To address this, we introduce a novel *balanced accuracy* (BACC) *surrogate loss* that provides a differentiable approximation of balanced accuracy, thereby directly encouraging the model to perform more evenly across classes, as shown in Figure 3.

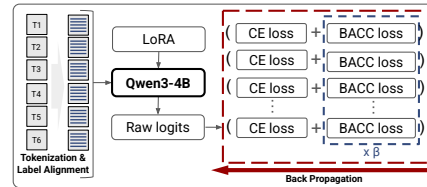


Figure 3: A multi-task training pipeline using Qwen-3 (4B) with LoRA, combining cross-entropy and balanced accuracy (BACC) losses weighted per task for joint optimization through backpropagation.

For a given task t , the standard cross-entropy loss is defined as:

$$\mathcal{L}_{\text{CE}}^{(t)} = -\frac{1}{N_t} \sum_{i=1}^{N_t} \log p(y_i | x_i), \quad (1)$$

where N_t is the number of samples for task t , and $p(y_i | x_i)$ is the predicted probability for the true label y_i , obtained via the softmax over model logits.

To approximate balanced accuracy (BACC), we first compute a *margin* for each class c :

$$m_{i,c} = z_{i,c} - \log \left(\sum_{k \neq c} \exp(z_{i,k}) \right), \quad (2)$$

where $z_{i,c}$ denotes the model logit of sample i for class c . This margin quantifies the relative confidence of the model in predicting class c compared to all other classes.

We then apply a sigmoid function with sharpness parameter α to obtain a soft correctness score:

$$s_{i,c} = \sigma(\alpha \cdot m_{i,c}), \quad (3)$$

where $\sigma(\cdot)$ is the sigmoid function.

Using these soft scores, the true positive rate (TPR) for each class c is estimated as:

$$\widehat{\text{TPR}}_c = \frac{1}{|I_c|} \sum_{i \in I_c} s_{i,c}, \quad (4)$$

where I_c is the set of indices such that $y_i = c$. If no samples belong to class c , the corresponding $\widehat{\text{TPR}}_c$ is defined as 0.

The surrogate BACC loss is then defined as:

$$\mathcal{L}_{\text{BACC}}^{(t)} = 1 - \frac{1}{\sum_c \gamma_c} \sum_c \gamma_c \cdot \widehat{\text{TPR}}_c, \quad (5)$$

where γ_c is a class scaling factor that adjusts the relative importance of each class.

The total multi-task loss is computed as a weighted combination of the cross-entropy loss and the surrogate BACC loss:

$$\mathcal{L}_{\text{total}} = \sum_t \lambda_t \left(\mathcal{L}_{\text{CE}}^{(t)} + \beta \cdot \mathcal{L}_{\text{BACC}}^{(t)} \right), \quad (6)$$

where λ_t is the task-specific weight and β is a fixed hyperparameter that controls the trade-off between cross-entropy and BACC optimization. In this formulation:

- λ_t : task-specific weight (externally specified),
- β : fixed trade-off parameter between \mathcal{L}_{CE} and $\mathcal{L}_{\text{BACC}}$,
- γ_c : class scaling factor,
- α : sigmoid sharpness parameter controlling sensitivity to the margin.

This approach encourages fairness across tasks and classes by directly optimizing a smooth surrogate of balanced accuracy. Unlike standard CE loss, it explicitly penalizes performance gaps between majority and minority classes, improving consistency in multi-task fine-tuning. An ablation study was conducted to compare different training strategies. Specifically, we evaluated the impact of single-task versus multi-task fine-tuning approaches for mental health prediction.

3.3 ON-DEVICE DEPLOYMENT

We developed an on-device mental health evaluation system for iOS using llama.cpp¹ as the inference backend, leveraging GGUF V3-quantized models optimized for mobile hardware. The models

¹<https://github.com/ggml-org/llama.cpp>

were exported from PyTorch to the GGUF format via the transformers-to-gguf conversion pipeline, enabling 4-bit quantization (Q4_K_M) for efficient CPU and GPU execution. The app was implemented in Swift/SwiftUI and integrates llama.cpp through a C++ bridge layer, allowing direct inference within the iOS runtime without external servers.

At runtime, inference is accelerated using Apple Metal for matrix operations, with automatic CPU fallback when GPU utilization is saturated. Thread-level parallelism (up to 8 threads) is managed through llama.cpp’s built-in thread scheduling API, ensuring optimal performance across Apple’s efficiency and performance cores. Model weights are memory-mapped to minimize RAM usage and enable lazy loading of tensors.

For stability, the system uses a 4,096-token context window with efficient KV-cache management and batched inference to avoid fragmentation during long input processing. Long social media posts are truncated to 6,000 characters. We deployed and evaluated three models locally, including our fine-tuned Menta (2.33 GB), Phi-4 Mini (2.40 GB), and Qwen-3 (2.30GB), on an iPhone 15 Pro Max (A17 Pro chip, 8GB RAM) with open-sourced code². The deployment framework also supports other quantization families (Q2_K–Q8_0), allowing trade-offs between latency and accuracy.

4 EXPERIMENTS

4.1 DATASETS AND TASK DEFINITIONS

For our fine-tuning model Menta, we adopted four high-quality corpora for task-specific detection of depression, stress, and suicidal ideation collected from the social media platform Reddit, where disorder-related samples were excluded from the texts and the remaining data were annotated by domain experts. The datasets used in this study include the following:

Dreaddit Turcan & McKeown (2019): From a multi-domain stress corpus, 3,500 segments were manually annotated via Amazon Mechanical Turk as either “stress” or “not stress”, ensuring coverage across diverse domains.

DepSeverity Naseem et al. (2022): This dataset is aimed at accurately identifying users’ depression severity levels, annotated by CLEF eRisk organizers based on DSAS and clinical standards, categorizing the dataset into four severity levels: Minimal, Mild, Moderate, and Severe.

SDCNL Haque et al. (2021): This dataset of 1,895 Reddit posts was collected from communities related to depression and suicidal ideation. Posts from suicide-related communities were labeled as “suicidal ideation,” and posts from depression communities were labeled as “non-suicidal ideation,” subsequently employed unsupervised labeling methods to correct potential misclassifications.

CSSRS-Suicide Gaur et al. (2019): This dataset contains samples related to depressive suicidal ideation and behaviors. Five hundred users were randomly selected, and domain experts annotated the data using the Columbia-Suicide Severity Rating Scale (Brown et al., 2020). The resulting labels fall into five categories: supportive, indicator, ideation, behavior, and attempt.

As shown in Figure 4, the Dreaddit dataset is used for Task 1 (binary stress classification), a subset of 715 posts was used, with 48.4% labeled as stressful with an average post length of 114 ± 42 tokens. We utilized the DepSeverity dataset for Tasks 2 and 3, which supports two post-level binary classification in Task 2 and four-level severity of depression tasks in Task 3). Both Task 2 and 3 have 3,553 posts with an average length of 114 ± 43 tokens. Specifically, the class distribution for Task 3 is: 72.8% Minimal, 8.2% Mild, 11.1% Moderate, and 7.9% Severe.

The SDCNL dataset addresses binary suicide ideation classification (labeled [1/0]) in Task 4 on 379 posts (49.1% flagged), with average text length of 105 ± 12 tokens. The CSSRS-Suicide dataset provides user-level data for binary (Task 5) and five-level (Task 6) suicide risk classification, with 500 users and an average of $1,783 \pm 2,178$ tokens per user. For Task 6, the class distribution is: 21.6% Supportive, 19.8% Indirect, 34.2% Ideation, 15.4% Behavior, and 9.0% Attempt. The levels of severity and classifications for the six mental health tasks are as follows:

- Task 1: Not stressed [0]; Stressed [1].

²<https://anonymous.4open.science/r/Menta-6CAF>

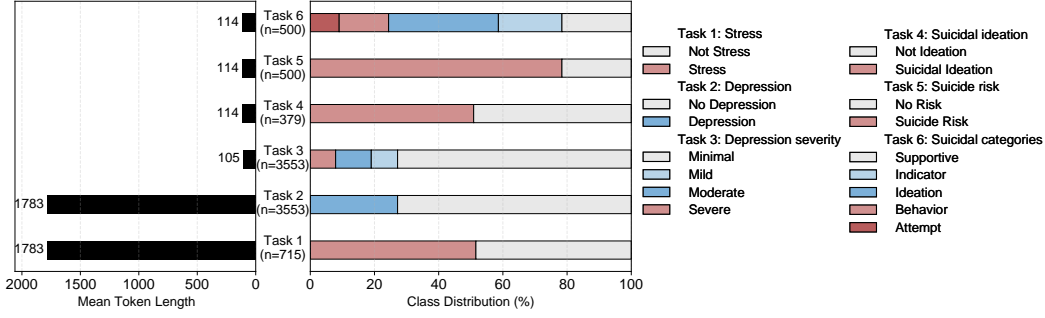


Figure 4: Class and token length distributions across six mental health classification tasks, highlighting label imbalance and diversity in annotation types.

- Task 2: No Depression [0] = Minimal; Depression [1] = Mild + Moderate + Severe.
- Task 3: Minimal [0]; Mild [1]; Moderate [2]; Severe [3].
- Task 4: No suicidal ideation [0]; Suicidal ideation [1].
- Task 5: No indicator of suicide risk [0] = Supportive; Indicator of suicide risk [1] = Indicator + Ideation + Behavior + Attempt.
- Task 6: Supportive [1] = Emotional support but no risk signals; Indicator [2] = Indirect signs of vulnerability or concern; Ideation [3] = Explicit suicidal thoughts without action; Behavior [4] = Suicide-related behaviors short of attempts; Attempt [5] = Evidence of actual suicide attempts.

For data pre-processing, raw datasets were first loaded from task-specific CSV files, after which categorical labels were mapped to standardized numeric formats, with any unmappable entries removed. To preserve class distributions, stratified splitting was applied, using a 72/8/20 ratio for training, validation, and testing.

4.2 BASELINE EXPERIMENTS

For each model and classification task, we conducted experiments under five distinct prompting configurations: zero-shot, one-shot, two-shot, three-shot, and four-shot learning. Each configuration was evaluated independently, and for robustness, we performed five separate runs per setting. This setup enables a systematic comparison of performance across different levels of in-context learning while maintaining consistency across tasks and models.

4.3 SMALL VS. LARGE LANGUAGE MODELS

In this study, we evaluate a diverse set of SLMs limited to under 7B parameters. **Phi-3 (3.8B)** and **Phi-3.5 (3.8B)** are lightweight Microsoft models designed for reasoning and efficiency (Abdin et al., 2024). **LLaMA 3.2 (3B)** is Meta’s compact release optimized for low-resource deployment (Dubey et al., 2024). **Gemma-3 (1B and 4B)** are Google’s instruction-tuned models emphasizing accessibility and safety (Team et al., 2025). **Qwen-2.5 (3B)** and **Qwen-3 (4B)** are Alibaba’s models tailored for multilingual and domain adaptability (Yang et al., 2025). **Phi-4 Mini (3.8B)** extends Microsoft’s efficiency line (Abouelenin et al., 2025). **TinyLLaMA (1.1B)** (Zhang et al., 2024) and **Falcon (1.3B)** (Almazrouei et al., 2023) are streamlined open-source transformer models for edge deployment. **StableLM (3B)** is Stability AI’s open family prioritizing lightweight deployment and transparency (Pinnaparaju et al., 2024).

