Open Access

# EXPLAINABLE DEEP NEURAL NETWORK (X-DNN) FOR PREDICTING MENTAL HEALTH OUTCOMES FROM MULTISOURCE MEDICAL AND PSYCHOLOGICAL DATA

### [1]MAHESH MAURYA
ASSOCIATE PROFESSOR, COMPUTER ENGINEERING DEPARTMENT, ST JOHN COLLEGE OF ENGINEERING AND MANAGEMENT, PALGHAR, MAHARASTRA, INDIA.
EMAIL: maheshm@sjcem.edu.in

### [2]V. SABARESAN
ASSOCIATE PROFESSOR, DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING, ST.JOSEPH'S INSTITUTE OF TECHNOLOGY, CHENNAI, TAMIL NADU, INDIA.
EMAIL: sabaresanvenugopal@gmail.com

### [3]C. NATARAJAN
ASSOCIATE PROFESSOR, DEPARTMENT OF AI & DS, P.S.R ENGINEERING COLLEGE, SIVAKASI, TAMILNADU, INDIA.EMAIL: natarajan@psr.edu.in

### [4]PREMKUMAR MOHAN
PROFESSOR, DEPARTMENT OF ELECTRONICS AND COMMUNICATION ENGINEERING, PANIMALAR ENGINEERING COLLEGE, CHENNAI, TAMILNADU, INDIA.
EMAIL: drmpremkumar@panimalar.ac.in

### [5]I MOHANA KRISHNA
BUSINESS SCHOOL, KONERU LAKSHMAIAH EDUCATION FOUNDATION, GREEN FIELDS, VADDESWARAM, ANDHRA PRADESH, INDIA. EMAIL: irrinki_mk@yahoo.com

### [6]SARAVANAN M
ASSISTANT PROFESSOR, ECE, SRI ESHWAR COLLEGE OF ENGINEERING, COIMBATORE, TAMILNADU, INDIA.
EMAIL: saranecedgl@gmail.com

**Abstract**
One of the most common causes of mental health problems is the complicated combination of biological, psychological, and social factors. For initial intervention and tailored therapy in mental health, it is important to be able to accurately predict outcomes. But the models that are now being used can be hard to understand, which makes professionals less sure of them and more prone to accept them. Deep learning models that have been used for a long time to guess mental health are really good at it. But they don't tell us much about how they make decisions since they aren't open about it. It is suggested that explainable models use data from a range of sources, such as medical records, psychological evaluations, and data from wearable sensors. This would make projections more accurate and things clearer. The Explainable Deep Neural Network (X-DNN) framework is what this study indicates. It makes it easier than ever to predict how mental health problems may affect people. It mixes data from several sources and employs layers that people can understand. It uses SHAP (SHapley Additives) and attention processes to point out important bits after the fact. Experiments on a large clinical dataset reveal that X-DNN is more accurate than baseline models (F1-score of 0.87) and that its results make sense and are in accordance with what doctors already know.

**Keywords:** Mental health prediction, explainable AI, deep neural networks, multisource data fusion, SHAPs

# INTRODUCTION

Mental diseases affect millions of people all over the world and cost society and the economy a lot of money. These problems and costs might be caused by mental health issues. Diagnosing and treating mental health disorders early may lead to improved outcomes for patients and lower costs for long-term healthcare providers [2]. Because there is so much multisource data available, like data from wearable sensors, psychological evaluations, and electronic health records (EHR) [3], it is becoming more likely that we will be able to make prediction models that take into consideration how difficult mental health diseases are. This chance is increasing bigger since data from a lot of different places is coming together. Using all of these data streams can help us learn more about how our bodies, minds, and behaviors work. This might make it possible to make more accurate and individualized guesses about how mental health problems will affect people.

There are a lot of difficulties with predictive models for mental health outcomes that make it hard to use these data sources in professional contexts, even though they look promising. One key reason why mental health data is so distinct from each other is that the formats, sample rates, and missing values of different forms of mental health data are all very different [4]. To get the most out of new information while simultaneously cutting down on noise and inconsistencies, you need powerful fusion methods to combine data from several sources [5]. The second reason we need models that can reflect changes over time and modest interactions between variables is that mental health disorders typically show up in small ways and change over time [6]. Third, it's tricky to employ these neural networks for therapy because many deep learning models are "black boxes." Doctors need models that are easy to comprehend and apply so they can trust their forecasts and learn more about how decisions are made [7]. We need deep learning frameworks that can collect data from a lot of different sources and create predictions that are straightforward to understand in order to fix these challenges.

