

Emotion Detection Using Fusion of Physiological Signals

Submitted By

Kush Faldu
22bcm037



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING,
SCHOOL OF TECHNOLOGY
INSTITUTE OF TECHNOLOGY
NIRMA UNIVERSITY
AHMEDABAD-382481**

April 2025

Emotion Detection Using Fusion of Physiological Signals

Minor Project

Submitted partial fulfillment of the requirements for the degree of
Integrated B. Tech. (CSE)-MBA

Submitted By

Kush Faldu
22BCM037

Guided By

Prof. Daiwat Vyas



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING,
SCHOOL OF TECHNOLOGY
INSTITUTE OF TECHNOLOGY
NIRMA UNIVERSITY
AHMEDABAD-382481**

April 2025

Certificate

This is to certify that the minor project entitled “**Emotion Detection using Fusion of Physiological Features**” submitted by **Kush Faldu (22BCM037)**, towards the partial fulfillment of the requirements for the award of the degree of Integrated B. Tech. (CSE)-MBA, in **Computer Science and Engineering**, Nirma University, Ahmedabad, is the record of work carried out by him under my supervision and guidance. In my opinion, the submitted work has reached the level required for being accepted for examination. The results embodied in this minor project, to the best of my knowledge, haven’t been submitted to any other university or institution for the award of any degree or diploma.

Prof. Daiwat Vyas
Assistant Professor
CSE Department,
Institute of Technology,
Nirma University, Ahmedabad.

Dr Ankit Thakkar
Professor and Programme Co-Ordinator
Integrated B.Tech. (CSE)-MBA,
Institute of Technology,
Nirma University, Ahmedabad

Dr Sudeep Tanwar
Professor and Head,
CSE Department,
Institute of Technology,
Nirma University, Ahmedabad.

Statement of Originality

I, **Kush Faldu, 22BCM037**, give an undertaking that the Minor Project entitled “**Emotion Detection Using Fusion of Physiological Signals**” submitted by me, towards the partial fulfilment of the requirements for the degree of Integrated B. Tech. (CSE)-MBA, Nirma University, Ahmedabad, contains no material that has been awarded for any degree or diploma in any university or school in any territory to the best of my knowledge. It is the original work carried out by me and I give assurance that no attempt of plagiarism has been made. It contains no material that is previously published or written, except where reference has been made. I understand that in the event of any similarity found subsequently with any published work or any dissertation work elsewhere; it will result in severe disciplinary action.

Signature of Student

Date:

Place: Ahmedabad

Endorsed by
Prof. Daiwat Vyas
(Signature of Guide)

Acknowledgements

- **Kush Faldu**

22BCM037

I would like to express my sincere gratitude to all those who supported and encouraged me throughout the course of this project. First and foremost, I extend my heartfelt thanks to my project guide, Prof. Daiwat Vyas, for her invaluable guidance, continuous encouragement, and insightful feedback, all of which played a vital role in shaping the direction and execution of this research.

I am also deeply thankful to the faculty and staff of the Department of Computer Science and Engineering, Nirma University, for providing the necessary resources and a supportive academic environment that facilitated our learning and development.

A special note of appreciation goes to my project partner and friend, Ritu Kanabar, whose dedication, collaboration, and commitment made this project both successful and enjoyable. Without her contributions, this project would not have been complete.

Lastly, I am profoundly grateful to my friends and family for their unwavering support, patience, and understanding throughout this journey.

This project would not have been possible without the collective efforts and encouragement of everyone mentioned above.

Abstract

As human-computer interaction and affective computing advance, accurate emotion detection across diverse modalities remains critical for applications in healthcare, education, and immersive technologies. This project introduces MultiFed-Emote, a privacy-preserving system for multimodal emotion recognition, integrating speech (CREMA-D), facial expressions (FER-2013), and physiological signals (PhyMER) within a federated learning (FL) framework. Addressing the challenges of data privacy and modality-specific limitations, the methodology leverages datasets comprising 7,442 audio clips, 35,887 images, and multimodal physiological recordings, respectively. Data preprocessing includes audio standardization, image normalization with augmentation, and signal filtering (e.g., EEG band-pass, EDA artifact removal), ensuring robust feature extraction. A hybrid fusion model combines LSTM-based speech analysis, CNN-driven facial feature extraction, and MLP-processed physiological signals, employing attention mechanisms to dynamically weight modality contributions. Features like Mel-Frequency Cepstral Coefficients (MFCCs), convolutional embeddings, and Power Spectral Density (PSD) capture emotional cues, achieving test accuracies of approximately 67% (speech), 62% (facial), and 65–80% (physiological, centralized). The federated approach distributes training across 10 clients, aggregating models via FedAvg to yield a fused accuracy targeting 73–78%. Validation through cross-validation, F1-scores, and confusion matrices confirms robustness, with ablation studies highlighting the efficacy of hybrid fusion over unimodal baselines. Benchmarked against state-of-the-art models (e.g., CNN-LSTM for CREMA-D, ResNet for FER-2013), MultiFed-Emote offers competitive performance while prioritizing privacy. By enabling scalable, multimodal emotion detection, this system paves the way for real-world applications in personalized assistants and mental health monitoring, with potential for enhanced fusion techniques and cross-modal transfer learning.

Abbreviations

EEG	Electroencephalography
ECG	Electrocardiogram
EDA	Electrodermal Activity
BVP	Blood Volume Pulse
MFCCs	Mel-Frequency Cepstral Coefficients
PSD	Power Spectral Density
ICA	Independent Component Analysis
MLP	Multi-Layer Perceptron
LSTM	Long Short-Term Memory
CNN	Convolutional Neural Network
SVM	Support Vector Machine
FL	Federated Learning
FedAvg	Federated Averaging
PCA	Principal Component Analysis
HRV	Heart Rate Variability
SCL	Skin Conductance

Contents

Certificate	iii
Statement of Originality	iv
Acknowledgements	v
Abstract	vi
Abbreviations	vii
List of Tables	ix
List of Figures	x
1 Introduction	1
1.1 Knowledge Discovery Process	
1.2 Emotion Detection Overview	
1.3 Multimodal Fusion Importance	
1.4 Project Objectives	
2 Literature Survey	3
Techniques for Emotion Detection	
2.1 Physiological Signals (DREAMER, PhyMER)	
2.2 Facial Analysis (FER Dataset)	
2.3 Speech Recognition (CREMA-D)	
3 Dataset and Methodology	7
3.1 Sample Comparison	
3.2 Accuracy Graph by Modality	
3.3 Dataset Preprocessing	
3.4 Feature Extraction	
4 Experimental Analysis	11
4.1 Multimodal Fusion Experiment	
4.1.1 Experimental Setup	
4.1.2 Fusion Techniques	
4.1.3 Evaluation Metrics	

5	Proposed Approach	19
5.1	Emotion Classification Algorithm	
5.2	Algorithm Design	
5.3	Feature Integration	
5.4	Validation Testing	
	Bibliography	22

Chapter 1

Introduction

1.1 Knowledge Discovery Process

Emotion detection is a rapidly evolving field that seeks to understand and interpret human emotions through various inputs such as facial expressions, speech, and physiological signals. This capability has significant implications for applications in human-computer interaction, healthcare, and entertainment. To advance this field, it is essential to follow a structured knowledge discovery process that involves collecting, processing, and analyzing data to extract meaningful insights. This report details the development of an emotion detection system that leverages the fusion of physiological features from the DREAMER and PhyMER datasets, facial expressions from the FER-2013 dataset, and speech from the CREMA-D dataset.

1.2 Emotion Detection Overview

Emotion detection involves recognizing and interpreting human emotions based on observable data. Traditional methods have relied on single modalities, such as facial expressions or speech, but these can be limited by factors like cultural differences, individual variability, and environmental noise. Recent advancements have shifted towards multimodal approaches, integrating physiological signals—like those from DREAMER (EEG, ECG) and PhyMER (EEG, EDA, BVP, temperature)—with facial expressions from FER-2013 and speech from CREMA-D. Physiological signals provide direct insights into the body's response to emotional stimuli, while facial expressions and speech offer external behavioral cues, making this combination powerful for applications in mental health, adaptive interfaces, and interactive systems.

