



# Transformer-Based Emotion Recognition Using Multimodal Data Fusion

Submitted By:  
**Shivendu Mishra (206124031)**

Guide:  
**Dr. R. Bala Krishnan**

DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING  
NATIONAL INSTITUTE OF TECHNOLOGY, TIRUCHIRAPPALLI



# 1. PROBLEM STATEMENT, OBJECTIVES, & MOTIVATION

## ➤ Problem Statement

- To develop a robust multimodal emotion recognition system using Transformer-based architecture
- Effectively models the cross-modal dependencies and temporal dynamics across the physiological signals

## ➤ Objectives

- To design a Transformer-based multimodal system
- To fuse the ECG, EDA, and acceleration features to predict the emotions
- To reduce the artifacts in the signals
- To improve the accuracy

## ➤ Motivation

- To enhance the human-computer interaction
- To bridge the gap between the laboratory validations and real-world deployment
- To support affective computing



## 2-A. LITERATURE SURVEY

Sl. No.	Paper Details	Techniques Used	Gaps Identified
1.	<b>Title:</b> UAED: Unsupervised Abnormal Emotion Detection Network Based on Wearable Mobile Device <b>Author:</b> J. Zhu et al. <b>Year:</b> 2023	<ul style="list-style-type: none"><li>• Gaussian Mixture VAE</li><li>• 2D-CNN with stacking</li><li>• Whitening distance anomaly scoring</li><li>• Unsupervised training</li></ul>	<ul style="list-style-type: none"><li>• Unimodal approach – Uses only ECG signals</li><li>• No cross-modal attention – Cannot leverage multiple physiological signals</li><li>• Limited emotion granularity – Only binary anomaly detection</li><li>• Fixed window processing – No adaptive temporal modeling</li></ul>
2.	<b>Title:</b> Self-supervised ECG Representation Learning for Emotion Recognition <b>Author:</b> P. Sarkar et al. <b>Year:</b> 2022	<ul style="list-style-type: none"><li>• Self-supervised learning</li><li>• Contrastive pre-training</li><li>• CNN-LSTM architecture</li><li>• Transfer learning fine-tuning</li></ul>	<ul style="list-style-type: none"><li>• Single modality dependency – ECG-only approach</li><li>• No multimodal fusion – Misses complementary signal information</li><li>• Supervised requirement – Needs labeled data for fine-tuning</li><li>• Short-term focus – Limited long-range temporal context</li></ul>



## 2-B. LITERATURE SURVEY CONTD.

Sl. No.	Paper Details	Techniques Used	Gaps Identified
3.	<b>Title:</b> Behavioral and Physiological Signals-Based Deep Multimodal Approach for Mobile Emotion Recognition <b>Author:</b> K. Yang et al. <b>Year:</b> 2023	<ul style="list-style-type: none"><li>• LSTM with attention</li><li>• Multimodal fusion</li><li>• Mobile deployment optimization</li><li>• Behavioral + physiological signals</li></ul>	<ul style="list-style-type: none"><li>• Limited temporal context – LSTM struggles with long sequences</li><li>• Sequential processing – Inherent RNN limitations</li><li>• Shallow fusion – Basic attention on temporal dimension only</li><li>• Limited physiological coverage – Focuses more on behavioral sensors</li></ul>
4.	<b>Title:</b> EmotionSense: An Adaptive Emotion Recognition System <b>Author:</b> Z. Wang et al. <b>Year:</b> 2020	<ul style="list-style-type: none"><li>• Random Forest classifiers</li><li>• Wearable sensor integration</li><li>• Context-aware recognition</li><li>• Personalization approach</li></ul>	<ul style="list-style-type: none"><li>• Shallow models – Limited feature learning capacity</li><li>• Traditional ML limitations – Poor scalability with complex data</li><li>• Feature engineering dependency – Manual feature extraction required</li><li>• No deep temporal modeling – Cannot capture the complex temporal patterns</li></ul>