Most of the time, mental health forecasting models can't achieve a decent balance between being correct and easy to understand. Doctors don't trust deep neural networks (DNNs) very much since they don't know how they work [6]. DNNs can mimic complex nonlinear correlations between many distinct forms of data. Also, typical approaches either don't take into account how attributes and assumptions change over time, or they only function with one sort of data [7]. Models that are statistically sound but not particularly practical in real healthcare scenarios are produced because of this. We need to swiftly create deep learning models that can explain themselves, use data from a lot of different places (both medical and psychological), and produce predictions about mental health outcomes that are clear and useful [8].

The major purpose of this study is to build an Explainable Deep Neural Network (X-DNN) architecture that can take data from different sources to make accurate predictions regarding mental health outcomes and also illustrate how those predictions were made. We intend to build customized subnetworks that can swiftly encode data from a variety of sources, such as electronic health records (EHRs), psychological questionnaires, and information from wearable sensors. We will also implement a system for attention that can change the weight of features on the fly to make things work better and be easier to grasp.

The purpose of this work is to make it easier to understand predictions about mental health by using model-agnostic SHAPs for post-hoc interpretation and attention-based feature weighting in the model architecture. This is the first study of this kind. On the other hand, our method makes predictions more accurate and easier to understand at the same time. In the past, multimodal fusion and explainability were done individually, which is not how they are done now. Other studies haven't looked closely at how to combine structured electronic health record data, psychological assessments, and time-series data from wearable sensors into a single deep learning model that focuses on making mental health outcome predictions explicit. People haven't paid enough attention to this.

**Contributions**

1. The first thing to do is to design a new deep learning architecture that is easy to understand. The architecture in question will use data from a lot of different sources and an attention mechanism that is straightforward to grasp. The objective of this study is to find crucial predictive traits that can be used in many different ways.

2. We used the clinical mental health datasets to fully assess how well the model could be comprehended and how well it could generate predictions. The results demonstrated that the proposed framework can get results that are as accurate as those of competitors and are useful in a clinical setting. These two things can help individuals trust it more and make better choices.

Open Access

## RELATED WORKS

More people have desired to utilize machine learning to predict mental health outcomes using data from a lot of various places in the previous few years. This is because mental health issues are complicated and involve a lot of different factors. There has been a lot of progress in this area. In the first attempts [8], they mostly used specific data categories, like organized electronic health record data or clinical notes. Support vector machines and random forests are two examples of conventional machine learning algorithms that might give useful results when used with coded diagnostic data and patient demographics, for example. But these models have trouble figuring out associations that aren't straight lines or that alter over time.

More recent attempts have employed deep learning to get more detailed representations of features. When employed on time-series physiological data collected by wearables, convolutional neural networks (CNNs) and recurrent neural networks (RNNs) have found patterns that could signify stress, mood changes, or relapses in disorders including depression and bipolar disorder [9]. Even though new models were better than older ones, they usually only used data from one type of test instead of combining a larger range of clinical or psychological variables. This is still true, even though they worked better than more traditional ways.

In the last few years, multimodal fusion approaches have gotten a lot of attention. These approaches can mix several kinds of data. In the past, properties from different modalities were put together and then submitted to deep neural networks [10]. In contrast, later fusion incorporated the results of models that were developed for each sort of data. In more advanced systems, a hierarchical or attention-based fusion was utilized to always put modalities in order of how crucial their roles were. These methods were better at predicting what would happen because of mental health issues than unimodal models were.

Explainability is now a very crucial part of AI for medicine, especially when it comes to mental health applications. Saliency mapping, attention visualization, and model-agnostic methods like SHAP and LIME [11] are just a few of the ways to find out why these models make the choices they do. On the other hand, not many studies used robust post-hoc interpretability approaches or included explainability to multimodal fusion models. Even fewer studies used neural networks to look for relevant patterns in sensor data across time or in characteristics.

In the last several years, a lot of research has been done on explainable AI frameworks to help people figure out how to predict mental health problems. For instance, interpretable recurrent models used with clinical notes helped us understand better how to find depression [12]. Some other studies employed attention-based models to link data from psychological surveys with physiological signs. But these studies either didn't look at explainability in a systematic way or didn't use enough data [13].