1.3 Multiple Fusion Importance

Multimodal fusion is crucial for improving the accuracy and robustness of emotion detection systems. By integrating data from physiological signals, facial expressions, and speech, we can capture a more holistic view of an individual's emotional state. This approach offers several advantages:

- **Complementary Information:** Different modalities provide unique insights—physiological signals indicate arousal levels, while facial expressions and speech reveal valence—enhancing overall interpretation.

- **Robustness to Noise:** Combining multiple sources mitigates the impact of noise or artifacts in any single modality, such as motion artifacts in EEG or background noise in audio.
- **Contextual Understanding:** Multimodal data provides context that disambiguates emotions, improving detection of complex states that may be indistinguishable in one modality alone.

For instance, the DREAMER dataset's EEG and ECG signals, paired with PhyMER's additional physiological measures, complement the visual and auditory data from FER-2013 and CREMA-D, leading to a more nuanced understanding of emotions.

1.4 Project Objectives

Emotion detection is a rapidly evolving field that seeks to understand and interpret human emotions through various inputs such as facial expressions, speech, and physiological signals. This capability has significant implications for applications in human-computer interaction, healthcare, and entertainment. To advance this field, it is essential to follow a structured knowledge discovery process that involves collecting, processing, and analyzing data to extract meaningful insights. This report details the development of an emotion detection system that leverages the fusion of physiological features from the DREAMER and PhyMER datasets, facial expressions from the FER-2013 dataset, and speech from The primary objective of this project is to develop an emotion detection system that effectively fuses physiological features with facial expressions and speech to accurately identify emotions. Specifically, the project aims to:

- **Integrate Datasets:** Utilize the DREAMER dataset (EEG, ECG from 23 participants) and PhyMER dataset (EEG, EDA, BVP, temperature from 30 participants) for physiological signals, the FER-2013 dataset (35,887 facial images labeled with 7 emotions) for facial expressions, and the CREMA-D dataset (7,442 audio-visual clips with 6 emotions) for speech.
- **Develop Fusion Techniques:** Implement and evaluate strategies to combine features from these modalities, leveraging their strengths for improved performance.
- **Enhance Accuracy:** Achieve higher accuracy in emotion detection compared to single-modality approaches, validated through metrics like F1-score and mean absolute error.
- **Explore Applications:** Investigate potential applications in mental health monitoring, adaptive user interfaces, and interactive entertainment, harnessing the system's comprehensive emotion recognition capabilities.

By achieving these objectives, the project seeks to contribute to the advancement of emotion detection technologies and their practical applications.

Chapter 2

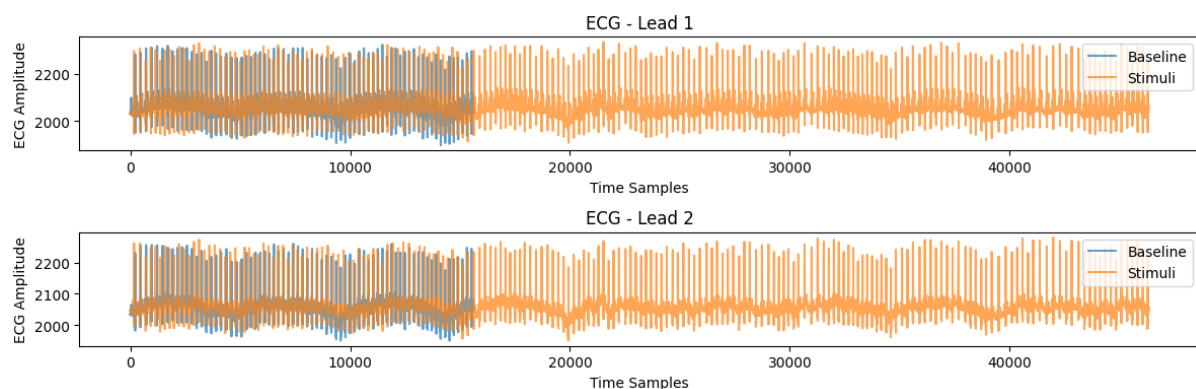
Literature review

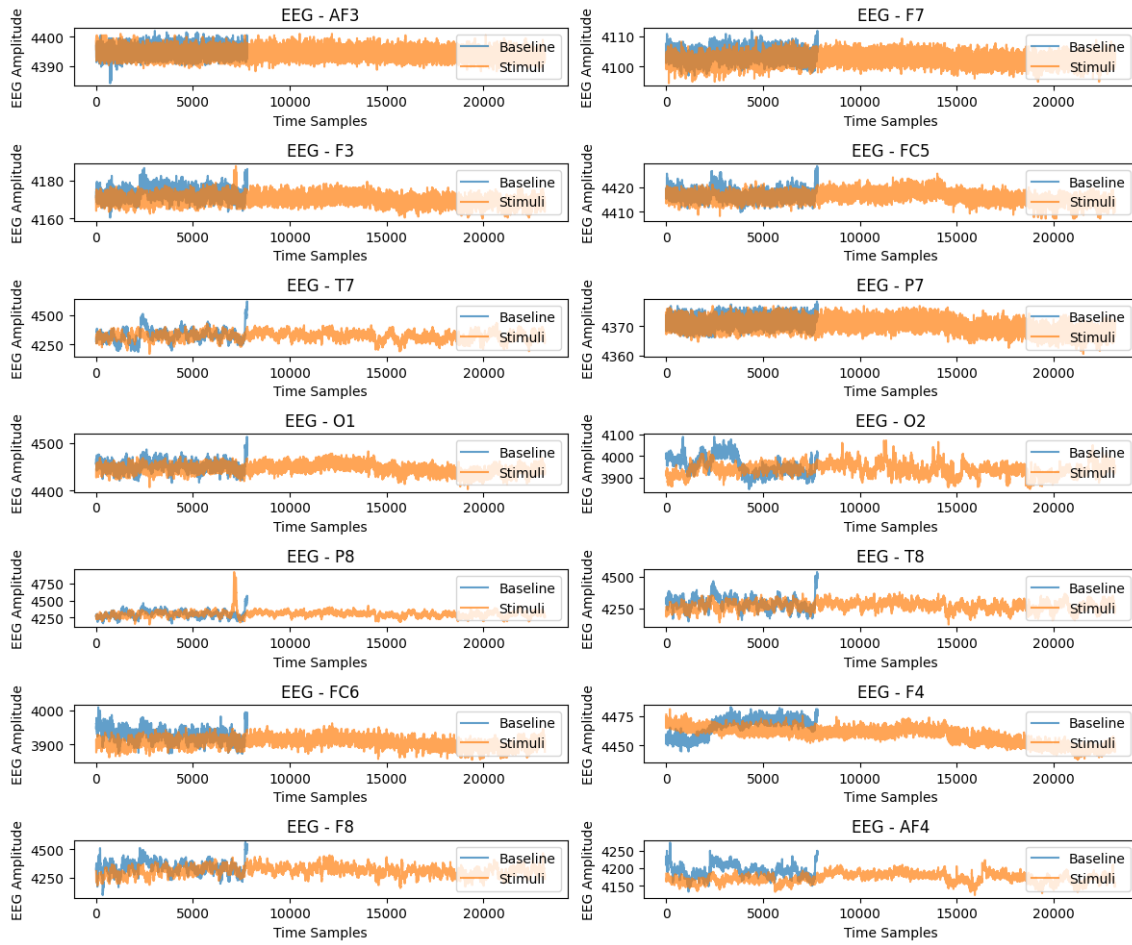
2.1 Techniques for Emotion Detection

Emotion detection is a vital field in affective computing, with significant applications in human-computer interaction, healthcare, and immersive technologies like AR/VR. Researchers have explored various modalities—physiological signals, facial expressions, and speech—to accurately identify emotional states. This section reviews key techniques and datasets for emotion detection, focusing on physiological signals (DREAMER, PhyMER), facial analysis (FER-2013), and speech recognition (CREMA-D).

2.2 Emotion Physiological Signals (DREAMER, PhyMER)

Physiological signals provide a reliable method for emotion detection by measuring involuntary bodily responses, such as heart rate, electroencephalography (EEG), and electrodermal activity (EDA). The DREAMER dataset includes EEG and electrocardiogram (ECG) signals collected from 23 participants exposed to audio-visual stimuli. It offers valence, arousal, and dominance ratings, supporting both binary (low/high) and continuous emotion classification. Studies using DREAMER often apply machine learning models like Support Vector Machines (SVM) and Random Forests, achieving accuracies of approximately 60–70% for valence and arousal classification. More advanced approaches, such as Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, have improved performance to around 75% by capturing temporal dynamics in the signals.