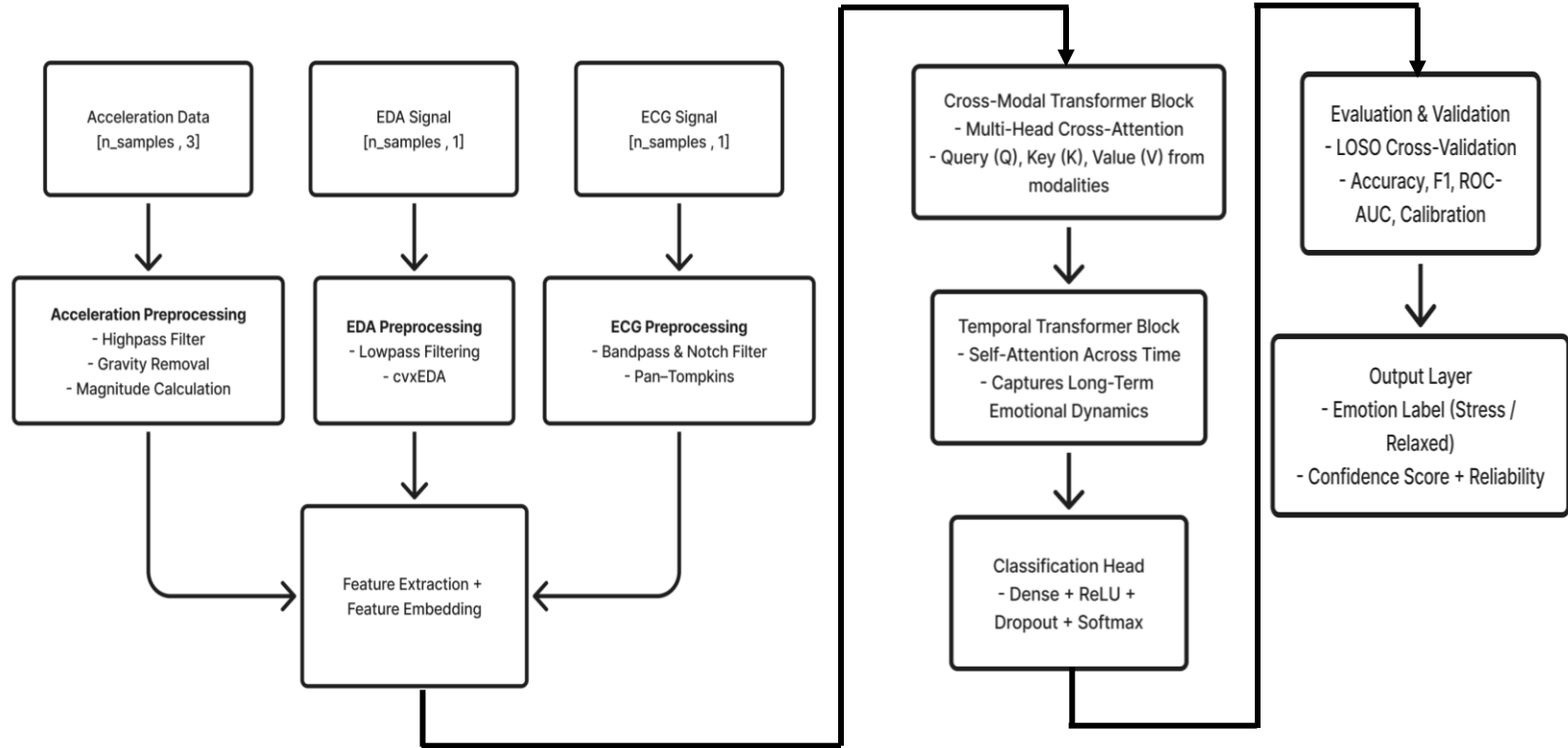


## 2-C. LITERATURE SURVEY CONTD.

Sl. No.	Paper Details	Techniques Used	Gaps Identified
5.	<b>Title:</b> Attention Is All You Need <b>Author:</b> A. Vaswani et al. <b>Year:</b> 2017	<ul style="list-style-type: none"><li>• Transformer Architecture</li><li>• Scaled Dot-Product Attention</li><li>• Multi-Head Attention</li><li>• Positional Encoding</li><li>• Encoder-Decoder Structure</li></ul>	<ul style="list-style-type: none"><li>• High Computational Complexity (<math>O(n^2)</math> with sequence length)</li><li>• No Native Multimodal Support – Designed for NLP</li><li>• Large Data Requirements for effective training</li><li>• No Domain Adaptation for physiological signals</li></ul>