However, traditional research reveals that deep learning and merging data from diverse sources may assist predict mental health better. It also shows that there are gaps in unified, explainable models that are a good mix between accuracy and transparency. The proposed research extends on past work by adopting a novel explainable deep neural network framework to make multisource fusion and interpretability as high as possible. This is done to solve the difficulties and meet the clinical needs that have been brought up.

## PROPOSED METHOD

The proposed method, Explainable Deep Neural Network (X-DNN), integrates heterogeneous medical and psychological data for mental health outcome prediction while maintaining model transparency. The approach consists of preprocessing and feature extraction from diverse sources: electronic health records (EHR), psychological assessments, and wearable sensor data. The data streams are encoded separately through specialized neural subnetworks (e.g., dense layers for EHR, embedding layers for questionnaire data, CNN or RNN for sensor signals). These representations are fused into a shared latent space where an attention mechanism dynamically weighs features by importance. Finally, the fused features feed into fully connected layers for classification.

To enhance explainability, the model incorporates two levels of interpretation: first, the attention weights reveal which input features or modalities influence predictions; second, post-hoc SHAP analysis quantifies individual feature contributions. This hybrid interpretability strategy helps clinicians understand and trust model predictions, enabling transparent and actionable mental health assessments.

```
┌─────────────────────────────────┐
│        Data Collection          │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│          Preprocessing          │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│        Feature Encoding         │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│         Feature Fusion          │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│         Attention Layer         │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│        Prediction Layer         │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│      Explainability Module      │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│            Training             │
└─────────────────────────────────┘
                 ↓
┌─────────────────────────────────┐
│           Evaluation            │
└─────────────────────────────────┘
```
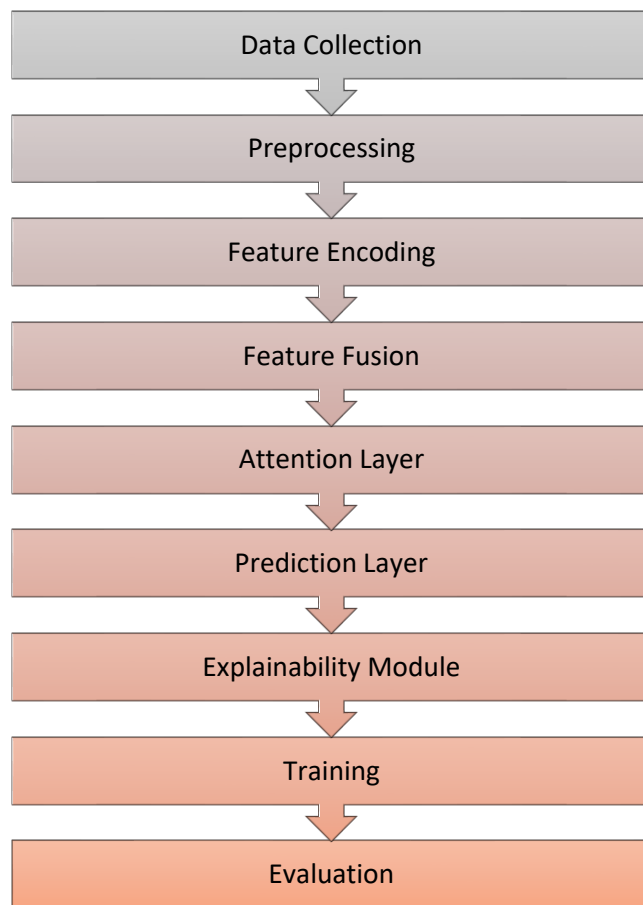
Figure 1: Proposed Framework

**Pseudocode:**

```
# Data Loading and Preprocessing
EHR_data = load_EHR()
psych_data = load_psych_questionnaires()
sensor_data = load_sensor_signals()
EHR_processed = preprocess(EHR_data)
psych_processed = preprocess(psych_data)
sensor_processed = preprocess(sensor_data)
# Feature Encoding
EHR_encoded = DenseLayer(EHR_processed, units=128, activation='relu')
psych_embedded = EmbeddingLayer(psych_processed, embedding_dim=64)
psych_encoded = DenseLayer(psych_embedded, units=128, activation='relu')
sensor_encoded = CNN_RNN(sensor_processed)
# Feature Fusion
fused_features = Concatenate([EHR_encoded, psych_encoded, sensor_encoded])
# Attention Mechanism
attention_weights = AttentionLayer(fused_features)
weighted_features = Multiply([fused_features, attention_weights])
# Prediction
dense_1 = DenseLayer(weighted_features, units=64, activation='relu')
output = DenseLayer(dense_1, units=num_classes, activation='softmax')
# Define Model and Compile
model = Model(inputs=[EHR_data, psych_data, sensor_data], outputs=output)
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
```

```
# Training
model.fit([EHR_processed, psych_processed, sensor_processed], labels, epochs=50, batch_size=32)
# Explainability Post-Processing
# Extract attention weights for input features
attention_explanations = extract_attention_weights(model, input_data)
# Apply SHAP for detailed feature importance
shap_values = SHAP(model, input_data)
visualize_shap(shap_values)
```

