IJCRT.ORG

ISSN: 2320-2882



INTERNATIONAL JOURNAL OF CREATIVE RESEARCH THOUGHTS (IJCRT)

An International Open Access, Peer-reviewed, Refereed Journal

A Survey On Privacy-Aware AI Approaches For Early Mental Health Detection And Support

¹ V. Kiruthiga, ²Dr. K. Lakshmi Priya ¹Ph.D. Scholar, ²Associate Professor ¹Department of Computer Science, ¹Karpagam Academy of Higher Education, Coimbatore, India.

Abstract: With the increasing burden of mental health disorders such as depression and anxiety, Artificial Intelligence (AI) is playing a significant role in early detection and intervention. However, deploying AI in this sensitive domain raises serious privacy concerns. This survey explores the landscape of AI-based systems for early mental health detection, with a strong focus on privacy-aware designs. It reviews existing models that analyse text, voice, and multimodal data, and evaluates them based on their privacy-preserving capabilities. The paper also identifies key challenges, existing gaps in the literature, and future research directions for developing secure, ethical, and user-centric AI solutions.

Index terms: Mental Health, AI, Privacy-Aware Systems, Multimodal AI, Differential Privacy, Federated Learning, Emotion Recognition

1. Introduction

Mental health disorders, including depression, anxiety, and stress-related conditions, are among the leading causes of disability worldwide. According to the World Health Organization (WHO), more than 280 million people globally suffer from depression, and these numbers continue to rise in the post-pandemic digital age. Despite growing awareness, early diagnosis and timely intervention remain significant challenges due to stigma, insufficient clinical access, and delays in recognizing symptoms. The surge in online activity, mobile phone usage, and digital communication presents a unique opportunity to harness Artificial Intelligence (AI) for real-time mental health screening and support.

Recent advances in machine learning (ML), deep learning (DL), and natural language processing (NLP) have enabled the development of AI systems that can detect early signs of mental health conditions using text, voice, and multimodal inputs. Such systems analyse behavioural cues, speech patterns, and social media expressions to infer emotional and cognitive states. While these AI-driven interventions offer scalability and accessibility, they also raise serious ethical concerns, primarily regarding user privacy, data ownership, and trustworthiness of AI-generated decisions.

Mental health data is inherently sensitive and personal. Mishandling or unauthorized use of such data could lead to discrimination, exploitation, or psychological harm. To address these concerns, researchers are incorporating privacy-preserving mechanisms into AI systems. Techniques such as federated learning, differential privacy, homomorphic encryption, and edge computing allow AI models to learn from user data without directly accessing it. These approaches not only enhance data security but also align with evolving legal and ethical standards.

This survey critically examines state-of-the-art AI models used in mental health detection and evaluates how privacy-aware design is integrated into their architecture. We categorize existing works based on input modalities (text, speech, multimodal), discuss the privacy techniques applied, and highlight datasets commonly used in the field. Furthermore, this paper identifies current research gaps and suggests future directions for developing ethically sound, privacy-compliant AI systems in mental health care.

2. Research Motivatio006E

The growing burden of mental health disorders has created an urgent demand for scalable and accessible screening systems. Several global trends underscore the need for AI-driven, privacy-preserving mental health tools:

- Digital Mental Health Revolution: The post-COVID era has seen a dramatic increase in the use of mental health apps and online therapy platforms. However, many of these services are underregulated and lack proper data protection mechanisms.
- AI's Diagnostic Potential: Studies have shown that AI can identify subtle linguistic, acoustic, and behavioural markers of depression and anxiety more effectively than traditional assessments, especially when trained on multimodal datasets.
- Privacy as a Barrier to Adoption: Surveys suggest that a significant portion of users, over 60% are hesitant to share sensitive emotional data with AI systems due to concerns about surveillance, data leaks, and misuse.
- Regulatory Pressure and Ethical Expectations: With the implementation of data protection laws like GDPR and HIPAA, researchers and developers must now consider ethical compliance, accountability, and transparency in model design.

Given these motivations, this survey focuses on evaluating privacy-aware AI methods for mental health monitoring that can balance diagnostic accuracy with ethical responsibility and legal compliance.

3. Literature Survey:

The literature on privacy-aware AI for mental health is organized under four key modalities:

3.1 Text-Based Mental Health Detection

Text-based approaches are among the most widely explored modalities for early mental health detection. These methods utilize written language from platforms such as Reddit, Twitter, clinical interviews, personal diaries, and chatbot conversations to infer users' psychological states. The linguistic structure, emotional tone, grammatical features, and syntactic complexity can provide significant insights into depression, anxiety, and suicidal ideation.