We developed a multi-task fine-tuned model called Menta employed a shared Qwen3-4B backbone with LoRA adapters, jointly trained on all six tasks. Training was balanced using task-specific sampling weights, and the proposed combined loss function (§3.2) was applied to each task. We compare the developed Menta model using the adjusted loss function with other selected fine-tuned SLMs, including Phi-4 Mini (3.8B), StableLM (3B) and Falcon-1.3B. In addition, we compare the SLMs against larger-scale baselines, Mental-Alpaca and Mental-FLAN-T5 Xu et al. (2024),

developed by Xu et al. through multi-task instruction fine-tuning of Alpaca (7B) Taori et al. (2023) and FLAN-T5 (11B) Chung et al. (2024) using the same datasets and tasks.

4.4 DEPENDANT MODEL VS. GENERAL MODEL

To provide a comprehensive evaluation, we conducted an ablation study comparing single-task and multi-task fine-tuning strategies. In the single-task fine-tuning setting, we trained six independent models, each corresponding to one of the mental health prediction tasks. Each model used Qwen3-4B as the backbone with LoRA adapters applied to enable efficient parameter updates. This setting represents a task-specialized approach, where each model is optimized exclusively for its own objective without considering cross-task transfer. This ablation disentangles the effects of multi-task training from the BACC surrogate loss, clarifying the design choice’s contribution.

4.5 EVALUATION METRICS

We evaluated model performance using accuracy (ACC) as the primary metric for both binary and multi-class classification tasks, selected for interpretability and broad applicability across tasks with differing label granularities and defined as:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

where TP , TN , FP , and FN denote true positives, true negatives, false positives, and false negatives, respectively.

To address label imbalance, we additionally reported balanced accuracy (BACC), which averages true positive rates across classes:

$$\text{BACC} = \frac{1}{C} \sum_{c=1}^C \frac{TP_c}{TP_c + FN_c}$$

where C is the total number of classes, and TP_c , FN_c are true positives and false negatives for class c . This metric is included to enable a fairer comparison among SLMs. Notably, prior work has not evaluated this metric for LLMs, and such evaluation falls outside the scope of this study.

To demonstrate the feasibility, usability, and lightweight nature of the SLMs for real-time, privacy-preserving mental health inference directly on mobile devices, the following metrics are evaluated for the deployment:

- Time to First Token (TTFT, sec): The measurement of how long it takes for the model to generate the first token in response to a prompt (response latency).
- Input Token Processing Speed (ITPS, tokens/sec): An indication of the rate at which input tokens are handled by the model (input throughput).
- Output Token Processing Speed (OTPS, tokens/sec): A representation of how quickly the model generates output tokens after beginning generation (output efficiency).
- Output Elapsed Time (OET, sec): The duration from the start of output generation to the end.
- End-to-End Latency (Total Time): Total elapsed time from prompt submission to final response (overall system efficiency).
- Memory Consumption (RAM, GB): Amount of system memory used by the model during inference.

5 RESULTS

We present the results of the experiments for zero-shot and few-shot learning baselines in §5.1, analyze the outcomes of fine-tuning tasks in §3.2, evaluate deployment considerations in §5.3, and demonstrate specific cases in §5.4.

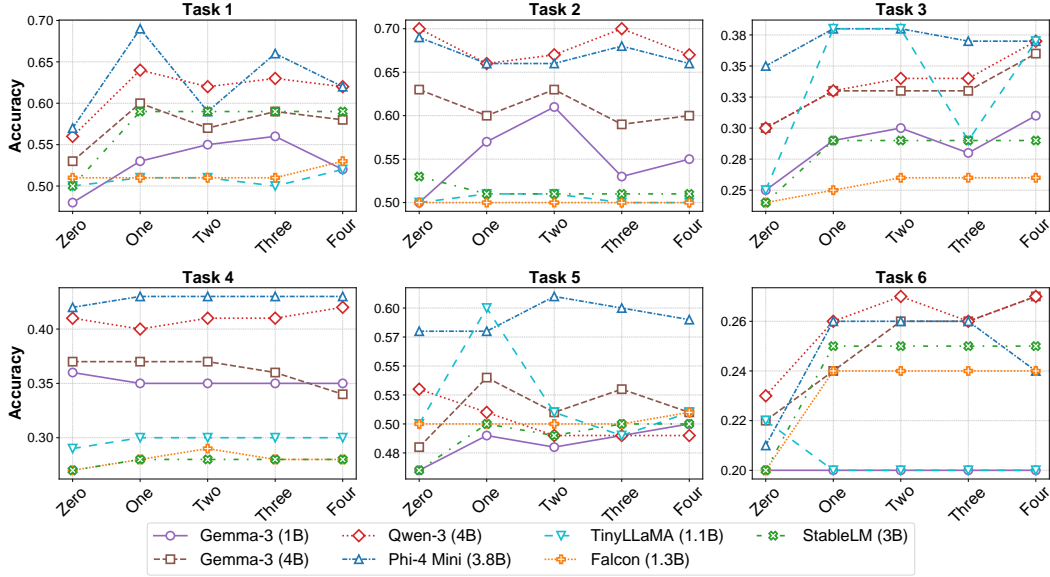


Figure 5: Performance of various SLMs across mental health tasks under zero-shot and few-shot settings. The x-axis indicates the number of shots used.

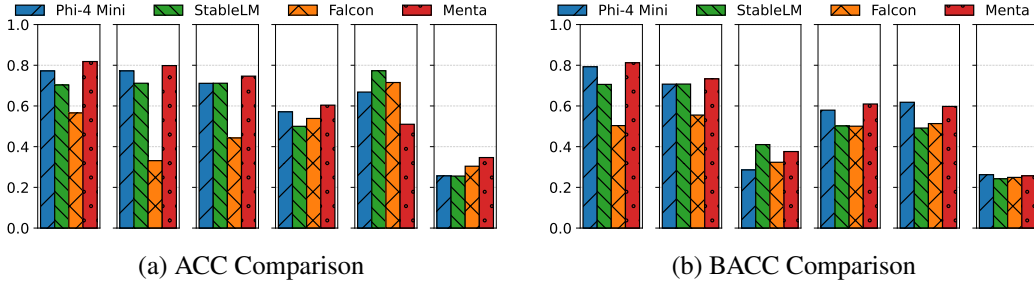


Figure 6: Accuracy (ACC) and Balanced Accuracy (BACC) scores for six tasks (left to right), comparing models Phi-4 Mini (blue), StableLM (green), Falcon (orange), and Menta (red), with Menta consistently achieving the highest performance.

5.1 SLM BASELINES

Evaluating the zero-shot and few-shot learning performance for nine SLMs, our results (Table 2) indicate that most SLMs possess only a basic capability to accurately classify social media posts related to depression, stress, and suicidal ideation without fine-tuning. Two models, Qwen-3 (4B) and Phi-4 Mini (3.8B), demonstrated the strongest performance, achieving average zero-shot accuracies of 45.5% and 47.0% across the six tasks respectively. Among the models, increasing the number of examples generally enhances performance despite exhibiting some variability across the six tasks, with the most notable improvement observed in Qwen-3 (4B) and Phi-4 Mini (3.8B), as shown in Figure 5. Despite modest improvements with few-shot examples, the overall accuracy remains insufficient for practical, real-world applications, which requires the fine-tuning models for further improvement.

Upon examining the model outputs, we observed that SLMs with approximately 1B parameters generally exhibit limited task comprehension and often fail to produce responses in the required format (e.g., binary outputs such as 1 for depression and 0 for no depression). Although their tendency to default to standard responses may result in seemingly acceptable evaluation metrics, these models fundamentally lack the robustness required for reliable real-world deployment.

Table 1: Performance of Menta in comparison with trained LLMs Across Tasks, with the best performance highlighted in bold and the second-best underlined.

Model	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6
Menta (4B) [Ours]	0.82 \pm 0.01	0.80 \pm 0.04	0.75 \pm 0.01	0.60 \pm 0.02	0.51 \pm 0.02	0.35 \pm 0.06
Mental-Alpaca (13B)	0.82 \pm 0.01	<u>0.78</u> \pm 0.01	0.75 \pm 0.01	0.72 \pm 0.00	<u>0.73</u> \pm 0.05	<u>0.40</u> \pm 0.03
Mental-FLAN-T5 (13B)	<u>0.80</u> \pm 0.00	0.76 \pm 0.00	<u>0.76</u> \pm 0.00	<u>0.68</u> \pm 0.01	0.87 \pm 0.01	0.48 \pm 0.01

5.2 FINE-TUNING TASKS

We demonstrate that the general Menta model consistently outperforms both zero-shot and few-shot learning settings, achieving an average improvement of 15.2% across all six tasks over Qwen-3. Notably, it achieves the largest improvements in Task 3, Task 1, Task 6 and Task 2 over the zero-shot setting, with gains of 44.4%, 25.4%, 11.5% and 10.2%, respectively. These results suggest that, compared to zero-shot and few-shot paradigms, fine-tuning with domain-specific mental health data enables SLMs to perform more effectively and reliably on these tasks.

5.2.1 FINE-TUNED SLMs

By fine-tuning various SLMs, our results show that across six mental health prediction tasks, our model Menta achieves the highest average performance (ACC=0.637), outperforming Phi-4 Mini by 1.2%, StableLM by 2.8%, and Falcon by nearly 15.4%, demonstrating its superior accuracy and robustness across diverse mental health tasks, as shown in Figure 6(a). Our results also demonstrate that incorporating the adjusted loss function led to improvements in BACC across all trained SLMs, Menta achieves the highest BACC (0.564) and outperforming Phi-4 Mini by 1.9%, StableLM by 3.1%, and Falcon by nearly 17.0%, as shown in Figure 6(b). These results demonstrate that Menta not only leads in overall accuracy but also offers more reliable performance across diverse task settings.

5.2.2 DEPENDANT MODEL VS. GENERAL MODEL

We further compared single-task fine-tuning Menta against the general Menta trained on cross-datasets and we show that the general Menta model consistently outperformed models trained on individual tasks (Table 3), suggesting robustness across varied mental health classification tasks and no multiple models needed for saving memory. Menta also exhibits lower standard deviations, underscoring its stability and consistency. We refer to the single-task variants as Menta-T1, Menta-T2, etc. As shown in Figure 7, while Menta-T1 and Menta-T2 perform strongly in tasks 1 and 2 respectively, they lack broader consistency across tasks. This ablation study results suggest that although single-task SLMs serve as effective specialists, Menta demonstrates superior versatility as a generalist, making it more practical and impactful for comprehensive mental health monitoring, particularly in mobile or edge deployment scenarios.

To further understand model behavior across tasks, we visualized training and validation loss curves over epochs for each of the six tasks (Figure 14). The plots reveal that most tasks exhibit optimization with a consistent decrease in training loss. However, the validation loss curves are flatter or plateau earlier, suggesting limited generalization. This trend highlights the constrained capacity of SLMs to transfer learned representations across diverse mental health tasks. Notably, the training loss drops most sharply within the first few hundred samples, after which improvements become incremental. This suggests that the model rapidly captures the core patterns of the task early in training, and additional data primarily contributes to fine-tuning and refinement rather than introducing fundamentally new information. The smooth decline in the general Menta model’s training loss, coupled with its balanced task-wise performance as reflected in BACC, indicates that the model achieves stable and generalized learning rather than overfitting to individual tasks.