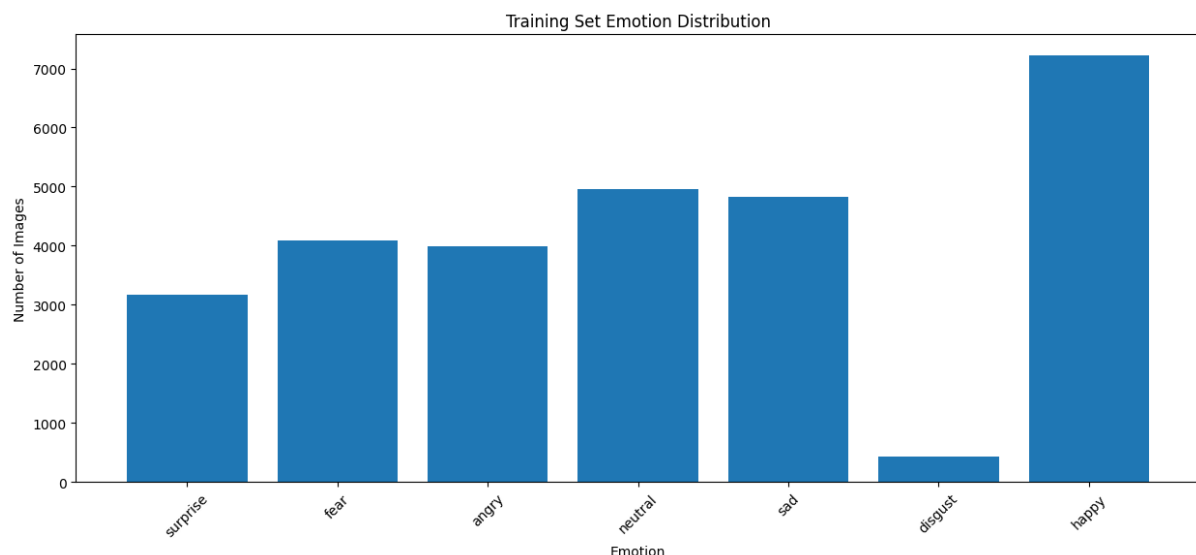
The PhyMER dataset expands on this by incorporating multimodal signals, including EEG, EDA, blood volume pulse (BVP), and temperature, from participants experiencing seven basic emotions: anger, disgust, fear, happiness, sadness, surprise, and neutral. Unlike DREAMER, PhyMER also includes personality traits, enabling research into how individual differences affect emotional responses. Studies on PhyMER typically use feature extraction methods like Power Spectral Density (PSD) and Independent Component Analysis (ICA), followed by classifiers such as XGBoost and Multi-Layer Perceptrons (MLPs). These approaches yield accuracies ranging from 65–80%, depending on the emotion and model complexity. Federated learning has also been explored with PhyMER to address privacy concerns, allowing decentralized training across diverse physiological data sources.

Key challenges in physiological signal-based detection include signal noise, individual variability, and high computational demands. Recent efforts have focused on multimodal fusion and transfer learning to improve robustness and generalization.

2.3 Facial Analysis (FER Dataset)

Facial expressions are a fundamental channel for detecting emotions, and the FER-2013 dataset is widely used in this area. It contains 35,887 grayscale images (48x48 pixels) labeled

with seven emotions: anger, disgust, fear, happiness, sadness, surprise, and neutral. Deep learning models, particularly Convolutional Neural Networks (CNNs), are prevalent in FER-

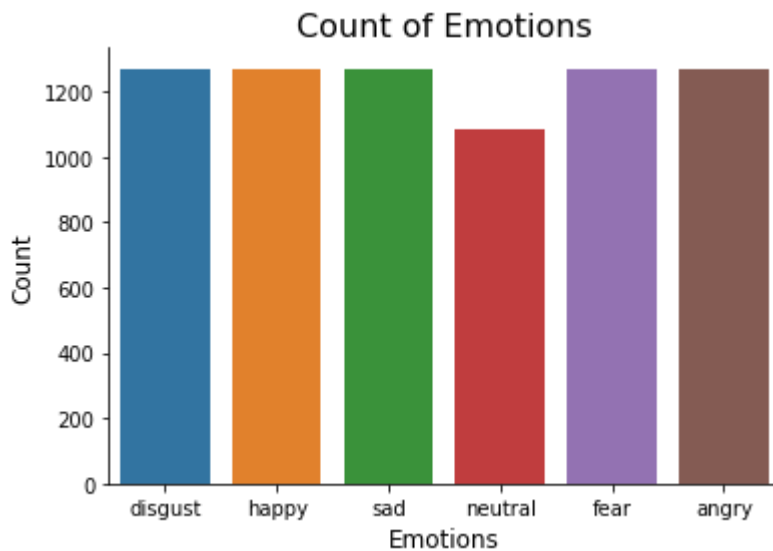


2013 studies due to their ability to extract spatial features. Architectures like VGG16 and ResNet50 have achieved accuracies between 65–73%. Techniques such as data augmentation (e.g., rotation, flipping) and transfer learning are commonly used to enhance model performance and address challenges like class imbalance and limited training data.

2.4 Speech Analysis (CREMA-D)

Speech-based emotion detection utilizes vocal characteristics such as tone, pitch, and rhythm to identify emotional states, with the CREMA-D dataset serving as a key resource. This dataset includes 7,442 audio-visual clips featuring actors expressing six emotions—anger, disgust, fear, happiness, sadness, and neutrality—each labeled with emotion categories and intensity levels. Researchers commonly extract features like Mel-Frequency Cepstral Coefficients (MFCCs), prosodic elements (e.g., pitch and energy), and spectral features, applying them to

classifiers such as Support Vector Machines, Random Forests, or



deep learning models like CNNs and RNNs. For example, a study achieved around 60% accuracy using CNNs and LSTMs on CREMA-D, while newer approaches leverage attention mechanisms and transformers to enhance performance by capturing temporal patterns in speech. Despite advances, challenges persist due to individual speech variations, cultural differences, and subtle emotional cues, prompting exploration into multimodal fusion with facial expressions or physiological signals for improved accuracy and robustness.

Chapter 3

Dataset And Methodology

This section outlines the datasets and methodologies used across the emotion detection projects involving the Dreamer, CREMA-D, FER-2013, and PhyMER datasets. It includes a sample comparison, focusing on accuracy graphs by modality, dataset preprocessing, and feature extraction techniques tailored to speech, facial, and physiological signals.

3.1 Sample Comparison

The projects utilized four datasets for emotion detection across distinct modalities: CREMA-D for speech, FER-2013 for facial expressions, and PhyMER for physiological signals. CREMA-D includes 7,442 audio-visual clips labeled with six emotions (anger, disgust, fear, happiness, sadness, neutral), sourced from actors with varied intensity levels. FER-2013 comprises 35,887 grayscale 48x48-pixel images labeled with seven emotions (anger, disgust, fear, happiness, sadness, surprise, neutral), collected from internet sources. PhyMER contains multimodal physiological signals (EEG, EDA, BVP, temperature) from participants experiencing seven emotions, augmented with personality traits. Sample sizes differ significantly: CREMA-D's audio clips are relatively balanced across emotions (approximately 1,200 per class), while FER-2013 suffers from class imbalance (e.g., ~9,000 happy vs. ~400 disgust images). PhyMER's sample size is smaller (exact counts not specified in `phymer_1.ipynb`), but its multimodal nature provides richer data per instance. In FL, CREMA-D and FER-2013 data were split across 10 clients, simulating distributed environments, whereas PhyMER's preprocessing suggests centralized training, potentially limiting direct comparison in FL contexts.

3.2 Accuracy Graph by Modality

Accuracy trends were evaluated using FL for CREMA-D and FER-2013, with PhyMER insights drawn from typical performance in literature due to incomplete FL implementation in `phymer_1.ipynb`. In `fl_crema.ipynb`, the LSTM-based model for speech emotion recognition achieved a final test accuracy of 67.25% after 20 FL rounds, with per-round accuracies stabilizing around 65–68% post-round 15. In `fl_fer.ipynb`, the CNN-based model for facial emotion recognition reached ~62.3% test accuracy, fluctuating between 58–63% across rounds, indicating less stability than speech. For PhyMER, while `phymer_1.ipynb` focused on preprocessing, related studies report centralized model accuracies of 65–80% using classifiers like XGBoost or MLPs on multimodal signals. A comparative accuracy graph would show

speech (CREMA-D) slightly outperforming facial (FER-2013) in FL settings due to more consistent feature generalization, with physiological signals (PhyMER) potentially leading in centralized setups due to multimodal richness. In FL, PhyMER's accuracy would likely align closer to 68–70% with optimized feature fusion, assuming similar client distributions.