## 3-A. BLOCK SCHEMATIC





## 4-A. ALGORITHMS

### **ALGORITHM:** Pan-Tompkins + Heart Rate Variability (HRV) Feature Extraction

#### ➤ **INPUT**

Raw ECG signal  $x(t)$  sampled at 700 Hz

#### ➤ **PROCESS**

**Step 1:** Apply a bandpass filter to remove the baseline wander and high-frequency noise

**Step 2:** Use the Pan-Tompkins algorithm to detect R-peaks

- Differentiate → Square → Moving window integration → Threshold
- Output is a list of detected R-peak indices

**Step 3:** Compute R–R intervals to estimate beat-to-beat timing

**Step 4:** Derive Heart Rate Variability (HRV) features, such as mean RR, SDNN, RMSSD, etc.

#### ➤ **OUTPUT**

Clean ECG signal with noise removed



## 4-B. ALGORITHMS CONT.

### **ALGORITHM:** Convex Optimization (cvxEDA Decomposition)

#### ➤ **INPUT**

Raw EDA signal  $y(t)$  collected from wrist sensor

#### ➤ **PROCESS**

**Step 1:** Apply a low-pass filter (cutoff = 2 Hz) to remove high-frequency fluctuations

**Step 2:** Use cvxEDA (Convex Optimization-based Decomposition) to split the signal

- Tonic component (T): Slow baseline skin conductance
- Phasic component (P): Fast changes linked to emotional arousal

**Step 3:** Detect Skin Conductance Response (SCR) peaks from the phasic component  $r$

**Step 4:** Extract statistical features — Mean, Standard Deviation, and Number of SCR Peaks

#### ➤ **OUTPUT**

Clean ECG signal with noise removed





## 4-C. ALGORITHMS CONT.

**ALGORITHM:** Convex Optimization (cvxEDA Decomposition)

➤ **INPUT**

3-axis acceleration signal  $a_x(t)$  ,  $a_y(t)$  ,  $a_z(t)$

➤ **PROCESS**

**Step 1:** Apply a high-pass filter (cutoff = 0.5 Hz) to remove the gravitational components

**Step 2:** Compute the magnitude of the resultant acceleration

**Step 3:** Segment the signal into small time windows (for example, 4 to 6 s)

**Step 4:** Extract the activity-level features such as mean, variance, energy, and entropy

➤ **OUTPUT**

Clean ECG signal with noise removed



## 4-D. ALGORITHMS CONTD..

### ALGORITHM: Cross-Modal Transformer with Multi-Head Attention (Fusion Stage)

#### ➤ INPUT

Embedded feature sequences for each modality:  $e_{ECG}$ ,  $e_{EDA}$ ,  $e_{ACC}$  obtained from CNN encoders or preprocessing

#### ➤ PROCESS

**Step 1:** Prepare attention inputs

- For each modality  $m$ , create: Query ( $Q_m$ ) , Key ( $K_m$ ) , Value ( $V_m$ )

**Step 2:** Compute attention weights

- Use scaled dot-product attention to decide how strongly one signal should focus on another

**Step 3:** Use multiple attention heads

- Several parallel “heads” learn different kinds of relationships

**Step 4:** All heads are concatenated and linearly combined

#### ➤ OUTPUT

Model outputs a single fused vector  $F_{fusion}$



## 4-E. ALGORITHMS CONTD..

### ALGORITHM: Temporal Transformer Block (Sequential Modeling)

#### ➤ INPUT

Sequence of fused embeddings  $F_{fusion,t}$  obtained from the Cross-Modal Transformer

#### ➤ PROCESS

**Step 1:** Add positional information

- Since Transformers don't know time order naturally, add positional encodings

**Step 2:** Compute self-attention across time

- For every time window, compute attention

**Step 3:** Residual + Layer Normalization

- Maintains stable gradients and keeps temporal context consistent

**Step 4:** A small network refines each time-step's contextual embedding, strengthening the temporal coherence

#### ➤ OUTPUT

Final time-dependent emotional embedding  $F_{temporal}$



## 5. EXPERIMENTS & PERFORMANCE METRICS

### ➤ Experimental Setup

- Dataset: WESAD (primary) and AffectiveROAD (supplementary)
- Training:
  - ECG/EDA CNN: 100 epochs  $\times$  32 batch size
  - Acceleration CNN: 50 epochs  $\times$  64 batch size
  - Optimizer: AdamW (Weight Decay = 0.01,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e-8$ )
  - Validation: LOSO Cross Validation