5.2.3 MENTA VS. LLMs

Comparing the performance of Menta and fine-tuned LLMs in previous studies, Table 1 shows that the Menta model is superior and comparable with fine-tuned LLMs with 81.8%, 79.8% and 75% accuracy for binary stress (Task 1) and binary depression (Task 2) classification and four-level

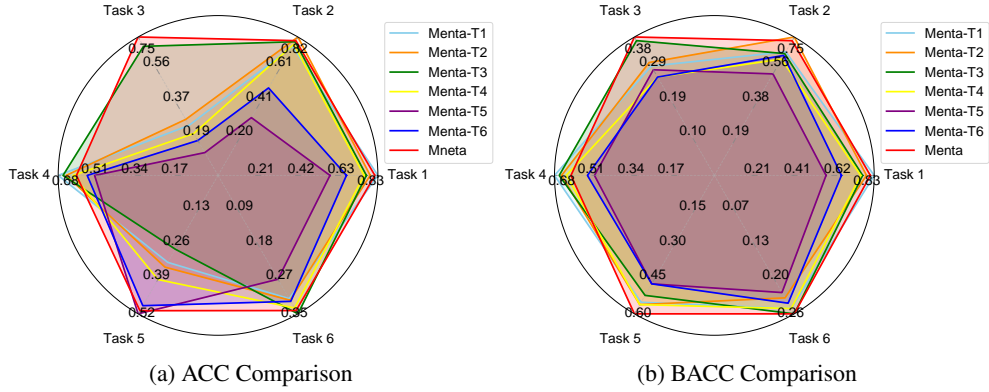


Figure 7: Model performance across six tasks for the individually fine-tuned models Menta-T1 to Menta-T6, compared with the general Menta model optimized using BACC-aware settings.

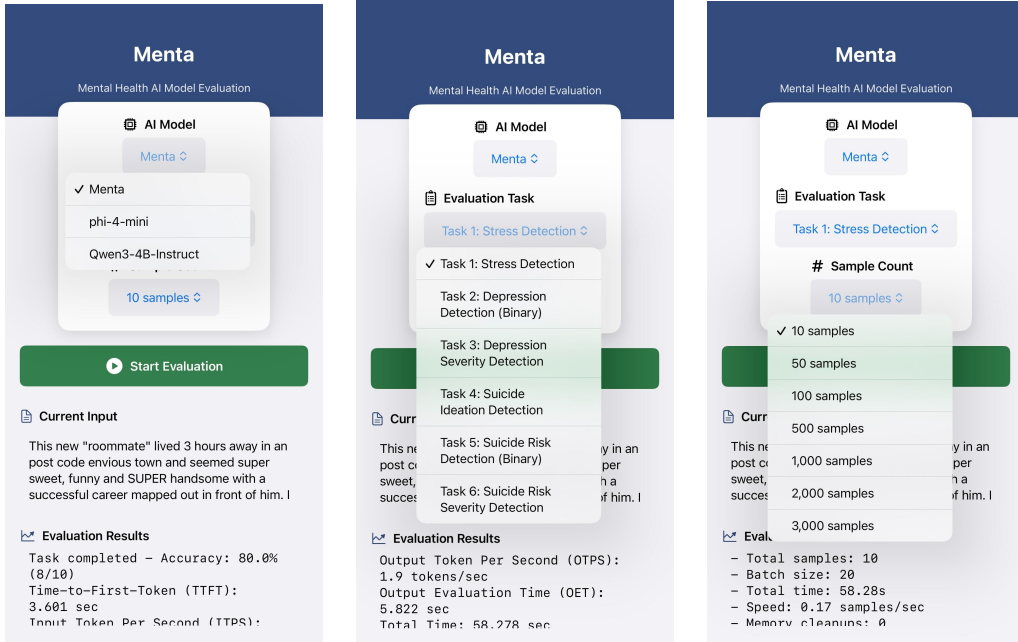


Figure 8: On-device deployment demonstration of Phi-4 Mini, Qwen-3, and Menta across six mental health tasks. The interface shows input social media posts and model predictions for model selection (left), task selection (middle), and sample count selection (right).

depression (Task 3) respectively. However, for the other three tasks, Menta underperforms with an average of 13% and 19% compared to the Mental-Alpaca and Mental-FLAN-T5 model. This gap emphasizes the ability of SLMs is not as good as LLMs for suicidal risk detection tasks.

Nonetheless, we found a trade-off between the size of the models and their performance. Although larger models (13B) achieve higher absolute performance across certain tasks, their performance improvements are not proportional to their parameter increase. When normalized by model size, the 4B Menta model delivers approximately three times higher performance-per-parameter compared to 13B models. Compared to 13B-parameter models, Mental-Alpaca and Mental-FLAN-T5, Menta is approximately 3.25 \times smaller, while delivering comparable or superior performance on Tasks 1 and 2, and maintaining competitive results across other tasks. This indicates that smaller models offer a more favorable trade-off between computational efficiency and accuracy, particularly in deployment-constrained mental health applications.

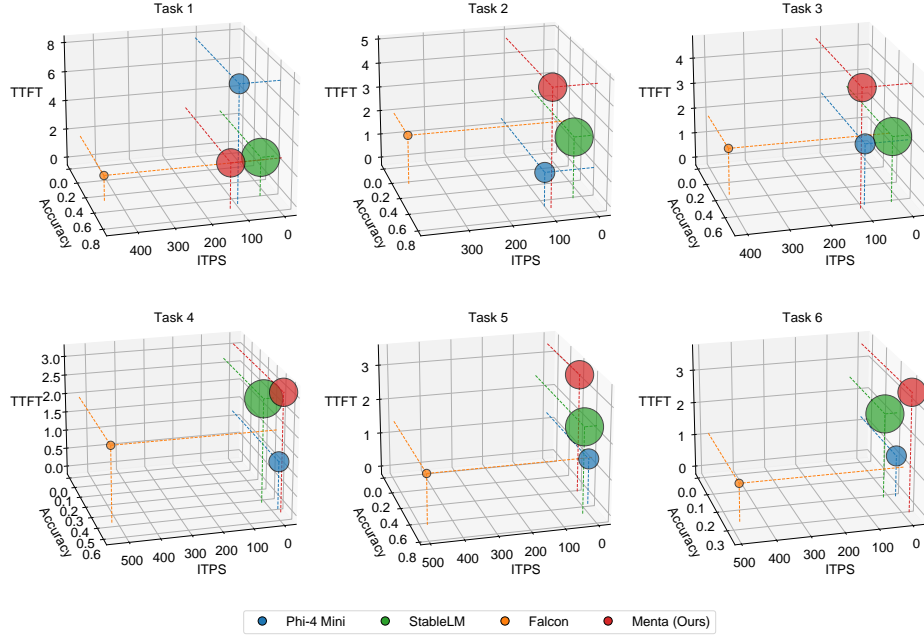


Figure 9: A four-dimensional performance presentation for models Phi-4 Mini, StableLM, Falcon and Menta across ITPS in x axis, Accuracy in y axis, TTFT in z axis, and RAM (bubble size) in the deployment setting on an iPhone 15pro max device.

5.3 DEPLOYMENT

The development of the on-device mental health evaluation system demonstrates the feasibility of running multi-task mental health prediction entirely on-device, without requiring internet connectivity or exposing sensitive user data. Comparing deployment results among the three language models, Phi-4 Mini, Qwen-3, and Menta (Table 4), we observe that Phi-4 Mini consistently shows the lowest TTFT across all tasks, suggesting better responsiveness (low latency) at the start of generation, using the least RAM (2.58–2.90 GB). Meanwhile, Qwen-3 generally achieves the highest ITPS, indicating faster token generation once decoding starts, though its OTPS is lower due to longer overall execution times in Task 5 and Task 6. The fine-tuned Menta offers a balanced performance, slightly slower in TTFT compared to Phi-4 Mini, but often competitive or superior in OTPS, especially in longer tasks (Tasks 5 and 6). Notably, Menta achieves up to 4× higher ITPS compared to existing SLMs such as Phi-4 Mini, while maintaining competitive output quality. It also achieves up to 4× faster total execution time compared to models like StableLM, with comparable memory usage. Additionally, Menta shows moderate RAM consumption (~3.0 GB), slightly higher than Phi-4 Mini but well within practical deployment limits.

As shown in Figure 8, the Menta interface supports multiple on-device classification tasks, including stress, depression (both binary and severity-level), and suicide risk (both binary and categorical). Upon clicking the Start Evaluation button, social media posts will be input directly into the selected model, and the system will display the predicted class labels, correctness, and associated evaluation results. Model performance was assessed across all six tasks, with dataset sizes ranging from 10 to 3,000 posts depending on task complexity and availability.

Furthermore, our analysis reveals a clear trade-off between model efficiency and predictive performance when comparing Menta with other SLMs. From a deployment perspective, Menta offers a highly balanced accuracy, inference speed, and latency, positioning it as an ideal model for real-world applications where both performance and responsiveness are critical.

As illustrated in Figure 9, the three axes represent ITPS, TTFT, and accuracy. Each circle corresponds to a model, and the bubble size reflects memory consumption (RAM). Falcon prioritizes speed with fast TTFT and high ITPS, but at the expense of accuracy, achieving the lowest predictive

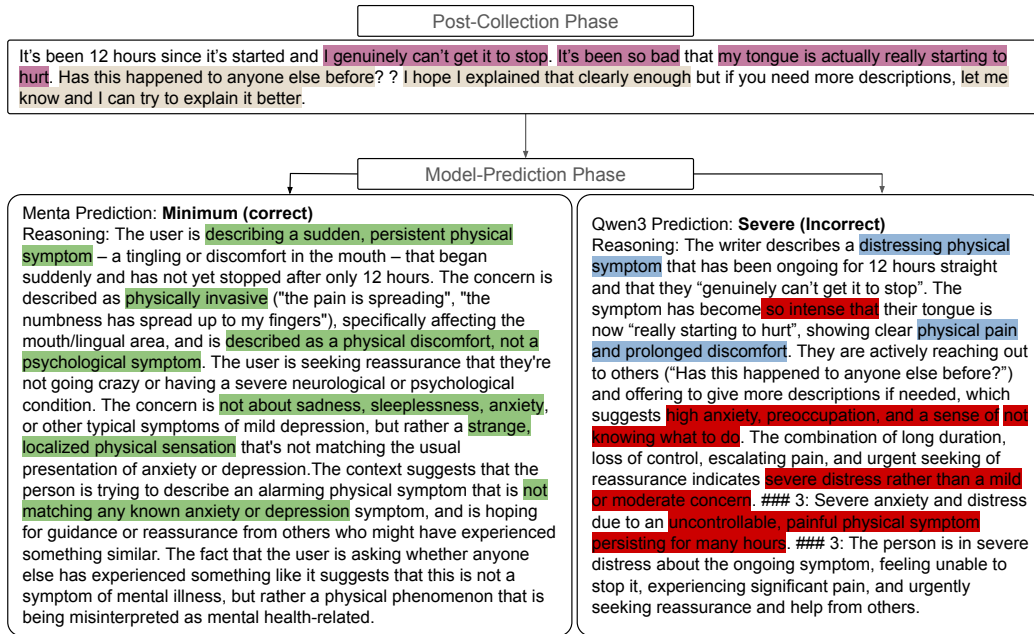


Figure 10: Case 1: A successful prediction by the Menta model on Task 3 (depression ordinal classification), in contrast of Qwen-3’s misclassification of the same post. In the original post, expression of personal distress is highlighted in pink, while help-seeking and uncertainty is highlighted in yellow. In the SLM reasoning, correct analysis is marked in green, correct analytical direction in blue, and overanalysis in red.

performance among all models. Conversely, StableLM attains higher accuracy but incurs substantial latency and resource usage, resulting in increased deployment costs and slower user response times.

In contrast, Menta maintains top-tier accuracy (with the exception of Task 5, where it is narrowly outperformed by StableLM and Phi-4 Mini) while simultaneously achieving relatively lower latency of TTFT and competitive throughput of ITPS. This indicates that Menta generates high-quality outputs efficiently and promptly, which is essential for scalable, user-facing systems. Overall, while all models exhibit trade-offs across evaluation metrics, Menta strategically occupies an optimal deployment-ready position, delivering strong performance with manageable computational requirements, thereby reducing infrastructure overhead without compromising user experience.

5.4 CASE STUDY

In this section, case studies are presented to underscore the key reasons in SLMs for nuanced mental health detection.