Early works applied traditional machine learning algorithms such as Support Vector Machines (SVMs), Logistic Regression, and Naïve Bayes with features like term frequency-inverse document frequency (TF-IDF) and sentiment scores. However, with the emergence of deep learning, models such as Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), and Transformer-based architectures like BERT (Bidirectional Encoder Representations from Transformers) have significantly improved classification performance.

For example, Morales et al. [3] used linguistic features to detect depression from written transcriptions, while Jin et al. [7] demonstrated the effectiveness of combining BERT with attention mechanisms in identifying emotional patterns in patient notes. More recently, generative models such as GPT have been explored for their potential to assess mental health from language context.

Despite the success of these models, text-based data carries inherent privacy risks. Written content may include names, events, or health details that can easily be linked back to individuals. Privacy-preserving strategies, such as anonymization and differential privacy, are therefore essential when dealing with large-scale text datasets used in mental health studies.

3.2 Voice-Based Emotion Recognition

Voice is one of the most expressive and naturally available modalities for detecting emotional and mental health states. Human speech conveys not just verbal content but also paralinguistic features such as pitch, tone, tempo, jitter, and pauses—elements that can be indicative of psychological conditions like depression and anxiety. Researchers have explored various acoustic features like Mel-frequency cepstral coefficients (MFCCs), prosody, and spectral energy distribution to model and interpret emotional speech signals.

Traditional classifiers like SVM and Random Forest have been used with handcrafted features. However, modern approaches leverage deep learning models such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and hybrid CNN-LSTM architectures to learn hierarchical acoustic features automatically. For instance, Tzirakis et al. [2] achieved significant accuracy in multimodal emotion recognition by integrating CNN with LSTM for speech and video data. Attention mechanisms have also been introduced to weight temporal voice patterns dynamically.

Datasets like DAIC-WOZ, AVEC, and EmoDB are widely used in this domain. The models trained on these datasets can detect not only emotion (e.g., sadness, fear) but also map these emotions to potential mental health states (e.g., mild or severe depression).

Despite high classification accuracy, voice data poses notable privacy risks. It can be reverse-engineered to reveal a speaker's identity, mood, or even health status. Techniques like anonymization, voice encryption, ondevice feature extraction, and privacy-aware model training using federated learning are increasingly adopted to minimize these risks. As voice interfaces become more prevalent in mobile health apps and smart speakers, ensuring secure voice-based inference is essential for ethical AI deployment.

3.3 Multimodal AI Systems

Multimodal AI systems combine two or more input sources—such as text, voice, facial expressions, physiological signals (like heart rate or skin temperature), and behaviour logs—to improve the accuracy and reliability of mental health detection. These systems aim to replicate the multifaceted diagnostic methods used by human clinicians, where emotion and mental state are inferred not from a single indicator but from a combination of cues.

For example, the DAIC-WOZ dataset includes audio, video, and text transcripts from clinical interviews and has been widely used to train multimodal depression classifiers. Other datasets like AVEC provide synchronized audio-visual data for emotion and affective state prediction.

Deep learning models such as multimodal transformers, attention-based fusion networks, and late-fusion CNN-RNN architectures have been used to learn patterns across modalities. These models show significantly higher predictive power compared to unimodal systems, particularly when data is incomplete or noisy in one modality.

However, integrating multiple data sources also increases the risk of user re-identification. A multimodal signature (e.g., someone's voice + typed responses + facial expressions) can be highly unique. Hence, privacy-preserving mechanisms become even more essential in this context. Techniques such as secure feature fusion, data obfuscation, and hybrid local-global federated learning are being explored to enhance the privacy of multimodal systems while retaining model performance.

3.4 Privacy-Preserving AI Techniques

As AI systems increasingly integrate into sensitive domains like mental health, preserving user privacy has become a central research challenge. Traditional centralized training models require user data to be uploaded to a cloud or server, raising concerns about data misuse, breaches, and regulatory non-compliance. In response, several privacy-enhancing techniques have been adopted in recent years to build secure, ethical, and legally compliant AI systems.

- Federated Learning: This decentralized training method allows AI models to be trained locally on a user's device without transferring raw data to a central server. Only model updates (gradients) are shared and aggregated to build a global model. Federated learning is especially suitable for mobile mental health apps, where data remains on the user's phone.
- Differential Privacy: This technique introduces calibrated noise to the input data or model outputs to protect individual data points from being reconstructed or identified. It provides a mathematical privacy guarantee and is often combined with other models like federated learning for stronger security.
- On-Device Inference: Rather than relying on remote servers, on-device inference allows models to run directly on a user's smartphone or wearable device. This ensures that data never leaves the local system, reducing the risk of interception or storage vulnerabilities.