### ➤ Performance Metrics

- Core Metrics: Accuracy, ROC-AUC, PR-AUC
- Calibration Metric: Expected Calibration Error (ECE) – Assesses the reliability of predictions
- Subject-Level Metric: Inter Subject Variability (ISV) – Measures the generalization

### ➤ Model Setup

- Modality specific CNN encoders for ECG, EDA, and Acceleration
- Cross-Modal Transformer Fusion + Temporal Transformer for sequence modeling



## 6. IMPLEMENTATION ENVIRONMENT

### ➤ Hardware

- **Processor:** Intel i7 / AMD Ryzen 7 (8-core, 3.0GHz+)
- **RAM:** 16 GB minimum
- **Storage:** 500 GB NVMe SSD
- **GPU:** NVIDIA RTX 4060 (12 GB VRAM) or equivalent CUDA-enabled GPU
- **Operating System:** Windows 11 (64-bit)

### ➤ Software & Libraries

- **Language:** Python 3.10.x
- **IDE:** Visual Studio Code
- **Core Framework:** PyTorch (with CUDA support)
- **Deep Learning Models**
  - NeuroKit2/SciPy Essential for robust ECG R-peak detection, cvxEDA decomposition, and HRV analysis
  - Transformers (for implementing attention mechanisms)
- **Data Handling & Processing:** NumPy, Pandas, OpenCV, Open3D
- **Visualization Tools:** Matplotlib, TensorBoard, Seaborn



## 7. EXPECTED MILESTONE / WORK PLAN

Milestone	Period	Works to be Completed
MS0	14 Aug 2025 – 26 Oct 2025	Literature survey, tool & language selection, dataset setup (WESAD/AffectiveROAD), sensor level study, bio-signal study, algorithm selection, and project planning
MS1	29 Oct 2025 – 23 Nov 2025	Implement preprocessing for ECG, EDA, and acceleration (filtering, HRV, SCR, activity features) and begin 1D-CNN feature extraction
MS2	26 Nov 2025 – 21 Dec 2025	Integrate feature embeddings and Transformer-based fusion, temporal modeling, and perform preliminary training and validation
MS3	24 Dec 2025 – 31 Dec 2025	Conduct full LOSO evaluation, performance analysis and finalize document



## 8. REFERENCES

- [1] J. Zhu, J. Jiang, and X. Zhao, “UAED: Unsupervised Abnormal Emotion Detection Network Based on Wearable Mobile Device,” *IEEE Transactions on Affective Computing*, vol. PP, no. 99, pp. 1–12, 2023.
- [2] K. Yang, X. Huang, and Y. Liu, “Behavioral and Physiological Signals-Based Deep Multimodal Approach for Mobile Emotion Recognition,” *IEEE Transactions on Mobile Computing*, vol. PP, no. 99, pp. 1–10, 2023.
- [3] P. Sarkar, S. Dey, and A. Dutta, “Self-Supervised ECG Representation Learning for Emotion Recognition,” *IEEE Sensors Journal*, vol. 22, no. 18, pp. 17456–17464, 2022.
- [4] Z. Wang, X. Chen, and Y. Zhang, “EmotionSense: An Adaptive Emotion Recognition System,” *IEEE Internet of Things Journal*, vol. 7, no. 9, pp. 8550–8563, 2020.
- [5] D. Pollreisz and N. TaheriNejad, “Detection and Removal of Motion Artifacts in PPG Signals,” *Mobile Networks and Applications*, vol. 27, pp. 728–738, Aug. 2019.
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention Is All You Need,” *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 5998–6008, 2017.
- [7] A. Greco, G. Valenza, M. Nardelli, M. Bianchi, L. Citi, and E. P. Scilingo, “Force–Velocity Assessment of Caress-Like Stimuli Through the Electrodermal Activity Processing: Advantages of a Convex Optimization Approach,” *IEEE Transactions on Human-Machine Systems*, vol. 47, no. 1, pp. 91–100, Feb. 2017.



**THANK YOU**