Case 1: Depression Level Detection. In this instance, the Reddit post labeled as ‘Minimum’, describes a user’s reflections on an involuntary and persistent action or sensation involving the tongue, as shown in Figure 10. The user expresses their concern and seek help toward a community, and they are willing to provide more details indicating a cooperative and communicative stance despite the discomfort.

The Menta model correctly classifies this post as ‘Minimum’ and offers a rationale consistent with that label. It identifies the user’s measured tone and composed approach to describing a physical issue, recognizing the absence of impulsive or emotionally escalated language as indicative of low psychological distress. This reflects a nuanced understanding that not all discomfort signals depression.

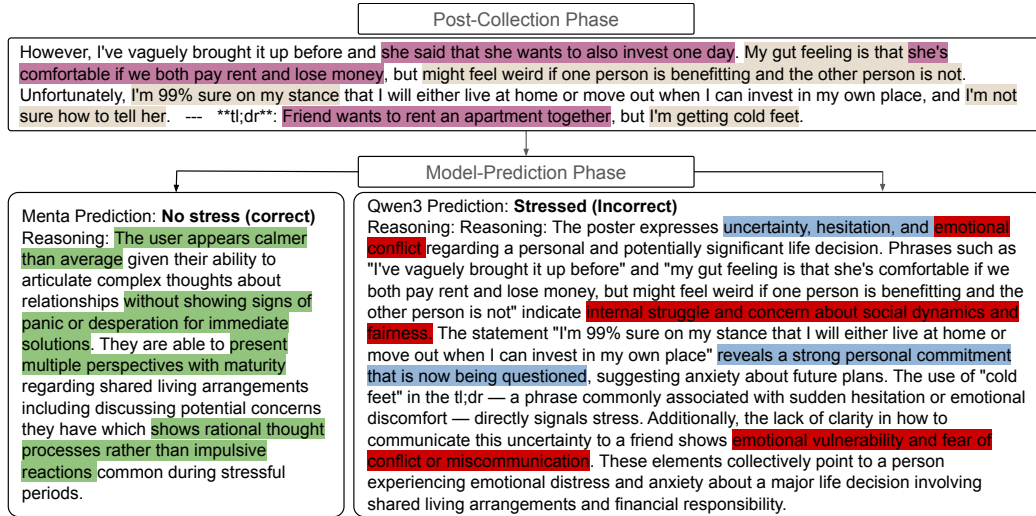


Figure 11: Case 2: A successful prediction by the Menta model on Task 1 (stress classification), in contrast of Qwen-3’s misclassification of the same post. In the original post, cognitive or emotional conflict and decision-making stress are highlighted in pink, while interpersonal or social strain involving external relational dynamics is highlighted in yellow. In the SLM’s reasoning, correct analysis is marked in green, correct analytical direction in blue, and overanalysis in red.

In contrast, Qwen-3 misclassifies the post as ‘Severe’, interpreting the intensity of physical symptoms and expressions of discomfort as proxies for psychological severity. Its rationale conflates acute physical pain with anxiety and psychological distress, ultimately leading to an incorrect prediction. This failure stems from an overreliance on surface-level distress cues without differentiating somatic complaints from depressive features.

With fine-tuning, the Menta model has learned to focus on depression-specific markers and to align predictions with the annotation conventions of the dataset. This case study illustrates Menta’s ability to distinguish between contextual distress and clinical depression, highlighting its strength in predicting mental health states with greater precision.

Case 2: Stress Level Detection. A Reddit post labeled as ‘No stress’ describes a user’s reflections on a housing decision involving a friend. The user expresses thoughts calmly, lays out options rationally, and shows no indication of panic or overwhelming pressure.

The Menta model successfully classifies this post as ‘No stress’ and provides a rationale aligned with this judgment. It highlights the user’s measured tone and maturity in navigating a potentially sensitive decision, correctly interpreting the absence of impulsive or anxious language as indicative of low stress. In contrast, Qwen-3 incorrectly predicts ‘Stressed’, over-indexing on phrases such as ‘cold feet’ and ‘99% sure’, and misinterpreting thoughtful deliberation as internal turmoil. According to the reasoning, Qwen-3 highlights the user’s stress, emotional discomfort, and indecision, which are psychological signals a well-functioning model should consider. However, it emphasizes external relational dynamic, social tension or potential relational mismatch, often linked to conflict-avoidant behavior or suppressed distress (highlighted red in Figure 11) and concludes with an incorrect prediction. Overall, this case study underscores the Menta model’s ability to effectively predict mental health status with training.

Case 3: Suicidal Risk Level Detection. This case highlights a misclassification where a Reddit post with supportive level of suicidal risk was predicted as behavior by the Menta model (Figure 12). The primary cause was the excessive input length, which included a long multi-speaker thread with formatting noise (e.g., ‘>’, list syntax, and HTML artifacts) and placeholder tokens such as ‘Hyperactive behavior’. These tokens acted as spurious cues, biasing the model toward behavior-

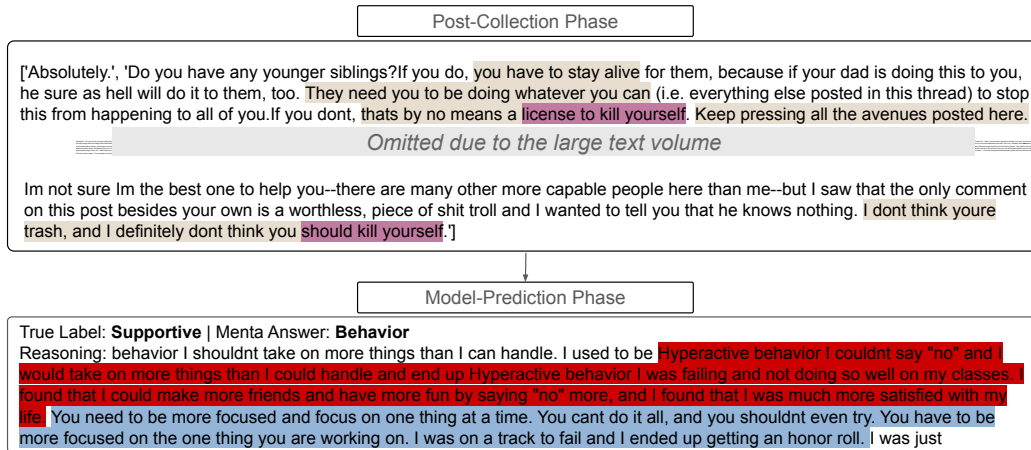


Figure 12: Case 3: A failure case of the Menta model on Task 6 (suicide risk classification), where supportive or encouraging language is highlighted in yellow and misinterpreted phrases in pink in the original post. Repeated “behavior” noise is highlighted in red, while truly supportive content (advice, encouragement, and personal experience) highlighted in blue.

related labels while ignoring clear supportive content such as ‘Please don’t kill yourself’ or ‘Go to the police’.

The model’s training distribution includes both single-speaker inputs and multi-user interactions. In this case, the post aggregates interleaved content from multiple speakers without clear boundaries, increasing contextual complexity and deviating from typical training patterns. This structural mismatch, along with the extended input length, can cause the model to fixate on early salient cues while overlooking the broader supportive intent due to long-context instability.

6 DISCUSSION

With task-specific supervision and systematic fine-tuning, our model Menta achieves strong performance across both binary and multi-level mental health classification tasks, including stress detection, depression severity, and suicidal ideation and risk assessment, capable of handling a diverse range of mental health outcomes within a unified, efficient framework. Our Menta model performs competitively with significantly larger LLMs, especially in depression and stress tasks, where it even surpasses larger models in certain settings. These results demonstrate the practical viability of well-optimized SLMs for continuous mental health monitoring, offering a cost-effective yet high-performing solution for early detection in digital mental health applications. In this section, we address three key questions regarding the use of SLMs and our fine-tuned Menta model in predicting mental health outcomes, while also outlining current limitations and directions for future research.

Can SLMs effectively predict mental health status? Our findings demonstrate that SLMs possess substantial capability in predicting mental health conditions from text-based social media data. Despite their smaller parameter sizes, the well-trained Menta successfully handled complex classification tasks with great performance, which highlights the potential of such models to deliver accurate mental health predictions, supporting their use in low-resource or deployment-constrained settings.

While a previous study (Jia et al., 2025) exploring SLMs in mental health focuses on SLMs in the 2–7B range, our experiments specifically target models strictly smaller than 7B and show that models below about 3.8B parameter scale (e.g., around 1B) fail to meaningfully perform, typically outputting blank or random results in zero- or few-shot settings. This suggests that scaling language models below a certain parameter threshold results in a catastrophic drop in task performance rather than a gradual degradation, indicating the presence of a sharper lower-bound capacity limit. Our work identifies this threshold explicitly within the mental health prediction domain. This observation

is consistent with prior findings (Sengupta et al., 2025; Subramanian et al., 2025), that performance degradation is not always smooth when size is reduced.

Our findings also align with prior work of MentalQLM (0.5B) (Shi et al., 2025), which demonstrated strong performance in predicting a range of mental health outcomes including depression, stress, and suicidal ideation. These results highlight the viability of fine-tuned SLMs for mental health prediction, offering a favorable trade-off between computational efficiency and predictive accuracy. Moreover, consistent with recent studies (Ren et al., 2024; Yang et al., 2024a), we observe that core reasoning and question-answering abilities are preserved after fine-tuning, enabling the resulting models to serve as multi-functional tools with superior performance on mental health tasks.

In summary, our experiments show that SLMs in the 3–7B parameter range offer a strong balance between efficiency and predictive performance with fine-tuning. While smaller models fail to generalize or recognize task constraints, models within this range effectively classify psychological signals in both zero-shot and few-shot settings. With targeted fine-tuning, our model Menta establishes a clear lower bound for model size, achieving strong performance on complex mental health tasks while maintaining computational efficiency.

Balancing model size, accuracy, safety, and real-world impact in mental health AI. Comparing the performance of SLMs and larger LLMs, fine-tuned SLMs show comparable or superior performance in depression and stress level classification tasks while lower performance in suicidal ideation prediction tasks. However, compared with non-trained SLMs, the performance improves in a large degree as well as retaining a more balanced accuracy among prediction tasks and classes. Our findings show that that well-optimized lightweight SLMs can retain core reasoning and contextual understanding necessary for mental health inference.

Our findings were observed in previous work (Jia et al., 2025) that SLMs often came within 2% of LLMs’ F1 scores in binary classification settings under zero-shot evaluation, highlighting the potential of SLMs as resource-efficient and privacy-conscious alternatives in clinical or low-resource environments. Other prior works have concentrated primarily on LLMs or hybrid models in mental health domains. For instance, Kim et al. (Kim et al., 2025) compares zero-shot LLMs and supervised classifiers built on LLM embeddings, showing that while LLMs generalize well in binary depression classification, they struggle in fine-grained severity classification tasks, where our findings collaborate the difficulty in finer labels. Recent surveys (Ge et al., 2025; Garg et al., 2025) highlight the growing promise of SLMs in digital health applications, noting their efficiency, reduced computational demands, and capacity for local deployment, which offer significant advantages over LLMs by being more privacy-preserving and easier to integrate into real-world clinical or mobile settings. Our work complements these surveys by providing empirical evidence for the operational boundary of this trade-off, showing precisely when SLMs transition from viable to non-viable performance levels for such sensitive mental health applications. Also, our compact Menta model demonstrates that the on-device SLM deployment is essential for protecting user privacy, a primary ethical concern in digital psychiatry (Patel et al., 2018; Jin et al., 2025). In resource-limited clinical settings, particularly in low- and middle-income countries where mental-health workforce shortages persist (Moitra et al., 2022), SLMs enable accessible, cost-efficient, and scalable mental-health screening tools. Moreover, smaller models are easier to interpret and audit, aligning with recent calls for transparent and accountable AI in healthcare (Doshi-Velez & Kim, 2017; Tonekaboni et al., 2019).