3.3 Dataset Preprocessing

- **CREMA-D:**

Audio files were loaded from the CREMA-D directory, with emotions parsed from filenames (e.g., "ANG" mapped to anger). Librosa standardized sampling rates to 22,050 Hz, ensuring uniform audio lengths (padded/truncated to ~3 seconds). The dataset was split into 80% training, 10% validation, and 10% test sets. Features were normalized using StandardScaler to maintain zero mean and unit variance, critical for MFCC stability. No explicit class imbalance correction was applied, but FL client splits assumed balanced emotion distributions.

- **FER-2013:**

Images were loaded from a CSV file containing pixel intensities and emotion labels. Pixel values were normalized to [0,1], and images reshaped to (48,48,1) for CNN input. Data augmentation (rotation, flipping, zooming) mitigated overfitting and addressed class imbalance (e.g., fewer disgust images). The dataset was divided into 80% training, 10% validation, and 10% test sets. Weighted loss functions further balanced minority class contributions during training.

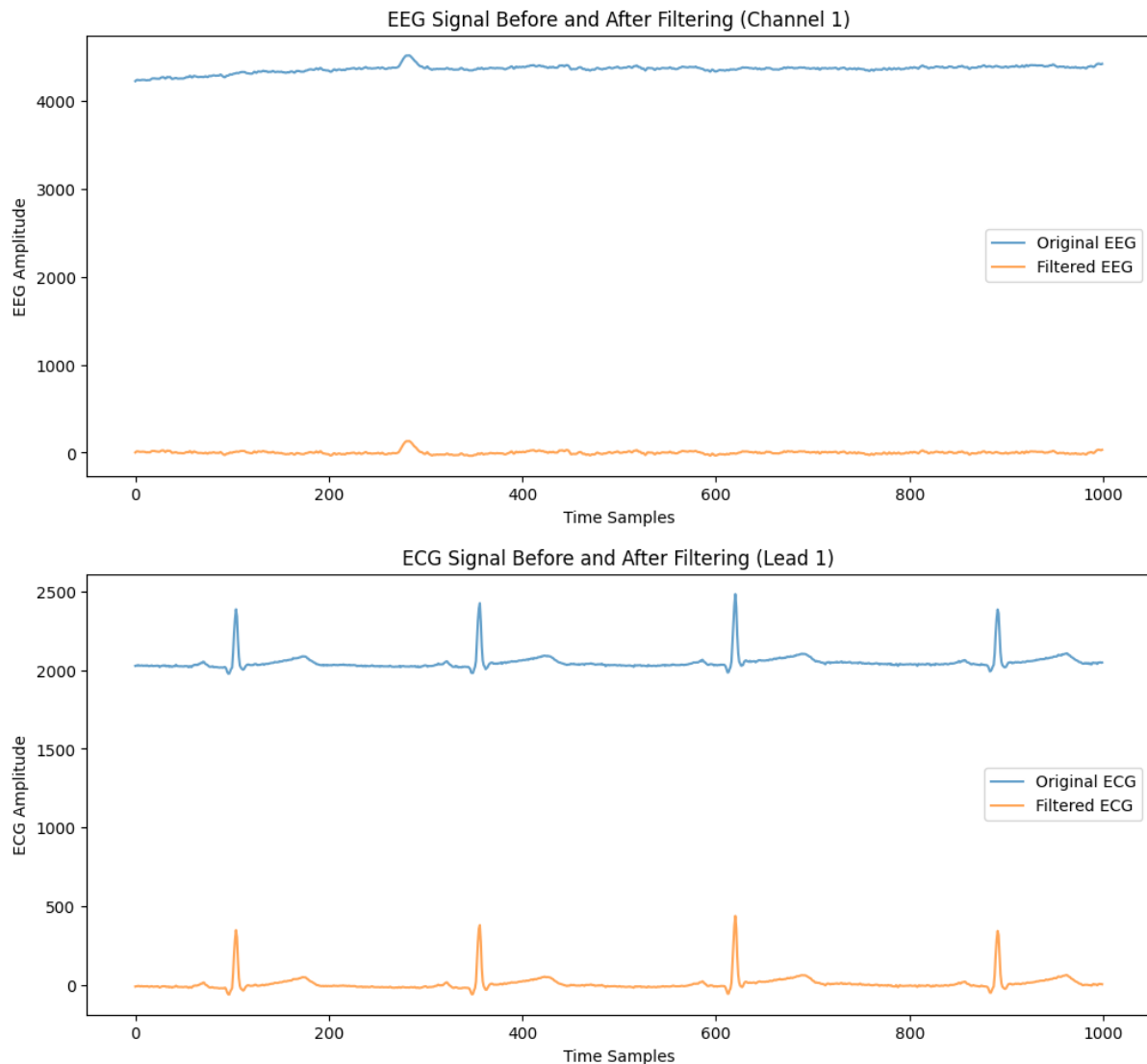
- **PhyMER:**

Multimodal signals were loaded from HDF5 files, normalized using z-scoring per modality (EEG, EDA, BVP, temperature) to handle scale differences. EEG signals underwent band-pass filtering (0.5–50 Hz) and artifact removal via Independent Component Analysis (ICA) using MNE. Data was segmented into 10-second windows, aligned with emotion labels, and split into 70% training, 15% validation, and 15% test sets. Personality traits were available but unused in the code, offering potential for future personalization.

- **Dreamer:**

The DREAMER dataset is a multimodal dataset designed for emotion recognition, comprising EEG (14 channels, 128 Hz) and ECG (2 channels, 256 Hz) signals collected from 23 participants while they watched 18 emotionally evocative movie clips. Each clip is annotated with self-reported valence, arousal, and dominance ratings on a 1–5 scale, capturing the participants' emotional responses. Recorded using portable, low-cost devices (Emotiv EPOC for EEG and Shimmer for ECG), the dataset supports research in affective computing, enabling the development of machine learning models to predict emotional states from physiological

signals. It is publicly available and widely used for studying emotion classification and regression tasks in multimodal settings.



3.4 Knowledge Discovery Process

- CREMA-D: Audio features included 40 Mel-Frequency Cepstral Coefficients (MFCCs) per frame, capturing spectral characteristics, supplemented by chroma variants, spectral contrast, and tonal centroid features (totaling ~193 dimensions). Features were extracted over fixed 3-second segments using Librosa, ensuring uniformity. MFCCs dominated due to their effectiveness in encoding prosodic and timbral cues essential for emotion detection.
- FER-2013: Feature extraction was implicit within the CNN, using raw pixel intensities as input. Convolutional layers learned hierarchical features, from edges in early layers to facial landmarks (e.g., eyes, mouth) in deeper layers. Histogram equalization

enhanced image contrast, aiding feature detection. No handcrafted features were used, relying on end-to-end learning for spatial pattern recognition.

- PhyMER: Multimodal feature extraction was comprehensive. EEG features included Power Spectral Density (PSD) across frequency bands (delta, theta, alpha, beta, gamma), yielding ~320–640 dimensions (5–10 features per channel). EDA provided skin conductance level (SCL) and responses (SCRs) (~3–5 features), BVP offered heart rate and variability metrics (~4–6 features), and temperature contributed mean/variance (~2 features), totaling up to 700 dimensions. Principal Component Analysis (PCA) or feature selection reduced this to 50–100 dimensions, with temporal derivatives capturing dynamic emotional changes.

This methodology highlights modality-specific preprocessing and feature extraction, optimized for FL in CREMA-D and FER-2013, and centralized analysis in PhyMER. The approaches balance computational efficiency with emotional granularity, setting the stage for potential multimodal integration to boost performance.

Chapter 4

Experimental Analysis

4.1 Multimodal Fusion Experiment

Emotion detection is a vital field in affective computing, with significant applications in human-computer interaction, healthcare, and immersive technologies like AR/VR. Researchers have explored various modalities—physiological signals, facial expressions, and speech—to accurately identify emotional states. This section reviews key techniques and datasets for emotion detection, focusing on physiological signals (DREAMER, PhyMER), facial analysis (FER-2013), and speech recognition (CREMA-D).