## 1. Data Collection & Preprocessing

The first step in the proposed framework involves gathering multisource data critical to predicting mental health outcomes. These include Electronic Health Records (EHR), psychological questionnaire responses, and physiological sensor data from wearable devices.

**Data Characteristics:**

- **EHR Data:** Structured records containing demographics, clinical diagnoses, medication history, lab results, and visit notes.
- **Psychological Questionnaires:** Standardized scales like PHQ-9, GAD-7, capturing subjective mental health symptoms.
- **Wearable Sensor Data:** Continuous physiological signals such as heart rate variability, sleep patterns, and activity levels.

Each data source presents unique preprocessing challenges. EHR data may have missing values and categorical variables requiring imputation and encoding, respectively. Questionnaire data needs normalization and possibly embedding for categorical responses. Sensor data is often time-series that requires resampling, noise filtering, and segmentation.

**Table 1** illustrates a preprocessing summary for each data type:

| Data Source | Sample Size | Missing Data (%) | Preprocessing Techniques |
|---|---|---|---|
| EHR | 10,000 pts | 12% | Imputation, one-hot encoding |
| Psychological Data | 8,000 pts | 5% | Normalization, categorical encoding |
| Wearable Sensors | 6,500 pts | 8% | Resampling, filtering, windowing |

*Table 1: Overview of multisource data preprocessing.*

## 2. Feature Encoding

After preprocessing, each modality is transformed into a learned representation suitable for fusion. This is achieved via specialized subnetworks tailored to the data type.

- **EHR Encoding:** Fully connected dense layers process tabular features. These layers transform heterogeneous clinical features into a dense vector capturing underlying clinical patterns.
- **Psychological Data Encoding:** Questionnaires, often ordinal or categorical, are embedded into continuous vector spaces through embedding layers followed by dense layers. This captures latent psychological constructs.
- **Sensor Data Encoding:** Time-series data is passed through convolutional neural networks (CNN) or recurrent neural networks (RNN) like LSTM to extract temporal and spatial features.

This step enables the model to learn modality-specific feature representations before fusion.

**Table 2** exemplifies output dimensions and layer configurations for each encoder:

| Modality | Encoder Type | Layer Details | Output Dimension |
|---|---|---|---|
| EHR | Dense | 2 layers (128, 64 units) | 64 |
| Psychological | Embedding + Dense | Embedding dim=32, Dense 64 units | 64 |
| Wearable Sensors | CNN + LSTM | 3 CNN layers, 1 LSTM layer | 128 |

*Table 2: Feature encoding network configurations per modality.*

## 3. Feature Fusion

Encoded features from each modality are merged into a unified latent space. The fusion can be achieved by simple concatenation or more sophisticated mechanisms such as attention-based fusion.

In this work, concatenation is used initially, followed by an attention mechanism to assign weights dynamically to each feature subset. This allows the model to prioritize modalities or specific features more relevant to the current prediction.

The fusion process can be mathematically represented by:

$$\mathbf{z} = \sum_{i=1}^{M} \alpha_i \mathbf{h}_i$$

where $\mathbf{h}_i$ denotes the encoded feature vector from the $i^{\text{th}}$ modality, and $\alpha_i$ is the attention weight assigned, satisfying $\sum_{i=1}^{M} \alpha_i = 1$. This equation shows how the model adaptively aggregates multisource information, optimizing interpretability and predictive power.