The competitiveness of SLMs in this domain primarily emerges because mental-health language tasks often rely on localized lexical, affective, and syntactic cues, such as expressions of self-reference, emotional polarity, and temporal framing, which can be effectively captured by moderate-sized models without requiring the broad world knowledge encoded in LLMs. As shown in prior work (Ji et al., 2021; Shen et al., 2018; Benton et al., 2017), fine-tuning on domain-specific corpora substantially narrows the performance gap between smaller and larger architectures. Our findings extend this evidence by demonstrating that parameter-efficient fine-tuning (e.g., LoRA) combined with adaptive, weighted loss functions enables SLMs to learn discriminative representations for mental-health prediction at a fraction of the computational cost.

While LLMs generally outperform across a broader range of tasks due to their larger capacity and instruction-following capabilities, SLMs demonstrate competitive and even superior performance on specific tasks, particularly depression and stress prediction tasks. For the remaining tasks, although SLMs show notable improvements, their performance still underperforms compared with

LLMs (Xu et al., 2024). This suggests that SLMs can be highly effective when the task structure and linguistic signals are well-aligned with their capacity. However, a critical limitation remains: SLMs struggle with long input sequences due to both token-length constraints and reduced capability in handling extended contextual dependencies (Kumar, 2025). This restricts their applicability in scenarios requiring nuanced understanding of lengthy or complex user narratives, where LLMs maintain a distinct advantage (Lu et al., 2024).

Overall, our findings support a paradigm shift from size-centric scaling to purpose-centric optimization in mental-health AI. As efficiency, transparency, safety and user privacy (Zhang et al., 2021) become increasingly prioritized in clinical machine learning, the balance between performance and practicality will define the next generation of digital mental-health tools. SLMs address emerging demands in digital mental health for in-situ, low-latency, and privacy-preserving early detection (Harari et al., 2016; Wang et al., 2018), while also aligning with growing calls for transparent and accountable AI in healthcare settings (Wies et al., 2021; Smith et al., 2023). By showing that fine-tuned SLMs with adaptive, cross-dataset strategies can match or surpass LLMs in mental-health outcome prediction, this study provides empirical evidence that smaller, well-aligned models may offer a more sustainable and ethically sound path forward for real-world mental-health monitoring.

What tasks are most suitable for SLMs? In our work, we demonstrate that cross-dataset training is not only feasible but beneficial for predicting depression, stress, and suicidal ideation, when combined with an adaptive loss design. This finding resonates with observations from Yao et al. (Benton et al., 2017), who reported that shared latent representations across related affective tasks can promote positive transfer. However, unlike their multi-task architecture, our approach achieves comparable cross-task gains within a parameter-efficient fine-tuning paradigm (LoRA), thereby reducing computational cost and memory footprint. This distinction highlights that even lightweight fine-tuning frameworks can leverage dataset diversity effectively, provided the optimization objective is properly calibrated.

Also, we found that the fine-tuned model Menta shows better performance in depression and stress detection tasks compared with suicidal risk detection. We found that fine-tuned and non-fine-tuned SLMs, and even larger language models, tend to achieve stronger performance on depression and stress detection tasks compared to suicide risk classification. For example, a systematic review (Bauer et al., 2024) further notes that language models often misclassify nuanced expressions of suicidality due to class imbalance and the reliance on contextual or temporal dependencies, which are not captured in single posts. Our findings also collaborate with another study (Lamichhane, 2023), which reported F1-scores of 0.86 for depression detection and only 0.37 for suicidality detection on social media data. These findings reflect that depression and stress often manifest with clearer linguistic patterns and more abundant data, whereas suicidality may be rarer, more context dependent, and linguistically subtler, making it harder for models to learn reliably. As a consequence, SLMs (and even larger models) may be better suited for screening tasks (e.g., detecting depression or stress) than for nuanced stratification of suicide risk without additional multimodal or temporal data.

Limitation and future directions. Despite demonstrating that SLMs such as Qwen-3 (4B) can effectively perform mental health prediction tasks, several limitations constrain the generalizability and interpretability of our findings.

While our evaluation used diverse datasets covering depression, stress, and suicidal ideation, all inputs were textual social media posts. This focus excludes multimodal features such as linguistic style dynamics, user metadata, or interaction patterns, that often play a critical role in understanding real-world mental health contexts. Incorporating multimodal signals could improve robustness and ecological validity in future work (Yazdavar et al., 2020; Khoo et al., 2024). Also, one key limitation of this study lies in the restricted diversity of datasets and mental health tasks explored. Our evaluation primarily focused on three task categories depression, stress, and suicidal ideation, each derived from publicly available social media datasets. While these benchmarks are widely adopted in affective computing, they still represent a narrow subset of the psychological spectrum. Real-world mental health expressions often encompass more nuanced and overlapping conditions such as anxiety, bipolar disorder, and emotional dysregulation, which are underrepresented in our current corpus (Cao et al., 2024). Additionally, all analyses were conducted on publicly available, anonymized datasets rather than in clinical environments. While this design ensures privacy and reproducibility, it limits conclusions about clinical safety, ethical compliance, and longitudinal reliability (Seyed-

salehi & Fazel, 2024). Future directions include the test of federated or on-device fine-tuned SLMs in clinical settings. Deploying such models as assistive rather than diagnostic tools, under clinician supervision, could provide a safe bridge between algorithmic insights and responsible mental health support.

7 CONCLUSION

In conclusion, this work demonstrates the feasibility and effectiveness of SLMs for digital mental health prediction from social media text. By introducing Menta, the first optimized and compact SLM fine-tuned across six mental health tasks, we show that careful architectural and training design enables high performance without sacrificing efficiency or deployability. Menta not only achieves strong predictive accuracy, especially for binary stress and suicidality detection, but also operates in real time on resource-constrained devices such as smartphones. These findings highlight that SLMs can offer scalable, private, and interpretable mental health support, bridging the gap between clinical relevance and technological accessibility.

These findings addresses critical challenges of cost, accessibility, and privacy that limit traditional approaches. While our results provide encouraging evidence, further work should evaluate multimodal and cross-lingual extensions, clinical validation, and interpretability to ensure responsible translation into practice. Overall, Menta demonstrates that small, task-optimized language models can deliver accurate, efficient, and interpretable mental health predictions, while enabling scalable and privacy-preserving real-world deployment and shows that adaptive fine-tuning can bridge the gap between efficiency and performance in next-generation digital mental-health systems.

REFERENCES

- Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J Hewett, Mojan Javaheripi, Piero Kauffmann, et al. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024.
- Abdelrahman Abouelenin, Atabak Ashfaq, Adam Atkinson, Hany Awadalla, Nguyen Bach, Jianmin Bao, Alon Benhaim, Martin Cai, Vishrav Chaudhary, Congcong Chen, et al. Phi-4-mini technical report: Compact yet powerful multimodal language models via mixture-of-loras. *arXiv preprint arXiv:2503.01743*, 2025.
- Ankit Aich, Avery Quynh, Pamela Osseyi, Amy Pinkham, Philip Harvey, Brenda Curtis, Colin Depp, and Natalie Parde. Using llms to aid annotation and collection of clinically-enriched data in bipolar disorder and schizophrenia. *arXiv preprint arXiv:2406.12687*, 2024.
- Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Mérouane Debbah, Étienne Goffinet, Daniel Hesslow, Julien Launay, Quentin Malartic, et al. The falcon series of open language models. *arXiv preprint arXiv:2311.16867*, 2023.
- Patricia A Areán, Kien Hoa Ly, and Gerhard Andersson. Mobile technology for mental health assessment. *Dialogues in clinical neuroscience*, 18(2):163–169, 2016.
- Thushari Atapattu, Mahen Herath, Charitha Elvitigala, Piyanjali de Zoysa, Kasun Gunawardana, Menasha Thilakaratne, Kasun De Zoysa, and Katrina Falkner. Emoment: An emotion annotated mental health corpus from two south asian countries. *arXiv preprint arXiv:2208.08486*, 2022.
- Brian Bauer, Raquel Norel, Alex Leow, Zad Abi Rached, Bo Wen, and Guillermo Cecchi. Using large language models to understand suicidality in a social media-based taxonomy of mental health disorders: Linguistic analysis of reddit posts. *JMIR mental health*, 11:e57234, 2024.
- Peter Belcak, Greg Heinrich, Shizhe Diao, Yonggan Fu, Xin Dong, Saurav Muralidharan, Yingyan Celine Lin, and Pavlo Molchanov. Small language models are the future of agentic ai. *arXiv preprint arXiv:2506.02153*, 2025.
- Loris Belcastro, Riccardo Cantini, Fabrizio Marozzo, Domenico Talia, and Paolo Trunfio. Detecting mental disorder on social media: a chatgpt-augmented explainable approach. *Online Social Networks and Media*, 48:100321, 2025.

- Adrian Benton, Margaret Mitchell, and Dirk Hovy. Multi-task learning for mental health using social media text. *arXiv preprint arXiv:1712.03538*, 2017.
- Prajval Bolegave and Pushpak Bhattacharya. A gold standard dataset and evaluation framework for depression detection and explanation in social media using llms. *arXiv preprint arXiv:2507.19899*, 2025.
- Diksha Brahmabhatt and William L Schpero. Access to psychiatric appointments for medicaid enrollees in 4 large us cities. *JAMA*, 332(8):668–669, 2024.
- Lily A Brown, Edwin D Boudreaux, Sarah A Arias, Ivan W Miller, Alexis M May, Carlos A Camargo Jr, Craig J Bryan, and Michael F Arney. C-ssrs performance in emergency department patients at high risk for suicide. *Suicide and Life-Threatening Behavior*, 50(6):1097–1104, 2020.
- John Bunyi, Kathryn E Ringland, and Stephen M Schueller. Accessibility and digital mental health: considerations for more accessible and equitable mental health apps. *Frontiers in digital health*, 3:742196, 2021.
- Yuchen Cao, Jianglai Dai, Zhongyan Wang, Yeyubei Zhang, Xiaorui Shen, Yunchong Liu, and Yexin Tian. Machine learning approaches for mental illness detection on social media: A systematic review of biases and methodological challenges. *arXiv preprint arXiv:2410.16204*, 2024.
- Stevie Chancellor and Munmun De Choudhury. Methods in predictive techniques for mental health status on social media: a critical review. *NPJ digital medicine*, 3(1):43, 2020.
- Xiangyong Chen and Xiaochuan Lin. Generating medically-informed explanations for depression detection using llms. *arXiv preprint arXiv:2503.14671*, 2025.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.
- Sarah Clement, Oliver Schauman, Tanya Graham, Francesca Maggioni, Sara Evans-Lacko, Nikita Bezborodovs, Craig Morgan, Nicolas Rüsch, June SL Brown, and Graham Thornicroft. What is the impact of mental health-related stigma on help-seeking? a systematic review of quantitative and qualitative studies. *Psychological medicine*, 45(1):11–27, 2015.
- Glen Coppersmith, Ryan Leary, Patrick Crutchley, and Alex Fine. Natural language processing of social media as screening for suicide risk. *Biomedical informatics insights*, 10:1178222618792860, 2018.
- Munmun De Choudhury, Michael Gamon, Scott Counts, and Eric Horvitz. Predicting depression via social media. In *Proceedings of the international AAAI conference on web and social media*, volume 7, pp. 128–137, 2013.
- Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.
- 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 e-prints*, pp. arXiv–2407, 2024.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- Muskan Garg, Shaina Raza, Shebuti Rayana, Xingyi Liu, and Sunghwan Sohn. The rise of small language models in healthcare: A comprehensive survey. *arXiv preprint arXiv:2504.17119*, 2025.
- 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*, pp. 514–525, 2019.