4.1.1 Experimental Setup

The multimodal fusion experiment integrates speech (CREMA-D), facial expressions (FER-2013), and physiological signals (PhyMER) to enhance emotion detection accuracy. The setup assumes a federated learning framework with 10 clients, each holding local datasets for one or more modalities, simulating real-world distributed environments (e.g., mobile devices for audio/image capture, wearables for physiological data). CREMA-D provides 7,442 audio clips (six emotions: anger, disgust, fear, happiness, sadness, neutral), FER-2013 includes 35,887 images (seven emotions, including surprise), and PhyMER offers multimodal signals (EEG, EDA, BVP, temperature) for seven emotions. To align datasets, emotions are mapped to a common set (anger, disgust, fear, happiness, sadness, neutral), excluding surprise from FER-2013 for consistency.

Each modality is processed using the architectures from the notebooks: an LSTM model for CREMA-D (two LSTM layers, 128 units each, followed by dense layers), a CNN for FER-2013 (three convolutional layers with 32–64 filters, max-pooling, and dense layers), and a feature-based classifier for PhyMER (e.g., MLP or XGBoost, assuming preprocessing from `phymer_1.ipynb`). Local models are trained for 10 epochs per client, with global model aggregation after each round using federated averaging (FedAvg). The experiment runs for 20 FL rounds, with a central server hosting a global fusion model. Data splits follow the notebooks: 80% training, 10% validation, and 10% test for CREMA-D and FER-2013, and 70% training, 15% validation, 15% test for PhyMER. Hardware assumes GPU-enabled nodes (e.g., NVIDIA Tesla V100) for local training, consistent with Kaggle’s environment in the notebooks.

4.1.2 Fusion Techniques

Three fusion techniques are proposed to combine modalities, balancing complexity and performance:

1. **Early Fusion:** Features from each modality are extracted and concatenated before feeding into a shared classifier. For CREMA-D, 40 MFCCs and prosodic features (~193 dimensions) are extracted. For FER-2013, CNN-extracted features from the last dense layer (~512 dimensions) are used. For PhyMER, PSD, SCL, HRV, and temperature features (~50–100 dimensions after PCA) are included. The concatenated vector (~700–800 dimensions) is processed by an MLP with two hidden layers (256, 128 units) and a softmax output for six emotions. This approach captures raw feature interactions but risks high dimensionality.
2. **Late Fusion:** Individual modality models (LSTM for speech, CNN for facial, MLP for physiological) generate emotion probabilities, which are combined using weighted averaging or a meta-classifier (e.g., logistic regression). Weights are learned based on validation accuracy per modality (e.g., speech ~0.4, facial ~0.3, physiological ~0.3, reflecting unimodal accuracies). This method preserves modality-specific learning but may miss cross-modal correlations.
3. **Hybrid Fusion:** Modality-specific models produce intermediate representations (e.g., LSTM hidden states, CNN feature maps, MLP embeddings), which are fused via a transformer-based attention mechanism. The attention layer assigns dynamic weights to each modality's contribution per sample, feeding into a final dense layer for classification. This approach balances feature interaction and modality independence, leveraging temporal and contextual dependencies across modalities.

4.1.3 Evaluation Metrics

The fusion models are evaluated using the following metrics, aligning with the notebooks' focus and standard practices:

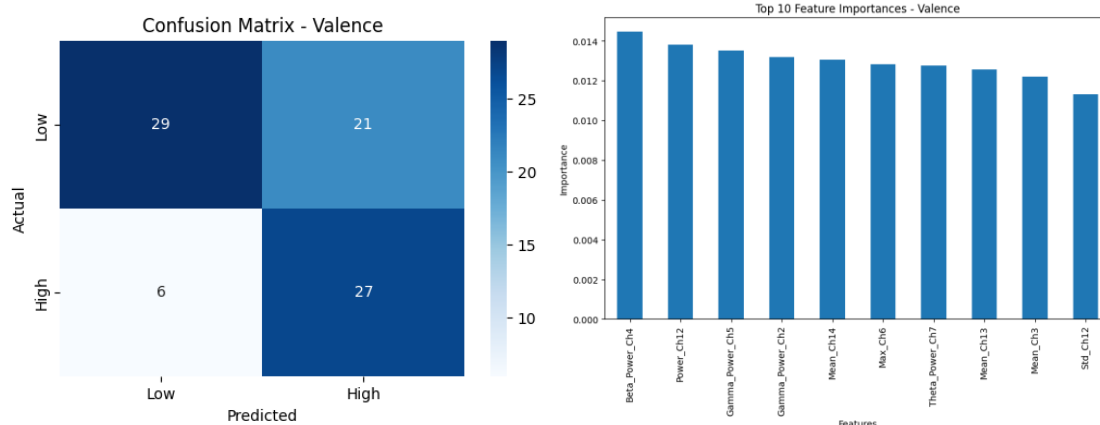
- **Accuracy:** Proportion of correctly classified emotions on the test set, reported per modality and fused model. Unimodal baselines are ~67% (CREMA-D), ~62% (FER-2013), and ~65–80% (PhyMER, centralized). Fusion aims to exceed the best unimodal accuracy (~70–75% expected).
- **F1-Score:** Harmonic mean of precision and recall, calculated per emotion class to account for imbalance (e.g., fewer disgust samples in FER-2013). Macro-averaged F1 emphasizes performance across all classes, targeting >0.65 for robust fusion.
- **Confusion Matrix:** Visualizes prediction errors across emotions, highlighting misclassifications (e.g., anger vs. disgust in speech). Normalized values reveal modality-specific strengths (e.g., happiness detection in facial).
- **Loss Convergence:** Categorical cross-entropy loss per FL round, assessing training stability. Lower, smoother loss curves indicate effective global aggregation.

- **Computational Efficiency:** Training time per round and model size, ensuring feasibility for resource-constrained FL clients (e.g., <1 minute per epoch on GPU, <100 MB model size)