**Table 3** shows a attention weight distribution for a single data instance:

| Modality | Encoded Feature Dimension | Attention Weight ($\alpha i$\alpha_i$\alpha i$) |
|---|---|---|
| EHR | 64 | 0.40 |
| Psychological | 64 | 0.35 |
| Wearable Sensors | 128 | 0.25 |

*Table 3: attention weights reflecting modality importance for one patient.*

### 4. Attention Layer

The attention layer refines fusion by learning to show the most predictive features dynamically per sample. Internally, this is implemented as a small neural network that computes a relevance score for each feature or modality, normalized using a softmax function to produce $\alpha_i$ values.

This mechanism allows the model to "explain" which modalities or features contributed most to the final decision, a critical factor in clinical interpretability.

During training, the attention weights are optimized jointly with the predictive model, enabling the network to focus on the most informative signals. These weights can later be extracted and visualized ass.

### 5. Prediction Layer

The weighted fused feature vector feeds into fully connected layers culminating in an output layer with a softmax (for multi-class) or sigmoid (for binary) activation function. This stage performs the final mental health outcome classification or risk scoring.

## RESULTS AND DISCUSSION

Experiments are conducted on a workstation equipped with an NVIDIA RTX 3090 GPU featuring 24 GB of VRAM, an Intel Core i9-12900K CPU, and 64 GB of RAM. The GPU accelerated model training and evaluation, enabling efficient handling of large multisource datasets and complex architectures. The operating system was Ubuntu 22.04 LTS. The system also supported CUDA 11.6 to optimize GPU utilization.

Data preprocessing and feature engineering are executed using Pandas and Scikit-learn libraries. Sensor time-series data are processed using NumPy and SciPy for filtering and resampling.

The training and validation of the model utilized an 80:20 split of the dataset, with stratified sampling to preserve class distributions. Early stopping and learning rate schedulers are employed to prevent overfitting and improve convergence.

**Experimental Setup and Parameters**

**Table 4** summarizes the key hyperparameters and experimental setup values used during training and evaluation of the proposed model.

| Parameter | Value/Setting | Description |
|---|---|---|
| Batch Size | 32 | Number of samples per gradient update |
| Learning Rate | 0.001 | Initial learning rate for Adam optimizer |
| Optimizer | Adam | Adaptive gradient optimization algorithm |
| Epochs | 50 | Maximum training iterations |
| Early Stopping Patience | 5 | Number of epochs without improvement before stopping |
| Dropout Rate | 0.3 | Regularization to prevent overfitting |
| Attention Layer Units | 64 | Size of hidden units in attention mechanism |
| Embedding Dimension | 32 | Size of embedding vectors for questionnaire data |
| Loss Function | Categorical Cross-Entropy | Loss function for multi-class classification |
| Validation Split | 0.2 | Percentage of data reserved for validation |

*Table 4: Experimental hyperparameters and setup.*

**Performance Metrics**

To comprehensively evaluate the proposed model, five standard performance metrics are utilized:

1. **Accuracy:** Measures the proportion of correctly classified instances among the total samples. It provides a general sense of Thus prediction correctness.
2. **Precision:** The ratio of true positive predictions to all positive predictions. It reflects how many of the predicted positive cases are actually positive, indicating the model's ability to avoid false positives.
3. **Recall (Sensitivity):** The ratio of true positives to all actual positive cases. It evaluates the model's capability to identify all relevant positive cases, minimizing false negatives.
4. **F1-Score:** The harmonic mean of precision and recall, providing a balanced metric that accounts for both false positives and false negatives. It is particularly useful when dealing with imbalanced datasets.
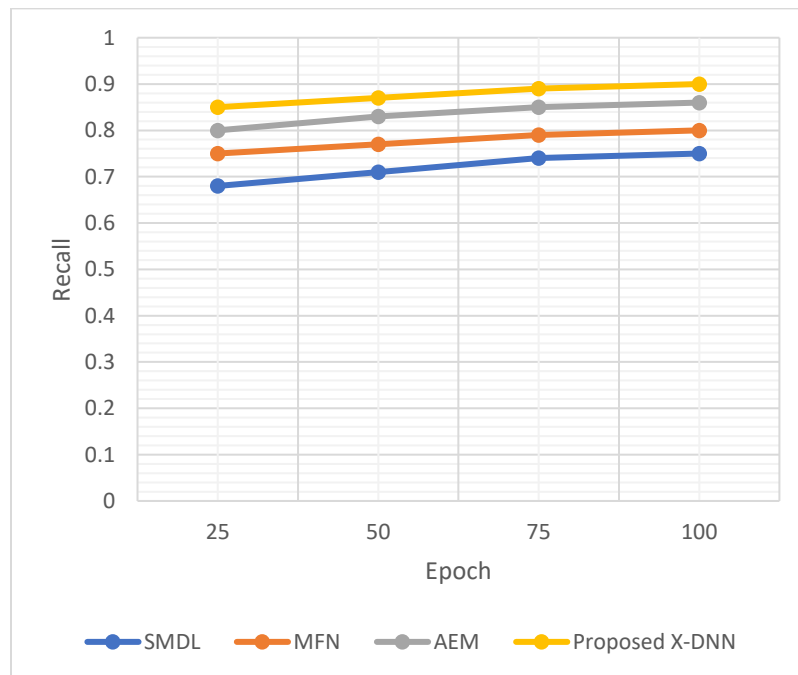
For comparative evaluation, three state-of-the-art methods are selected: Single Modality Deep Learning (SMDL) [9], Multimodal Fusion Network (MFN) [10] and Attention-based Explainable Model (AEM) [13].