- Zhuohan Ge, Nicole Hu, Darian Li, Yubo Wang, Shihao Qi, Yuming Xu, Han Shi, and Jason Zhang. A survey of large language models in mental health disorder detection on social media. In *2025 IEEE 41st International Conference on Data Engineering Workshops (ICDEW)*, pp. 164–176. IEEE, 2025.
- Lovedeep Gondara, Jonathan Simkin, Graham Sayle, Shebnum Devji, Gregory Arbour, and Raymond Ng. Small or large? zero-shot or finetuned? guiding language model choice for specialized applications in healthcare. *arXiv preprint arXiv:2504.21191*, 2025.
- Renee D Goodwin, Lisa C Dierker, Melody Wu, Sandro Galea, Christina W Hoven, and Andrea H Weinberger. Trends in us depression prevalence from 2015 to 2020: the widening treatment gap. *American Journal of Preventive Medicine*, 63(5):726–733, 2022.
- Jessica L Hamilton, Melissa J Dreier, Bianca Caproni, Jennifer Fedor, Krina C Durica, and Carissa A Low. Improving the science of adolescent social media and mental health: Challenges and opportunities of smartphone-based mobile sensing and digital phenotyping. *Journal of Technology in Behavioral Science*, 10(2):301–319, 2025.
- Ayaan Haque, Viraaj Reddi, and Tyler Giallanza. Deep learning for suicide and depression identification with unsupervised label correction. In *Artificial Neural Networks and Machine Learning—ICANN 2021: 30th International Conference on Artificial Neural Networks, Bratislava, Slovakia, September 14–17, 2021, Proceedings, Part V 30*, pp. 436–447. Springer, 2021.
- Gabriella M Harari, Nicholas D Lane, Rui Wang, Benjamin S Crosier, Andrew T Campbell, and Samuel D Gosling. Using smartphones to collect behavioral data in psychological science: Opportunities, practical considerations, and challenges. *Perspectives on Psychological Science*, 11(6):838–854, 2016.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuezhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.
- Muhammad Hussain, Caikou Chen, Muzammil Hussain, Muhammad Anwar, Mohammed Abaker, Abdelzahir Abdelmaboud, and Iqra Yamin. Optimised knowledge distillation for efficient social media emotion recognition using distilbert and albert. *Scientific Reports*, 15(1):30104, 2025.
- Shaoxiong Ji, Celina Ping Yu, Sai-fu Fung, Shirui Pan, and Guodong Long. Supervised learning for suicidal ideation detection in online user content. *Complexity*, 2018(1):6157249, 2018.
- Shaoxiong Ji, Tianlin Zhang, Luna Ansari, Jie Fu, Prayag Tiwari, and Erik Cambria. Mentalbert: Publicly available pretrained language models for mental healthcare. *arXiv preprint arXiv:2110.15621*, 2021.
- Hong Jia, Shiya Fu, Feng Xia, Vassilis Kostakos, and Ting Dang. Beyond scale: Small language models are comparable to gpt-4 in mental health understanding. *arXiv preprint arXiv:2507.08031*, 2025.
- Zheng Ping Jiang, Sarah Ita Levitan, Jonathan Zomick, and Julia Hirschberg. Detection of mental health from reddit via deep contextualized representations. In *Proceedings of the 11th international workshop on health text mining and information analysis*, pp. 147–156, 2020.
- Yu Jin, Jiayi Liu, Pan Li, Baosen Wang, Yangxinyu Yan, Huilin Zhang, Chenhao Ni, Jing Wang, Yi Li, Yajun Bu, et al. The applications of large language models in mental health: Scoping review. *Journal of Medical Internet Research*, 27:e69284, 2025.
- Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang. Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks. *Advances in Neural Information Processing Systems*, 36:48573–48602, 2023.
- Lin Sze Khoo, Mei Kuan Lim, Chun Yong Chong, and Roisin McNaney. Machine learning for multimodal mental health detection: a systematic review of passive sensing approaches. *Sensors*, 24(2):348, 2024.

- Dae-young Kim, Rebecca Hwa, and Muhammad Mahbubur Rahman. mhgpt: A lightweight generative pre-trained transformer for mental health text analysis. *arXiv preprint arXiv:2408.08261*, 2024.
- Jina Kim, Jieon Lee, Eunil Park, and Jinyoung Han. A deep learning model for detecting mental illness from user content on social media. *Scientific reports*, 10(1):11846, 2020.
- Samuel Kim, Oghenemaro Imieye, and Yunting Yin. Interpretable depression detection from social media text using llm-derived embeddings. *arXiv preprint arXiv:2506.06616*, 2025.
- Akshi Kumar. From large to small: The rise of small language models (slms) in text analytics. 2025.
- Bishal Lamichhane. Evaluation of chatgpt for nlp-based mental health applications. *arXiv preprint arXiv:2303.15727*, 2023.
- Jooyoung Lee, Fan Yang, Thanh Tran, Qian Hu, Emre Barut, Kai-Wei Chang, and Chengwei Su. Can small language models help large language models reason better?: Lm-guided chain-of-thought. *arXiv preprint arXiv:2404.03414*, 2024.
- Pierre Lepagnol, Thomas Gerald, Sahar Ghannay, Christophe Servan, and Sophie Rosset. Small language models are good too: An empirical study of zero-shot classification. *arXiv preprint arXiv:2404.11122*, 2024.
- Zhenyan Lu, Xiang Li, Dongqi Cai, Rongjie Yi, Fangming Liu, Xiwen Zhang, Nicholas D Lane, and Mengwei Xu. Small language models: Survey, measurements, and insights. *arXiv preprint arXiv:2409.15790*, 2024.
- Hengyu Luo, Peng Liu, and Stefan Esping. Exploring small language models with prompt-learning paradigm for efficient domain-specific text classification. *arXiv preprint arXiv:2309.14779*, 2023.
- Aishik Mandal, Tanmoy Chakraborty, and Iryna Gurevych. Towards privacy-aware mental health ai models: Advances, challenges, and opportunities. *arXiv preprint arXiv:2502.00451*, 2025.
- Ritesh Maurya, Nikhil Rajput, MG Diviit, Satyajit Mahapatra, and Manish Kumar Ojha. Exploring the potential of lightweight large language models for ai-based mental health counselling task: a novel comparative study. *Scientific Reports*, 15(1):22463, 2025.
- Kristina McMahan, Karli M Martin, Melissa J Greenfield, Pamela Hay, Madison Bates Redwine, Rachel Fargason, and Kristine Lokken. Using a tele-behavioral health rapid intake model to address high demand for psychotherapy at an academic medical center during covid-19. *Frontiers in Psychiatry*, 13:989838, 2022.
- Modhurima Moitra, Damian Santomauro, Pamela Y Collins, Theo Vos, Harvey Whiteford, Shekhar Saxena, and Alize J Ferrari. The global gap in treatment coverage for major depressive disorder in 84 countries from 2000–2019: a systematic review and bayesian meta-regression analysis. *PLoS medicine*, 19(2):e1003901, 2022.
- 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*, pp. 2563–2572, 2022.
- Yang Ni and Fanli Jia. A scoping review of ai-driven digital interventions in mental health care: mapping applications across screening, support, monitoring, prevention, and clinical education. In *Healthcare*, volume 13, pp. 1205. MDPI, 2025.
- World Health Organization. Who special initiative for mental health (simh). <https://www.who.int/initiatives/who-special-initiative-for-mental-health>, 2025. Accessed: YYYY-MM-DD.
- World Health Organization et al. Depression and other common mental disorders: global health estimates. In *Depression and other common mental disorders: global health estimates*. 2017.
- Vikram Patel, Shekhar Saxena, Crick Lund, Graham Thornicroft, Florence Baingana, Paul Bolton, Dan Chisholm, Pamela Y Collins, Janice L Cooper, Julian Eaton, et al. The lancet commission on global mental health and sustainable development. *The lancet*, 392(10157):1553–1598, 2018.

- Branislav Pecher, Ivan Srba, and Maria Bielikova. Comparing specialised small and general large language models on text classification: 100 labelled samples to achieve break-even performance. *arXiv preprint arXiv:2402.12819*, 2024.
- Thang M Pham, Phat T Nguyen, Seunghyun Yoon, Viet Dac Lai, Franck Dernoncourt, and Trung Bui. Slimlm: An efficient small language model for on-device document assistance. *arXiv preprint arXiv:2411.09944*, 2024.
- Xiangdan Piao, Jun Xie, and Shunsuke Managi. Continuous worsening of population emotional stress globally: universality and variations. *BMC Public Health*, 24(1):3576, 2024.
- Nikhil Pinnaparaju, Reshith Adithyan, Duy Phung, Jonathan Tow, James Baicoianu, Ashish Datta, Maksym Zhuravinskyi, Dakota Mahan, Marco Bellagente, Carlos Riquelme, et al. Stable code technical report. *arXiv preprint arXiv:2404.01226*, 2024.
- Xuan Ren, Biao Wu, and Lingqiao Liu. I learn better if you speak my language: Understanding the superior performance of fine-tuning large language models with llm-generated responses. *arXiv preprint arXiv:2402.11192*, 2024.
- Shailik Sarkar, Abdulaziz Alhamadani, Lulwah Alkulaib, and Chang-Tien Lu. Predicting depression and anxiety on reddit: a multi-task learning approach. In *2022 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pp. 427–435. IEEE, 2022.
- SM Sarwar. Fedmentalcare: towards privacy-preserving fine-tuned llms to analyze mental health status using federated learning framework. *arXiv preprint arXiv:2503.05786*, 2025.
- Ayan Sengupta, Siddhant Chaudhary, and Tanmoy Chakraborty. Compression laws for large language models. *arXiv preprint arXiv:2504.04342*, 2025.
- Aida Seyedsalehi and Seena Fazel. Suicide risk assessment tools and prediction models: new evidence, methodological innovations, outdated criticisms. *BMJ mental health*, 27(1), 2024.
- Shahid Munir Shah, Syeda Anshrah Gillani, Mirza Samad Ahmed Baig, Muhammad Aamer Saleem, and Muhammad Hamzah Siddiqui. Advancing depression detection on social media platforms through fine-tuned large language models. *Online Social Networks and Media*, 46:100311, 2025.
- Tiancheng Shen, Jia Jia, Guangyao Shen, Fuli Feng, Xiangnan He, Huanbo Luan, Jie Tang, Thanassis Tiropanis, Tat Seng Chua, and Wendy Hall. Cross-domain depression detection via harvesting social media. *International Joint Conferences on Artificial Intelligence*, 2018.
- Jiayu Shi, Zexiao Wang, Jiandong Zhou, Chengyu Liu, Poly ZH Sun, Erying Zhao, and Lei Lu. Mentalqlm: A lightweight large language model for mental healthcare based on instruction tuning and dual lora modules. *IEEE Journal of Biomedical and Health Informatics*, 2025.
- Jaemin Shin, Hyungjun Yoon, Seungjoo Lee, Sungjoon Park, Yunxin Liu, Jinho D Choi, and Sung-Ju Lee. Fedtherapist: Mental health monitoring with user-generated linguistic expressions on smartphones via federated learning. *arXiv preprint arXiv:2310.16538*, 2023.
- Katharine A Smith, Charlotte Blease, Maria Faurholt-Jepsen, Joseph Firth, Tom Van Daele, Carmen Moreno, Per Carlbring, Ulrich W Ebner-Priemer, Nikolaos Koutsouleris, Heleen Riper, et al. Digital mental health: challenges and next steps. *BMJ Ment Health*, 26(1), 2023.
- Dan J Stein, Steven J Shoptaw, Daniel V Vigo, Crick Lund, Pim Cuijpers, Jason Bantjes, Norman Sartorius, and Mario Maj. Psychiatric diagnosis and treatment in the 21st century: paradigm shifts versus incremental integration. *World Psychiatry*, 21(3):393–414, 2022.
- Shreyas Subramanian, Vikram Elango, and Mecit Gungor. Small language models (slms) can still pack a punch: A survey. *arXiv preprint arXiv:2501.05465*, 2025.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B Hashimoto. Stanford alpaca: An instruction-following llama model, 2023.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025.