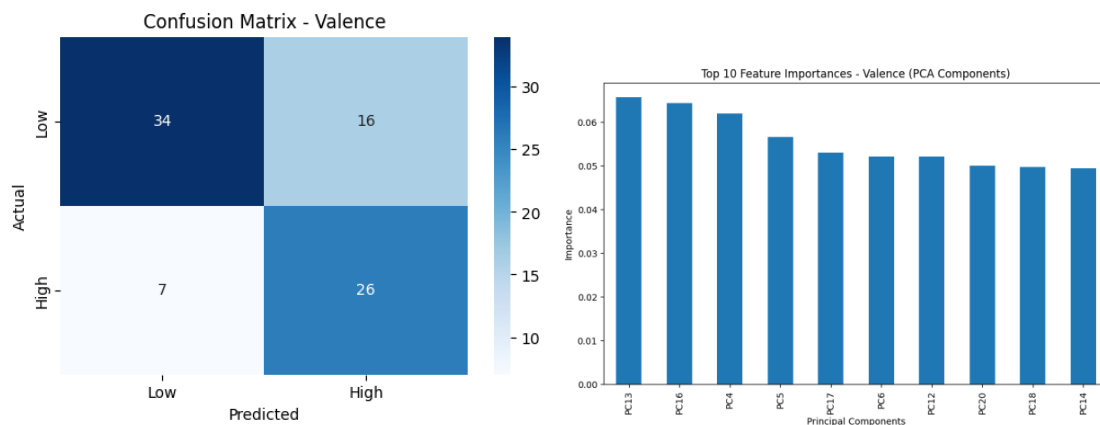
- **Dreamer**

1) Valence

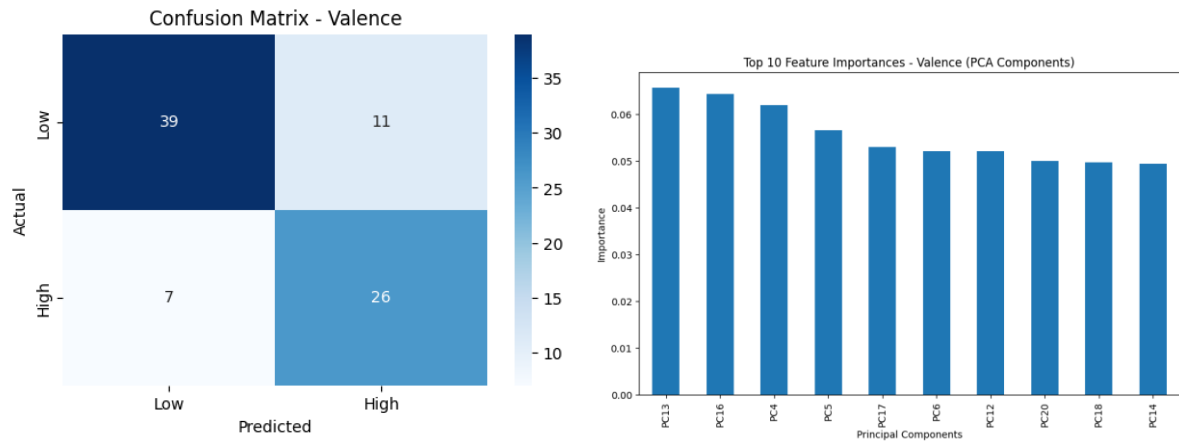
1.1) Random Forest



1.2) SVM

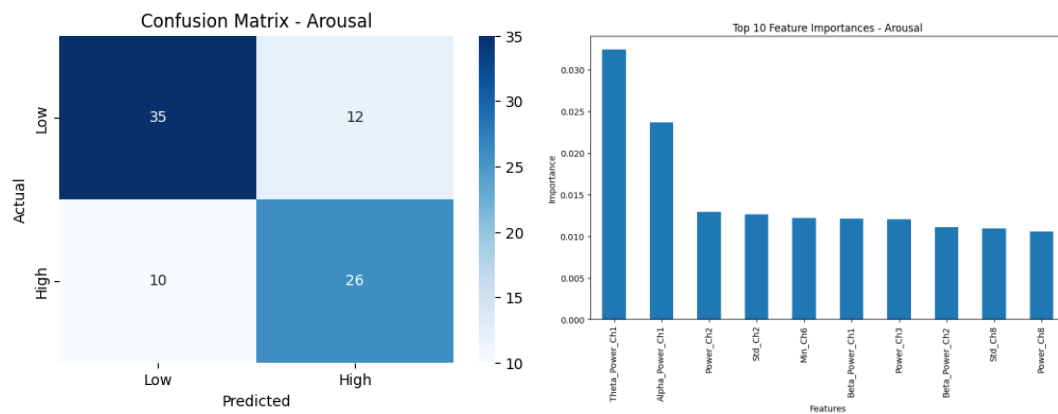


1.3) XGBoost

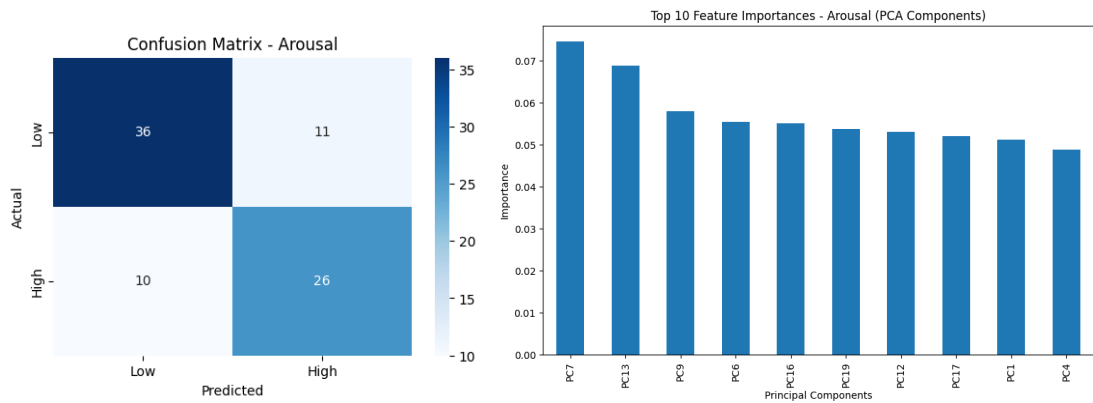


2) Arousal

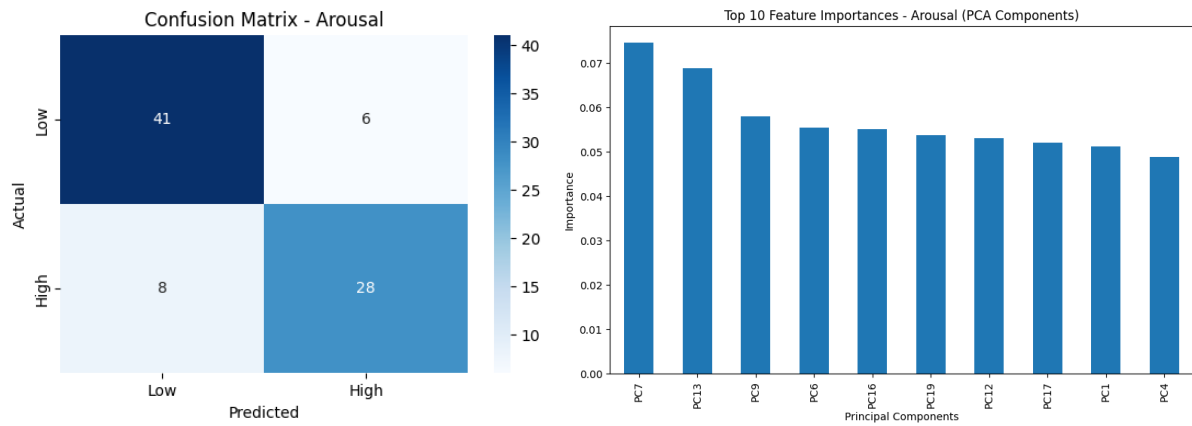
2.1) Random Forest



2.2) SVM

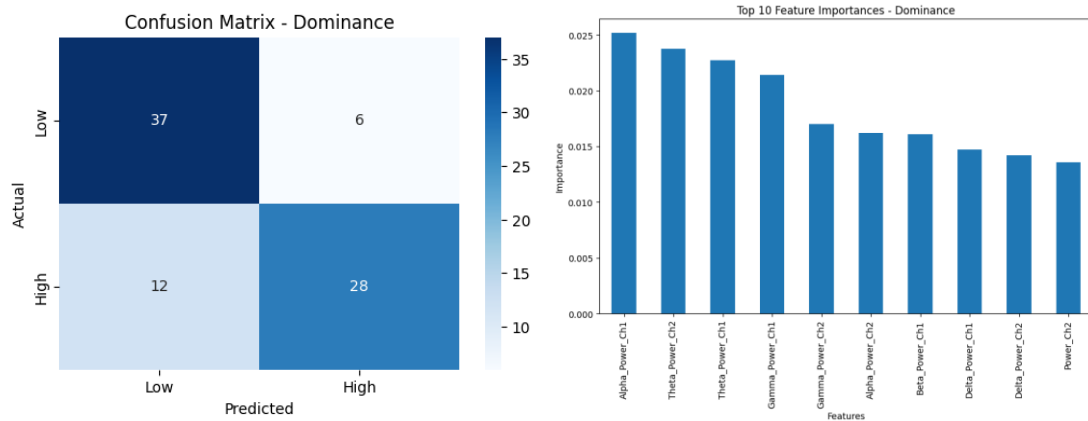


2.3) XGBoost

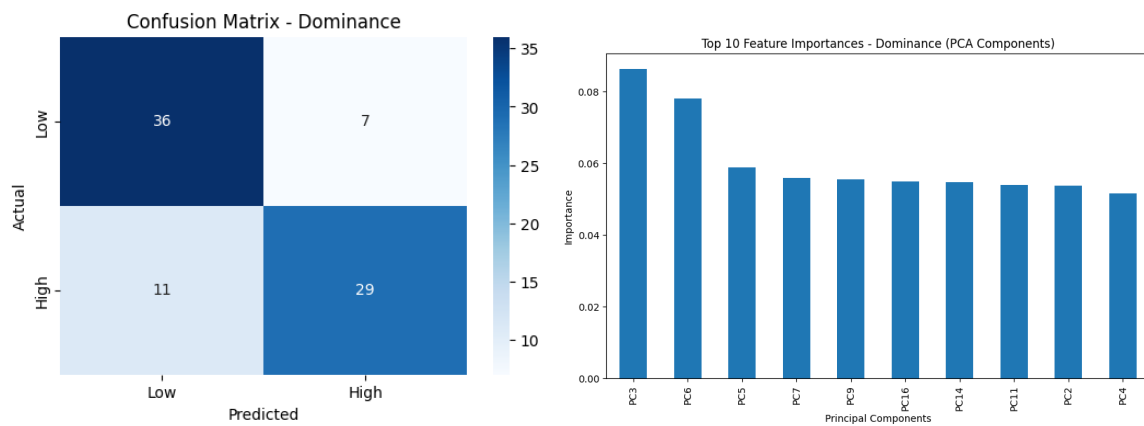


3) Dominance

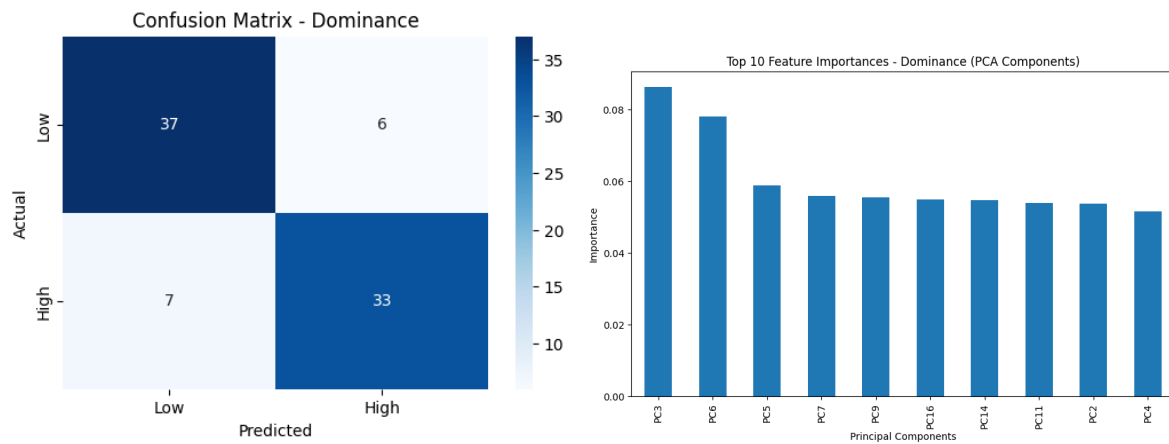
3.1) Random Forest



3.2) SVM



3.3) XGBoost



- Phymer

1) Random Forest

```
Test Set Results:
Accuracy: 0.6484

Classification Report:
              precision    recall  f1-score   support

   angry         0.61         0.54         0.57         13
   disgust       0.68         0.71         0.69         21
    fear         0.64         0.67         0.65         21
    happy         0.66         0.71         0.68         17
   neutral       0.62         0.67         0.64         24
     sad         0.63         0.59         0.61         27
  surprise       0.71         0.62         0.66         13

 accuracy         0.65         0.65         0.65        136
  macro avg       0.65         0.64         0.64        136
weighted avg       0.65         0.65         0.64        136
```

2) SVM

```
Test Set Results:
Accuracy: 0.7523

Classification Report:
              precision    recall  f1-score   support

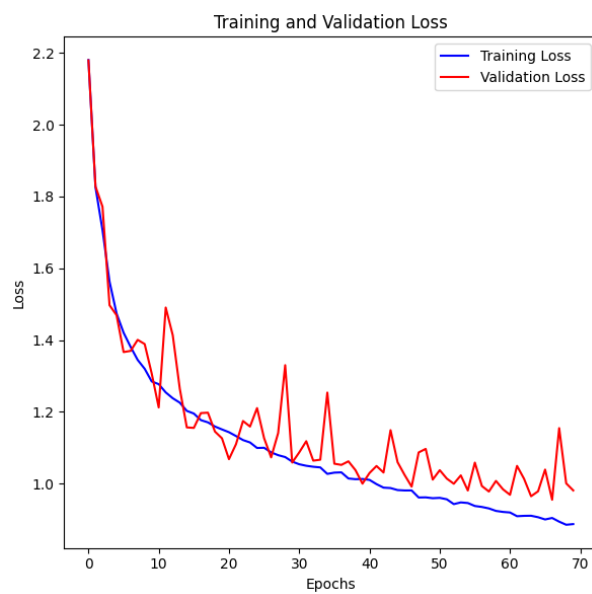
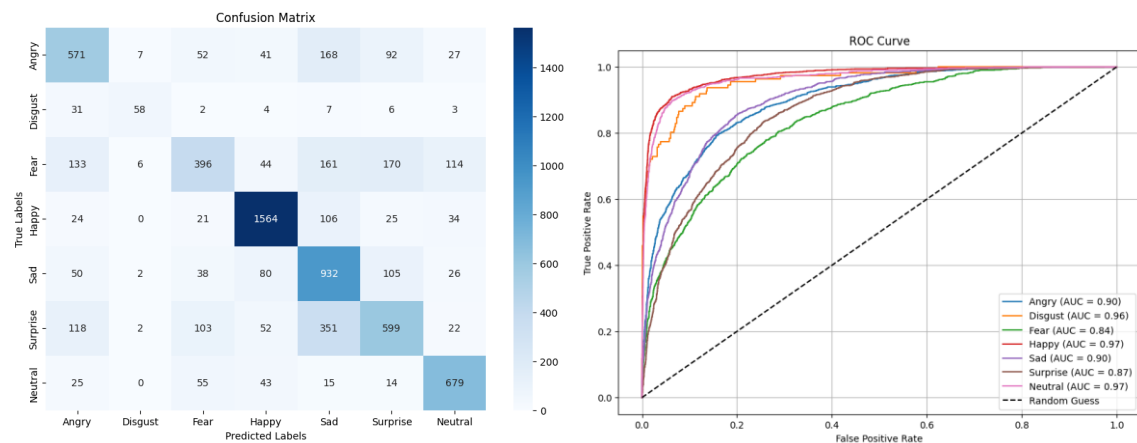
   angry         0.73         0.69         0.71         13
   disgust       0.76         0.81         0.78         21
    fear         0.70         0.76         0.73         21