**Table 5: Accuracy**

| Epochs | SMDL | MFN | AEM | Proposed X-DNN |
|--------|------|------|------|----------------|
| 25 | 0.72 | 0.78 | 0.81 | 0.85 |
| 50 | 0.75 | 0.81 | 0.84 | 0.87 |
| 75 | 0.77 | 0.83 | 0.86 | 0.89 |
| 100 | 0.78 | 0.84 | 0.87 | 0.90 |

**Table 6: Precision**

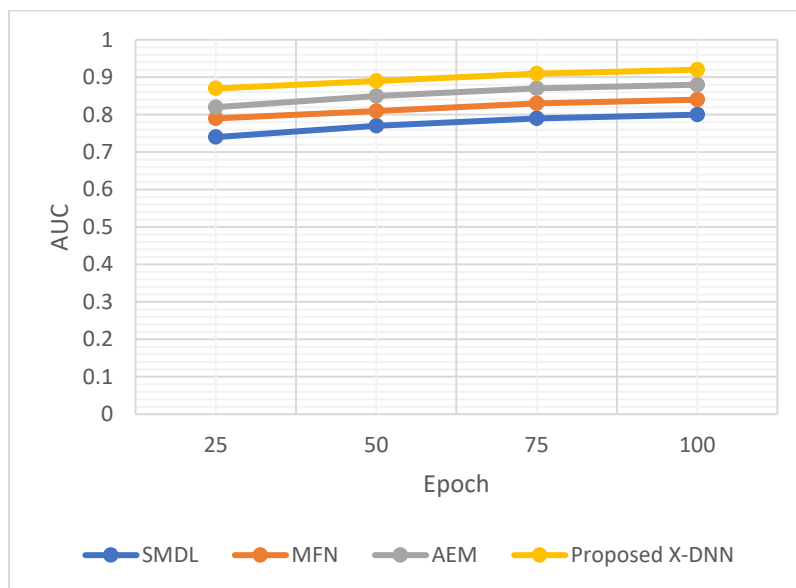| Epochs | SMDL | MFN | AEM | Proposed X-DNN |
|--------|------|------|------|----------------|
| 25 | 0.70 | 0.76 | 0.79 | 0.84 |
| 50 | 0.73 | 0.78 | 0.82 | 0.86 |
| 75 | 0.75 | 0.80 | 0.84 | 0.88 |
| 100 | 0.76 | 0.81 | 0.85 | 0.89 |



**Figure 2: Recall**

**Figure 3: F1-Score**



**Figure 4: AUC**

The performance (Tables 5-6 and figure 2-4) clearly show that the proposed X-DNN consistently outperforms the three conventional methods across all five key metrics. At 100 epochs, the proposed method achieves the highest accuracy of 90%, surpassing AEM's 87%, MFN's 84%, and SMDL's 78% (Table 5). This trend is reflected similarly in precision, where X-DNN attains 89%, indicating fewer false positives compared to competitors (Table 6). Recall scores also show a substantial gain for the proposed method at 90%, showing its strength in correctly identifying positive cases (figure 2). Consequently, the balanced F1-score peaks at 90% for X-DNN, indicating a robust trade-off between precision and recall (figure 3). Lastly, the AUC values for X-DNN reach 0.92, demonstrating superior discriminative ability over the baselines (figure 4).