- Homomorphic Encryption and Secure Multiparty Computation (SMPC): These cryptographic techniques allow computations on encrypted data. Homomorphic encryption permits model training and inference without ever decrypting the data, while SMPC enables collaborative model building across multiple users without data sharing.
- Synthetic Data Generation: Generating privacy-preserving synthetic datasets using generative models like GANs (Generative Adversarial Networks) or VAEs (Variational Autoencoders) can reduce the need to use identifiable real-world mental health records.

Together, these privacy-preserving approaches are shaping a new generation of ethical, secure AI systems in mental health applications, supporting trust, compliance, and long-term user engagement. Runs models locally, ensuring data never leaves the device.

4. Comparative Analysis

Study	Year	Input Modality	Algorithm	Privacy Technique	Accuracy	Dataset
Morales et al.	2018	Text & Audio	LSTM + SVM	None	74.5%	DAIC- WOZ
Tzirakis et al.	2017	Audio Video	CNN + LSTM	None	81.3%	AVEC
Danner et al.	2023	Text	GPT-based	Differential Privacy	85.1%	Simulated
Jin et al.	2023	Text & & Speech	BERT + Attention	Federated Learning	87.6%	Internal
Devlin et al.	2019	Text	BERT	None	90.0%	Reddit
Lee et al.	2021	Multimodal	Transformer Fusion	Differential Privacy	88.4%	AVEC+

5. Key Challenges and Gaps

Despite significant progress in AI-driven mental health analysis, several persistent challenges remain:

- Lack of Real-World, Privacy-Compliant Datasets: Most existing datasets are either synthetic or anonymized post-collection, which limits real-world generalizability and robustness of models.
- Model Explainability: Black-box models (e.g., deep neural networks) lack interpretability, making it difficult for clinicians to trust or validate the predictions.
- Bias and Fairness: AI models often inherit biases from training data, resulting in inconsistent performance across age groups, genders, and cultural backgrounds.
- Cross-Platform Generalization: Systems trained on one platform (e.g., social media) often fail when applied to clinical or mobile app settings.
- Trade-Off Between Privacy and Accuracy: Privacy-preserving methods may reduce model performance, making it challenging to balance ethical design with predictive strength.
- Regulatory and Ethical Uncertainty: A lack of unified legal standards for AI in mental health leads to uncertainty in deployment and compliance.

6. Conclusion and Future Scope

This survey has reviewed the current landscape of privacy-aware AI systems for early mental health detection. Through analysis of text, voice, and multimodal models, and an evaluation of privacy-enhancing techniques, it is evident that the field is moving toward ethically conscious, user-centric design. However, substantial gaps remain, particularly in real-world testing, bias mitigation, and transparent model architectures.

Future research should focus on:

- Developing large-scale, privacy-first datasets representative of diverse populations
- Integrating layered privacy mechanisms, such as combining differential privacy with federated learning to enhance data protection while maintaining model performance
- Exploring Explainable AI (XAI) and Human-in-the-Loop systems
- Designing compliant frameworks that align with GDPR, HIPAA, and other global policies

By embedding privacy from the ground up, researchers and developers can ensure AI becomes a trusted ally in improving global mental health outcomes.

7. References

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proc. of NAACL-HLT*, Minneapolis, USA, 2019.
- [2] P. Tzirakis, G. Trigeorgis, M. A. Nicolaou, B. Schuller, and S. Zafeiriou, "End-to-End Multimodal Emotion Recognition using Deep Neural Networks," *IEEE Journal of Selected Topics in Signal Processing*, vol. 11, no. 8, pp. 1301–1309, Dec. 2017.
- [3] M. Morales, S. I. Levitan, and R. Levitan, "A Multimodal Approach for Depression Detection with DAIC-WOZ Dataset," in *Proc. of the 5th International Workshop on Audio/Visual Emotion Challenge*, 2015.
- [4] M. Abadi et al., "Deep Learning with Differential Privacy," in *Proc. ACM CCS*, 2016, pp. 308–318.
- [5] H. B. McMahan et al., "Communication-Efficient Learning of Deep Networks from Decentralized Data," in *Proc. of AISTATS*, 2017.
- [6] M. Danner, S. Kliem, and A. Böckler, "GPT-based Methods for Depression Detection: A Neural Approach for Emotional Assessment," *IEEE Access*, vol. 11, pp. 33121–33134, 2023.
- [7] K. W. Jin, S. Y. Lim, and J. H. Lee, "Artificial Intelligence in Mental Healthcare: A Comprehensive Review," *IEEE Reviews in Biomedical Engineering*, vol. 16, pp. 89–104, 2023.
- [8] M. Casu, L. Pulina, and G. Giuffrida, "AI Chatbots for Mental Health: A Review of Recent Advances," *IEEE Transactions on Emerging Topics in Computational Intelligence*, vol. 7, no. 3, pp. 546–558, 2023.