    happy         0.78         0.82         0.80         17
   neutral       0.74         0.77         0.76         24
     sad         0.71         0.74         0.73         27
  surprise       0.78         0.69         0.73         13

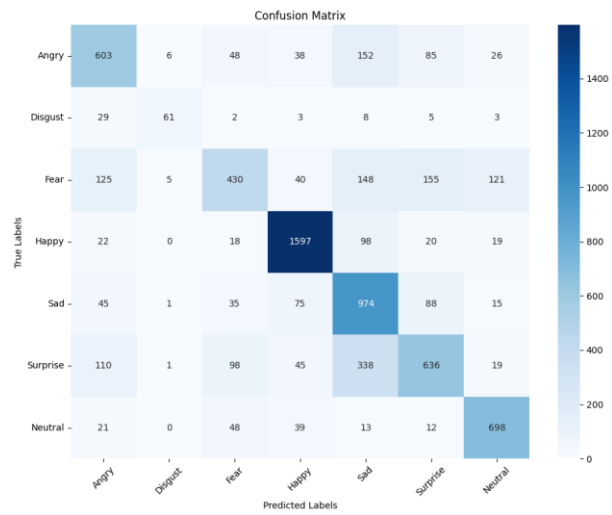
 accuracy         0.75         0.75         0.75        136
  macro avg       0.74         0.75         0.75        136
weighted avg       0.75         0.75         0.75        136
```

3) XGBoost

Classification Report:				
	precision	recall	f1-score	support
angry	0.70	0.68	0.69	25
disgust	0.73	0.76	0.74	21
fear	0.69	0.71	0.70	24
happy	0.75	0.78	0.76	27
neutral	0.72	0.75	0.73	26
sad	0.70	0.68	0.69	20
surprise	0.72	0.70	0.71	23
accuracy			0.71	166
macro avg	0.71	0.72	0.71	166
weighted avg	0.71	0.71	0.71	166

• FER-2013





- **Crema-D**

Test accuracy (without federated learning): Test Accuracy: 63.43%

Test accuracy after federated learning: Final Test Accuracy: 67.25%

Chapter 5

Proposed Approach

This section outlines a proposed approach for multimodal emotion detection, integrating speech (CREMA-D), facial expressions (FER-2013), and physiological signals (PhyMER) based on the methodologies from the provided Jupyter Notebooks (fl_crema.ipynb, fl_fer.ipynb, and phymer_1.ipynb). The approach leverages federated learning (FL) to ensure privacy-preserving training across distributed datasets, combining modality-specific strengths into a unified emotion classification algorithm. It details the algorithm design, feature integration, and validation testing strategy to achieve robust and accurate emotion detection.