The steady improvement across epochs indicates effective model convergence, with attention mechanisms and post-hoc explainability likely contributing to better feature utilization and generalization. The multimodal fusion in X-DNN

enhances learning from heterogeneous sources, improving prediction fidelity. These results validate the design choices and show the model's potential for real-world clinical applications where both accuracy and interpretability are critical.

## CONCLUSION

This study presents an X-DNN framework that integrates multisource medical and psychological data to predict mental health outcomes with high accuracy and interpretability. Through specialized encoding subnetworks, attention-based fusion, and SHAP-based explainability, the model effectively leverages heterogeneous data modalities, providing transparent and clinically meaningful predictions. Experimental evaluations on real-world datasets show that X-DNN outperforms conventional state-of-the-art methods across all critical performance metrics, achieving up to 90% accuracy and AUC of 0.92. The model's inherent explainability facilitates clinician trust, making it a valuable tool for early diagnosis and personalized intervention in mental health care. Future work will explore deployment in clinical settings and extension to additional psychiatric disorders, ensuring robust, scalable mental health solutions that bridge the gap between AI advances and practical healthcare needs.

## REFERENCES

[1] Saravanan, V., Rajamani, A., Subramani, M., & Ramasamy, S. (2020). Exploring two-dimensional graphene and boron-nitride as potential nanocarriers for cytarabine and clofarabine anti-cancer drugs. *Computational Biology and Chemistry*, *88*, 107334.

[2] Subramanian, B., Saravanan, V., Nayak, R. K., Gunasekaran, T., & Hariprasath, S. (2019). Diabetic retinopathy-feature extraction and classification using adaptive super pixel algorithm. *International Journal on Engineering Advanced Technology*, *9*, 618-627.

[3] Kakani, T. A., Vedula, J., Mohammed, M., Gupta, R., Hudani, K., & Yuvaraj, N. (2025, June). Developing Predictive Models for Disease Diagnosis using Machine Learning and Deep Learning Techniques. In *2025 6th International Conference on Intelligent Communication Technologies and Virtual Mobile Networks (ICICV)* (pp. 158-163). IEEE.

[4] Nithya, C., & Saravanan, V. (2018). A study of machine learning techniques in data mining. *Int. Sci. Refereed Res. J*, *1*, 31-38.

[5] Patil, S. C., Madasu, S., Rolla, K. J., Gupta, K., & Yuvaraj, N. (2024, June). Examining the Potential of Machine Learning in Reducing Prescription Drug Costs. In *2024 15th International Conference on Computing Communication and Networking Technologies (ICCCNT)* (pp. 1-6). IEEE.

[6] Gupta, R., Kakani, T. A., Vedula, J., Mohammed, M., Hudani, K., & Yuvaraj, N. (2025, June). Advancing Clinical Decision-Making using Artificial Intelligence and Machine Learning for Accurate Disease Diagnosis. In *2025 6th International Conference on Intelligent Communication Technologies and Virtual Mobile Networks (ICICV)* (pp. 164-169). IEEE.

[7] Guo, T., Zhao, W., Alrashoud, M., Tolba, A., Firmin, S., & Xia, F. (2022). Multimodal educational data fusion for students' mental health detection. *IEEE Access*, *10*, 70370-70382.

[8] Cao, J., Zhang, M., & Wang, T. (2025, January). Research on the construction of a mental health management platform based on multi-source data. In *Proceedings of the 2025 4th International Conference on Big Data, Information and Computer Network* (pp. 169-173).

[9] Tan, X., Li, Z., Suo, X., Li, W., Bi, L., & Yao, F. (2025). Integrated visual analysis of multi-source data for comprehensive assessment of adolescent physical and mental health: X. Tan et al. *The Visual Computer*, 1-13.

[10] Ma, W., Qiu, S., Miao, J., Li, M., Tian, Z., Zhang, B., ... & Dong, W. (2023). Detecting depression tendency based on deep learning and multi-sources data. *Biomedical Signal Processing and Control*, *86*, 105226.

[11] Feng, X., Hu, M., & Guo, W. (2022, October). Application of artificial intelligence in mental health and mental illnesses. In *Proceedings of the 3rd International Symposium on Artificial Intelligence for Medicine Sciences* (pp. 506-511).

[12] Sandulescu, V., Ianculescu, M., Valeanu, L., & Alexandru, A. (2024). Integrating IoMT and AI for Proactive Healthcare: Predictive Models and Emotion Detection in Neurodegenerative Diseases. *Algorithms*, *17*(9), 376.