- Qwen Team. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2024.
- Sana Tonekaboni, Shalmali Joshi, Melissa D McCradden, and Anna Goldenberg. What clinicians want: contextualizing explainable machine learning for clinical end use. In *Machine learning for healthcare conference*, pp. 359–380. PMLR, 2019.
- Elsbeth Turcan and Kathleen McKeown. Dreddit: A reddit dataset for stress analysis in social media. *arXiv preprint arXiv:1911.00133*, 2019.
- University of Maine. Social media statistics details. <https://umaine.edu/undiscoveredmaine/small-business/resources/marketing-for-small-business/social-media-tools/social-media-statistics-details/>, 2023. Accessed: 2025-09-08.
- Chien Van Nguyen, Xuan Shen, Ryan Aponte, Yu Xia, Samyadeep Basu, Zhengmian Hu, Jian Chen, Mihir Parmar, Sasidhar Kunapuli, Joe Barrow, et al. A survey of small language models. *arXiv preprint arXiv:2410.20011*, 2024.
- Fali Wang, Zhiwei Zhang, Xianren Zhang, Zongyu Wu, Tzuhao Mo, Qiuhao Lu, Wanjing Wang, Rui Li, Junjie Xu, Xianfeng Tang, et al. A comprehensive survey of small language models in the era of large language models: Techniques, enhancements, applications, collaboration with llms, and trustworthiness. *arXiv preprint arXiv:2411.03350*, 2024.
- Fali Wang, Minhua Lin, Yao Ma, Hui Liu, Qi He, Xianfeng Tang, Jiliang Tang, Jian Pei, and Suhang Wang. A survey on small language models in the era of large language models: Architecture, capabilities, and trustworthiness. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pp. 6173–6183, 2025a.
- Linyong Wang, Lianwei Wu, Shaoqi Song, Yaxiong Wang, Cuiyun Gao, and Kang Wang. Distilling structured rationale from large language models to small language models for abstractive summarization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 25389–25397, 2025b.
- Rui Wang, Weichen Wang, Alex DaSilva, Jeremy F Huckins, William M Kelley, Todd F Heatherton, and Andrew T Campbell. Tracking depression dynamics in college students using mobile phone and wearable sensing. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2(1):1–26, 2018.
- Xiaofei Wang, Hayley M Sanders, Yuchen Liu, Kennarey Seang, Bach Xuan Tran, Atanas G Atanasov, Yue Qiu, Shenglan Tang, Josip Car, Ya Xing Wang, et al. Chatgpt: promise and challenges for deployment in low-and middle-income countries. *The Lancet Regional Health–Western Pacific*, 41, 2023.
- Xin Wang, Ting Dang, Xinyu Zhang, Vassilis Kostakos, Michael J Witbrock, and Hong Jia. Healthslm-bench: Benchmarking small language models for mobile and wearable healthcare monitoring. *arXiv preprint arXiv:2509.07260*, 2025c.
- Nicole Davis Weaver, Gregory J Bertolacci, Emily Rosenblad, Sama Ghoba, Matthew Cunningham, Kevin S Ikuta, Madeline E Moberg, Vincent Mougin, Chieh Han, Eve E Wool, et al. Global, regional, and national burden of suicide, 1990–2021: a systematic analysis for the global burden of disease study 2021. *The Lancet Public Health*, 10(3):e189–e202, 2025.
- Blanche Wies, Constantin Landers, and Marcello Ienca. Digital mental health for young people: a scoping review of ethical promises and challenges. *Frontiers in digital health*, 3:697072, 2021.
- Courtney N Wiesepeape, Sarah E Queller Soza, and Laura A Faith. Behind the gaps: A narrative review of healthcare barriers for individuals with serious mental illness. In *Healthcare*, volume 13, pp. 2387. MDPI, 2025.
- World Health Organization. Depressive disorder (depression) — WHO fact sheet, 2023. URL <https://www.who.int/news-room/fact-sheets/detail/depression>. Accessed: 2025-09-01.

- World Health Organization. Suicide — WHO fact sheet, 2025a. URL <https://www.who.int/news-room/fact-sheets/detail/suicide>. Accessed: 2025-09-01.
- World Health Organization. Over a billion people living with mental health conditions – services require urgent scale-up. <https://www.who.int/news/item/02-09-2025-over-a-billion-people-living-with-mental-health-conditions-services-require-urgent-scale-up>, 2025b. Accessed: 2025-10-30.
- Xuhai Xu, Bingsheng Yao, Yuanzhe Dong, Saadia Gabriel, Hong Yu, James Hendler, Marzyeh Ghassemi, Anind K Dey, and Dakuo Wang. Mental-llm: 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, 2024.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Haoran Yang, Yumeng Zhang, Jiaqi Xu, Hongyuan Lu, Pheng Ann Heng, and Wai Lam. Unveiling the generalization power of fine-tuned large language models. *arXiv preprint arXiv:2403.09162*, 2024a.
- Kailai Yang, Shaoxiong Ji, Tianlin Zhang, Qianqian Xie, Ziyang Kuang, and Sophia Ananiadou. Towards interpretable mental health analysis with large language models. *arXiv preprint arXiv:2304.03347*, 2023.
- Kailai Yang, Tianlin Zhang, Ziyang Kuang, Qianqian Xie, Jimin Huang, and Sophia Ananiadou. Mentallama: interpretable mental health analysis on social media with large language models. In *Proceedings of the ACM Web Conference 2024*, pp. 4489–4500, 2024b.
- Ou Yang and Yuting Zhang. Wait times for psychiatric specialist services in australia. *JAMA Network Open*, 8(2):e2461947–e2461947, 2025.
- Amir Hossein Yazdavar, Mohammad Saeid Mahdavi, Goonmeet Bajaj, William Romine, Amit Sheth, Amir Hassan Monadjemi, Krishnaprasad Thirunarayan, John M Meddar, Annie Myers, Jyotishman Pathak, et al. Multimodal mental health analysis in social media. *Plos one*, 15(4): e0226248, 2020.
- Dongsong Zhang, Jaewan Lim, Lina Zhou, and Alicia A Dahl. Breaking the data value-privacy paradox in mobile mental health systems through user-centered privacy protection: a web-based survey study. *JMIR Mental Health*, 8(12):e31633, 2021.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024.
- Assem Zhunis, Gabriel Lima, Hyeonho Song, Jiyoung Han, and Meeyoung Cha. Emotion bubbles: Emotional composition of online discourse before and after the covid-19 outbreak. In *Proceedings of the ACM Web Conference 2022*, pp. 2603–2613, 2022.

A APPENDIX

Task 1: "Is the poster of the post stressed? Classify the Users' post text into 0 and 1. 1 means stressed, 0 means not stressed."

Task 2: "Is the poster of the post experiencing suicidal ideation? Classify the following text as 0 or 1. 0 means user has no suicidal ideation, 1 means user has suicidal ideation."

Task 3: "Is the poster of the post depressed? Classify the Users' post text into 0 and 1. 1 means depressed, 0 means not depressed."

Task 4: "What level of depression is the person experiencing? Classify the Users' post text into 0, 1, 2 and 3. 3 means severe depression, 2 means moderate depression, 1 means mild depression, 0 means minimum depression."

Task 5: "Is the poster of the post at high risk or low risk for suicide? Classify the Users' post text into 0 and 1. 1 means high suicide risk, 0 means low suicide risk."

Task 6: "What is the suicide risk level of the poster of the post? Classify the Users' post text into 0, 1, 2, 3 and 4. 0 means Supportive, 1 means Ideation, 2 means Behavior, 3 means Attempt, 4 means Indicator."

Figure 13: Prompt Instruction for the Six Mental Health Tasks.

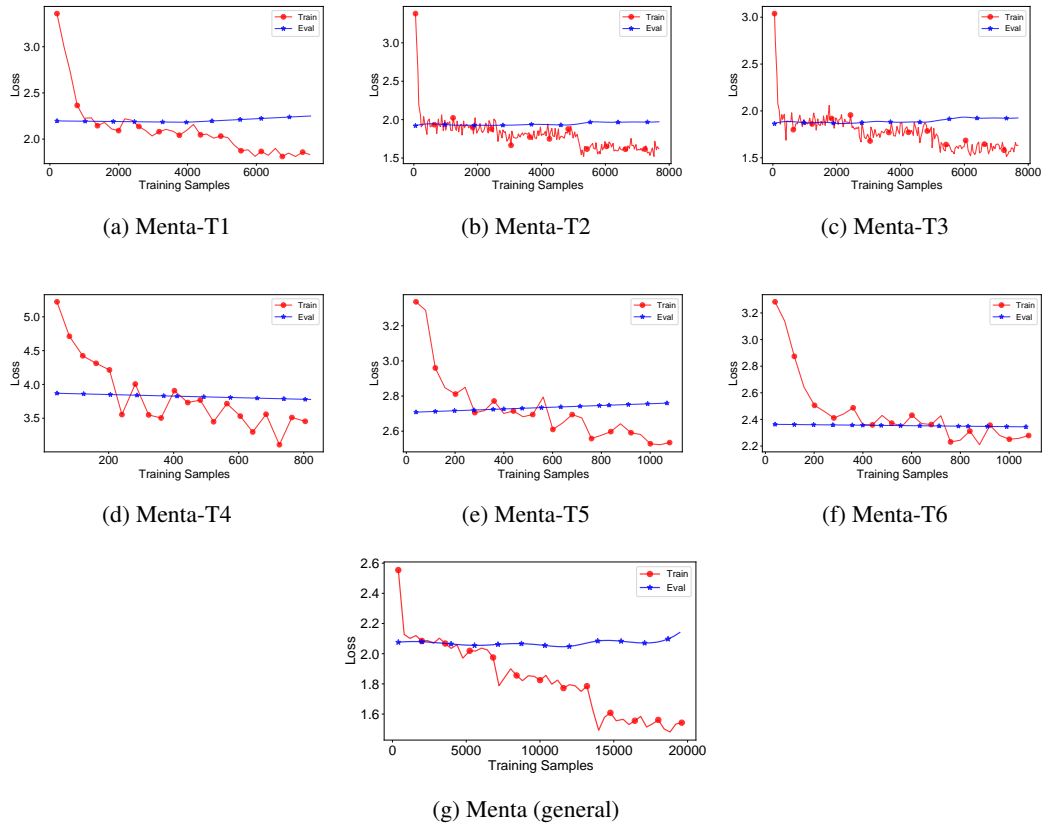


Figure 14: Training loss for the six individually fine-tuned models trained with Task 1-6 data respectively, compared with the Menta model trained on all six tasks.

Table 2: Accuracy of different models across tasks and few-shot settings (mean \pm standard deviation), with the best performance highlighted in bold and the second-best underlined.