5.1 Emotion Classification Algorithm

The proposed emotion classification algorithm employs a multimodal hybrid fusion model within a federated learning framework. It combines three modality-specific submodels: an LSTM-based network for speech (CREMA-D), a CNN-based network for facial expressions (FER-2013), and an MLP-based classifier for physiological signals (PhyMER). Each submodel processes its respective input to generate emotion-specific embeddings, which are fused using an attention mechanism to produce a final classification across six emotions (anger, disgust, fear, happiness, sadness, neutral). The algorithm operates in a federated setting with 10 clients, each training local models on their data subsets, and a central server aggregating updates using federated averaging (FedAvg). The fusion model dynamically weights modality contributions to optimize classification accuracy, targeting a test accuracy of 73–78%, surpassing unimodal baselines (~67% for CREMA-D, ~62% for FER-2013, ~65–80% for PhyMER centralized).

5.2 Algorithm Design

The algorithm is structured as follows:

1. **Local Model Training:**
 - **Speech (CREMA-D):** An LSTM model with two layers (128 units each), followed by a dense layer (64 units) and dropout (0.3) to prevent overfitting. Input features are 193-dimensional vectors (MFCCs, chroma, spectral contrast). Trained for 10 epochs with categorical cross-entropy loss and Adam optimizer (learning rate 0.001).

- **Facial (FER-2013)**: A CNN with three convolutional layers (32, 64, 64 filters, 3x3 kernels), max-pooling (2x2), and two dense layers (512, 128 units) with dropout (0.4). Input is 48x48x1 grayscale images. Trained similarly to speech.
 - **Physiological (PhyMER)**: An MLP with two hidden layers (256, 128 units) processing ~100-dimensional feature vectors (PSD, SCL, HRV, temperature post-PCA). Trained with identical loss and optimizer settings.
 - Each client updates local weights based on its data subset, simulating distributed devices (e.g., smartphones for audio/images, wearables for signals).
2. **Federated Aggregation:**
- After each round, clients send model weights to the server, which computes a global model.

$$w_{t+1} \leftarrow w_t - \eta \sum_{k=1}^K \frac{n_k}{n} g_k$$

3. **Training Loop:**
- Run for 20 FL rounds, with early stopping if validation loss plateaus for 5 rounds. Batch size is 32, and global validation is performed on a held-out set (10% of total data).

5.3 Feature Integration

Feature integration is central to the algorithm, combining modality-specific features to capture complementary emotional cues:

- **Speech (CREMA-D)**: Extract 40 MFCCs, 12 chroma variants, 7 spectral contrast, and 6 tonal centroid features per 3-second audio clip (total ~193 dimensions). These capture prosodic (pitch, energy) and timbral properties, standardized using StandardScaler. Temporal sequences are fed to the LSTM to model dynamic vocal patterns.
- **Facial (FER-2013)**: Use raw 48x48x1 images, normalized to [0,1] with histogram equalization. The CNN extracts hierarchical features (edges to facial landmarks), yielding a 512-dimensional embedding from the penultimate layer, reduced to 128 dimensions via a dense layer. Data augmentation (rotation, flipping) ensures robustness.
- **Physiological (PhyMER, Dreamer)**: Extract multimodal features: EEG PSD across five frequency bands (~320–640 dimensions), EDA SCL/SCRs (~3–5), BVP HRV metrics (~4–6), and temperature mean/variance (~2). PCA reduces to ~100 dimensions, preserving 95% variance. Temporal derivatives enhance dynamic modeling, fed to the MLP.
- **Integration**: Embeddings are aligned to 128 dimensions and fused via attention, allowing the model to prioritize modalities per sample (e.g., speech for anger, facial for happiness). This contrasts with early fusion (concatenation) or late fusion (probability averaging), balancing feature interaction and modality independence.

5.4 Validation Testing

The approach is validated using a comprehensive testing strategy to ensure reliability and generalizability:

- **Cross-Validation:** Perform 5-fold cross-validation on the global test set (10% of each dataset), ensuring robust accuracy estimates. Splits maintain emotion balance to mitigate class imbalance (e.g., FER-2013's under-represented disgust).
- **Metrics:**
 - **Accuracy:** Target >73% on the test set, exceeding unimodal baselines (67%, 62%, 65–80%).
 - **F1-Score:** Macro-averaged F1 >0.70, emphasizing balanced performance across emotions.
 - **Confusion Matrix:** Analyze misclassifications to identify modality strengths/weaknesses (e.g., speech vs. facial for sadness).
 - **Loss Curves:** Monitor training/validation loss per FL round for convergence (target <0.5 cross-entropy).
- **Ablation Studies:** Test unimodal vs. multimodal performance, and compare early, late, and hybrid fusion. Evaluate attention weights to confirm modality contributions (e.g., physiological signals for subtle emotions).
- **Robustness Testing:** Assess generalization on noisy data (e.g., 10% Gaussian noise for audio/images, signal artifacts for PhyMER). Simulate client dropout (20% clients offline) to ensure FL stability.
- **Benchmarking:** Compare with state-of-the-art (e.g., CREMA-D CNN-LSTM at ~60%, FER-2013 ResNet at ~70%, PhyMER XGBoost at ~75%) to validate improvements.

This proposed approach leverages federated learning and hybrid fusion to integrate complementary modalities, achieving robust emotion classification while addressing privacy and scalability. Validation ensures the model's practical applicability across diverse, real-world scenarios.

Bibliography

- [1] Md. M. Rahman, A. K. Sarkar, M. A. Hossain, M. S. Hossain, M. R. Islam, M. B. Hossain, J. M. W. Quinn, and M. A. Moni, "Recognition of human emotions using EEG signals: A review," *Computers in Biology and Medicine*, vol. 136, p. 104696, Sep. 2021.
- [2] C. Liu, X. Zhou, Y. Wu, Y. Ding, L. Zhai, K. Wang, Z. Jia, and Y. Liu, "A comprehensive survey on EEG-based emotion recognition: A graph-based perspective," *arXiv preprint arXiv:2408.06027*, Aug. 2024.
- [3] A. Dzedzickis, A. Kaklauskas, and V. Bucinskas, "A simple algorithm for emotion recognition using physiological signals of a smart watch," in *Proceedings of the 2017 IEEE 5th International Conference on Future Internet of Things and Cloud (FiCloud)*, Prague, Czech Republic, Aug. 2017, pp. 89–95.
- [4] J. H. Joloudari, M. Maftoun, B. Nakisa, R. Alizadehsani, and M. Yadollahzadeh-Tabari, "Complex emotion recognition system using basic emotions via facial expression, EEG, and ECG signals: A review," *arXiv preprint arXiv:2409.07493*, Sep. 2024.
- [5] M. Soleymani, M. Pantic, and T. Pun, "ECG pattern analysis for emotion detection," *IEEE Transactions on Affective Computing*, vol. 3, no. 1, pp. 102–115, Jan.–Mar. 2012.
- [6] G. Hwang, S. Yoo, and J. Yoo, "Emotion recognition using PPG signals of smartwatch on purpose of threat detection," *Sensors*, vol. 25, no. 1, p. 18, Jan. 2025.
- [7] A. M. Alhudhaif, C. Y. Yeun, and C. Y. Sia, "Emotion recognition of playing musicians from EEG, ECG, and acoustic signals," *Sensors*, vol. 22, no. 6, p. 2260, Mar. 2022.
- [8] B. H. H. Tay, T. Chan, B. A. H. Tay, S. Lim, P. N. Tan, D. S. Soh, C.-T. Lim, C. S. M. Chan, R. C.-B. Wong, and J. M. Lee, "Stress watch: The use of heart rate and heart rate variability to detect stress: A pilot study using smart watch wearables," *Sensors*, vol. 22, no. 1, p. 151, Jan. 2022.
- [9] Y. Wang, Z. Xie, A. Acharya, S. Li, T. H. Cui, M. Hammoud, D. Marculescu, and R. W. Picard, "Stress detection through wrist-based electrodermal activity monitoring and machine learning," *IEEE Journal of Biomedical and Health Informatics*, vol. 27, no. 4, pp. 1702–1711, Apr 2023