Task	Model-Size	Zero-shot	One-shot	Two-shot	Three-shot	Four-shot
Task 1	Gemma-3 (1B)	0.48 \pm 0.01	0.53 \pm 0.03	0.55 \pm 0.02	0.56 \pm 0.03	0.52 \pm 0.02
	Gemma-3 (4B)	0.53 \pm 0.03	0.60 \pm 0.05	0.57 \pm 0.03	0.59 \pm 0.02	0.58 \pm 0.03
	Qwen-3 (4B)	<u>0.56</u> \pm 0.01	<u>0.64</u> \pm 0.04	0.62 \pm 0.02	<u>0.63</u> \pm 0.02	0.62 \pm 0.02
	Phi-4 Mini (3.8B)	0.57 \pm 0.02	0.69 \pm 0.02	0.59 \pm 0.02	0.66 \pm 0.03	0.62 \pm 0.02
	TinyLLaMA (1.1B)	0.50 \pm 0.01	0.51 \pm 0.01	0.51 \pm 0.01	0.50 \pm 0.05	0.52 \pm 0.01
	Falcon (1.3B)	0.51 \pm 0.01	0.51 \pm 0.02	0.51 \pm 0.01	0.51 \pm 0.01	0.53 \pm 0.01
	StableLM (3B)	0.50 \pm 0.00	0.59 \pm 0.00	<u>0.59</u> \pm 0.00	0.59 \pm 0.00	<u>0.59</u> \pm 0.00
Task 2	Gemma-3 (1B)	0.50 \pm 0.00	0.57 \pm 0.06	0.61 \pm 0.08	0.53 \pm 0.08	0.55 \pm 0.03
	Gemma-3 (4B)	0.63 \pm 0.01	<u>0.60</u> \pm 0.08	0.63 \pm 0.03	0.59 \pm 0.01	0.62 \pm 0.05
	Qwen-3 (4B)	0.70 \pm 0.01	0.66 \pm 0.03	0.67 \pm 0.01	0.70 \pm 0.02	0.67 \pm 0.01
	Phi-4 Mini (3.8B)	<u>0.69</u> \pm 0.01	0.66 \pm 0.05	<u>0.66</u> \pm 0.04	<u>0.68</u> \pm 0.02	<u>0.66</u> \pm 0.01
	TinyLLaMA (1.1B)	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.01	0.50 \pm 0.00	0.50 \pm 0.01
	Falcon (1.3B)	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.02	0.50 \pm 0.01
	StableLM (3B)	0.53 \pm 0.01	0.51 \pm 0.01	0.51 \pm 0.01	0.51 \pm 0.01	0.51 \pm 0.01
Task 3	Gemma3 (1B)	0.25 \pm 0.01	0.29 \pm 0.03	0.30 \pm 0.03	0.28 \pm 0.01	0.28 \pm 0.01
	Gemma3 (4B)	<u>0.30</u> \pm 0.01	<u>0.32</u> \pm 0.02	0.31 \pm 0.01	0.30 \pm 0.01	0.31 \pm 0.02
	Qwen-3 (4B)	<u>0.30</u> \pm 0.00	0.33 \pm 0.02	0.34 \pm 0.01	<u>0.34</u> \pm 0.02	0.37 \pm 0.01
	Phi-4 Mini (3.8B)	0.35 \pm 0.01	0.38 \pm 0.01	0.38 \pm 0.00	0.37 \pm 0.01	0.37 \pm 0.01
	TinyLLama (1.1B)	0.25 \pm 0.00	0.29 \pm 0.00	0.29 \pm 0.00	0.29 \pm 0.00	0.29 \pm 0.00
	Falcon (1.3B)	0.24 \pm 0.01	0.25 \pm 0.00	0.26 \pm 0.01	0.26 \pm 0.01	0.26 \pm 0.01
	StableLM (3B)	0.24 \pm 0.01	0.29 \pm 0.00	0.29 \pm 0.00	0.29 \pm 0.00	<u>0.29</u> \pm 0.00
Task 4	Gemma-3 (1B)	0.36 \pm 0.01	0.35 \pm 0.02	0.35 \pm 0.02	0.36 \pm 0.01	0.34 \pm 0.01
	Gemma-3 (4B)	0.37 \pm 0.01	0.37 \pm 0.01	0.37 \pm 0.01	0.36 \pm 0.01	0.36 \pm 0.01
	Qwen-3 (4B)	<u>0.41</u> \pm 0.01	<u>0.40</u> \pm 0.01	<u>0.41</u> \pm 0.01	<u>0.41</u> \pm 0.01	<u>0.42</u> \pm 0.01
	Phi-4 Mini (3.8B)	0.42 \pm 0.01	0.43 \pm 0.01	0.43 \pm 0.01	0.43 \pm 0.01	0.43 \pm 0.01
	TinyLLama (1.1B)	0.29 \pm 0.00	0.29 \pm 0.00	0.30 \pm 0.00	0.30 \pm 0.00	0.30 \pm 0.01
	Falcon (1.3B)	0.27 \pm 0.01	0.28 \pm 0.00	0.29 \pm 0.01	0.29 \pm 0.00	0.28 \pm 0.00
	StableLM (3B)	0.27 \pm 0.00	0.28 \pm 0.00	0.28 \pm 0.00	0.28 \pm 0.00	0.28 \pm 0.00
Task 5	Gemma-3 (1B)	0.46 \pm 0.01	0.49 \pm 0.01	0.48 \pm 0.03	0.49 \pm 0.01	0.50 \pm 0.03
	Gemma-3 (4B)	0.48 \pm 0.00	<u>0.54</u> \pm 0.02	<u>0.51</u> \pm 0.02	<u>0.53</u> \pm 0.01	<u>0.54</u> \pm 0.02
	Qwen-3 (4B)	<u>0.53</u> \pm 0.01	0.51 \pm 0.01	0.49 \pm 0.03	0.49 \pm 0.01	0.49 \pm 0.01
	Phi-4 Mini (3.8B)	0.58 \pm 0.01	0.58 \pm 0.02	0.61 \pm 0.02	0.60 \pm 0.01	0.59 \pm 0.01
	TinyLLama (1.1B)	0.50 \pm 0.00	0.46 \pm 0.01	<u>0.51</u> \pm 0.02	0.49 \pm 0.02	0.48 \pm 0.01
	Falcon (1.3B)	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.00	0.50 \pm 0.00
	StableLM (3B)	0.46 \pm 0.02	0.50 \pm 0.01	0.49 \pm 0.02	0.50 \pm 0.01	0.51 \pm 0.02
Task 6	Gemma-3 (1B)	0.20 \pm 0.01	0.20 \pm 0.00	0.20 \pm 0.01	0.20 \pm 0.01	0.20 \pm 0.01
	Gemma-3 (4B)	<u>0.22</u> \pm 0.01	0.24 \pm 0.02	0.24 \pm 0.01	0.26 \pm 0.01	0.27 \pm 0.01
	Qwen-3 (4B)	0.23 \pm 0.12	0.26 \pm 0.02	0.27 \pm 0.02	0.26 \pm 0.02	0.27 \pm 0.01
	Phi-4 Mini (3.8B)	0.21 \pm 0.00	0.26 \pm 0.00	0.26 \pm 0.01	0.26 \pm 0.00	0.24 \pm 0.03
	TinyLLama (1.1B)	<u>0.22</u> \pm 0.02	0.21 \pm 0.01	0.21 \pm 0.01	0.21 \pm 0.01	0.21 \pm 0.01
	Falcon (1.3B)	0.20 \pm 0.00	0.20 \pm 0.01	0.19 \pm 0.01	0.20 \pm 0.01	0.21 \pm 0.02
	StableLM (3B)	0.20 \pm 0.00	<u>0.25</u> \pm 0.00	0.24 \pm 0.00	<u>0.24</u> \pm 0.01	<u>0.24</u> \pm 0.00

Table 3: Balanced Accuracy (BACC) results across six tasks with mean and standard deviation, with the best performance highlighted in bold and the second-best underlined.

Model	Task 1	Task 2	Task 3	Task 4	Task 5	Task 6
Menta-T1	0.83 \pm 0.02	0.67 \pm 0.04	0.30 \pm 0.10	0.68 \pm 0.04	0.55 \pm 0.01	0.25 \pm 0.01
Menta-T2	0.79 \pm 0.04	0.75 \pm 0.04	0.31 \pm 0.01	0.65 \pm 0.04	<u>0.56</u> \pm 0.02	0.23 \pm 0.03
Menta-T3	0.77 \pm 0.02	0.66 \pm 0.05	<u>0.37</u> \pm 0.06	<u>0.66</u> \pm 0.03	0.52 \pm 0.02	0.26 \pm 0.02
Menta-T4	0.76 \pm 0.01	0.63 \pm 0.04	0.27 \pm 0.04	0.63 \pm 0.03	<u>0.56</u> \pm 0.01	<u>0.25</u> \pm 0.02
Menta-T5	0.58 \pm 0.02	0.55 \pm 0.04	0.29 \pm 0.03	0.51 \pm 0.00	0.47 \pm 0.06	0.22 \pm 0.03
Menta-T6	0.66 \pm 0.05	0.65 \pm 0.05	0.27 \pm 0.03	0.54 \pm 0.02	0.47 \pm 0.06	0.24 \pm 0.04
Menta (General Model)	<u>0.81</u> \pm 0.02	<u>0.73</u> \pm 0.05	0.38 \pm 0.07	0.61 \pm 0.03	0.60 \pm 0.03	0.26 \pm 0.05

Table 4: Consolidated deployment performance across all tasks, with the best performance highlighted in bold and the second-best underlined.

Task	Model	TTFT (s)	ITPS (/s)	OET (s)	OTPS (/s)	Total Time (s)	RAM (GB)
Task 1	Phi-4 Mini	1.239	870.8	8.148	<u>5.4</u>	5834.33	2.58
	Qwen-3	3.705	<u>3516.5</u>	5.193	<u>1.9</u>	3718.19	<u>2.94</u>
	Menta	4.024	3587.0	<u>6.335</u>	1.8	<u>4535.86</u>	2.99
	StableLM	<u>2.599</u>	37.5	23.2	6.6	16541.60	4.77
Task 2	Phi-4 Mini	1.376	723.1	<u>8.174</u>	<u>4.4</u>	<u>29050.40</u>	2.58
	Qwen-3	4.202	2941.3	6.760	1.6	24025.04	<u>3.03</u>
	Menta	4.777	<u>2931.6</u>	10.885	1.9	38685.29	<u>3.04</u>
	StableLM	<u>2.486</u>	38.6	22.54	6.7	80084.62	3.89
Task 3	Phi-4 Mini	2.333	1055.1	<u>8.384</u>	<u>5.1</u>	<u>29668.79</u>	2.58
	Qwen-3	4.288	4337.2	6.834	1.6	24265.29	<u>3.03</u>
	Menta	5.083	<u>4275.4</u>	12.837	1.8	45622.70	<u>3.04</u>
	StableLM	<u>2.512</u>	38.3	22.288	6.8	79189.26	3.91
Task 4	Phi-4 Mini	1.236	1097.0	8.213	<u>3.4</u>	3120.94	2.90
	Qwen-3	3.042	2930.2	4.875	1.5	1852.50	<u>2.94</u>
	Menta	3.080	<u>2698.1</u>	4.910	1.4	1917.20	3.00
	StableLM	<u>2.708</u>	33.5	24.473	6.3	<u>9299.74</u>	4.11
Task 5	Phi-4 Mini	1.287	1041.5	8.656	<u>5.1</u>	4328.00	2.59
	Qwen-3	6.966	<u>5550.0</u>	25.428	1.3	12650.30	<u>2.95</u>
	Menta	6.005	6031.4	27.657	1.3	13828.50	3.04
	StableLM	<u>2.431</u>	37.4	<u>24.104</u>	6.3	<u>12052.00</u>	3.52
Task 6	Phi-4 Mini	1.227	1063.4	7.296	6.0	3648.00	2.59
	Qwen-3	6.904	<u>5416.7</u>	27.509	1.4	13754.50	3.04
	Menta	5.741	5851.4	<u>19.883</u>	1.3	<u>9941.50</u>	<u>2.94</u>
	StableLM	<u>2.51</u>	36.9	23.1	6.4	11500.00	3.53