

INTRODUCTION

- Emotion recognition is vital in affective computing and HCI, and eye movement has emerged as a promising modality.
- Prior work shows that blink patterns and pupil changes correlate with emotional states.
- Using eye data alone can offer a non-invasive, cost-effective alternative to EEG.
- This work proposes a deep model (DGCCA-AM) to classify emotions from eye movement alone.

OBJECTIVES

- Develop an interpretable deep model (DGCCA-AM) using only eye-tracking data for four-class emotion classification. We group eye features into semantically meaningful sets and apply DGCCA with an attention mechanism to align and fuse them for robust emotion prediction.

PROBLEM STATEMENT

- Most existing emotion classifiers rely on EEG or multimodal inputs, while purely eye-based methods remain under-explored.
- Standalone eye-tracking systems often lack sophisticated feature fusion.
- There is a need for models that explicitly group and fuse eye movement features (e.g. pupil, saccades, events) to fully exploit their emotional cues.

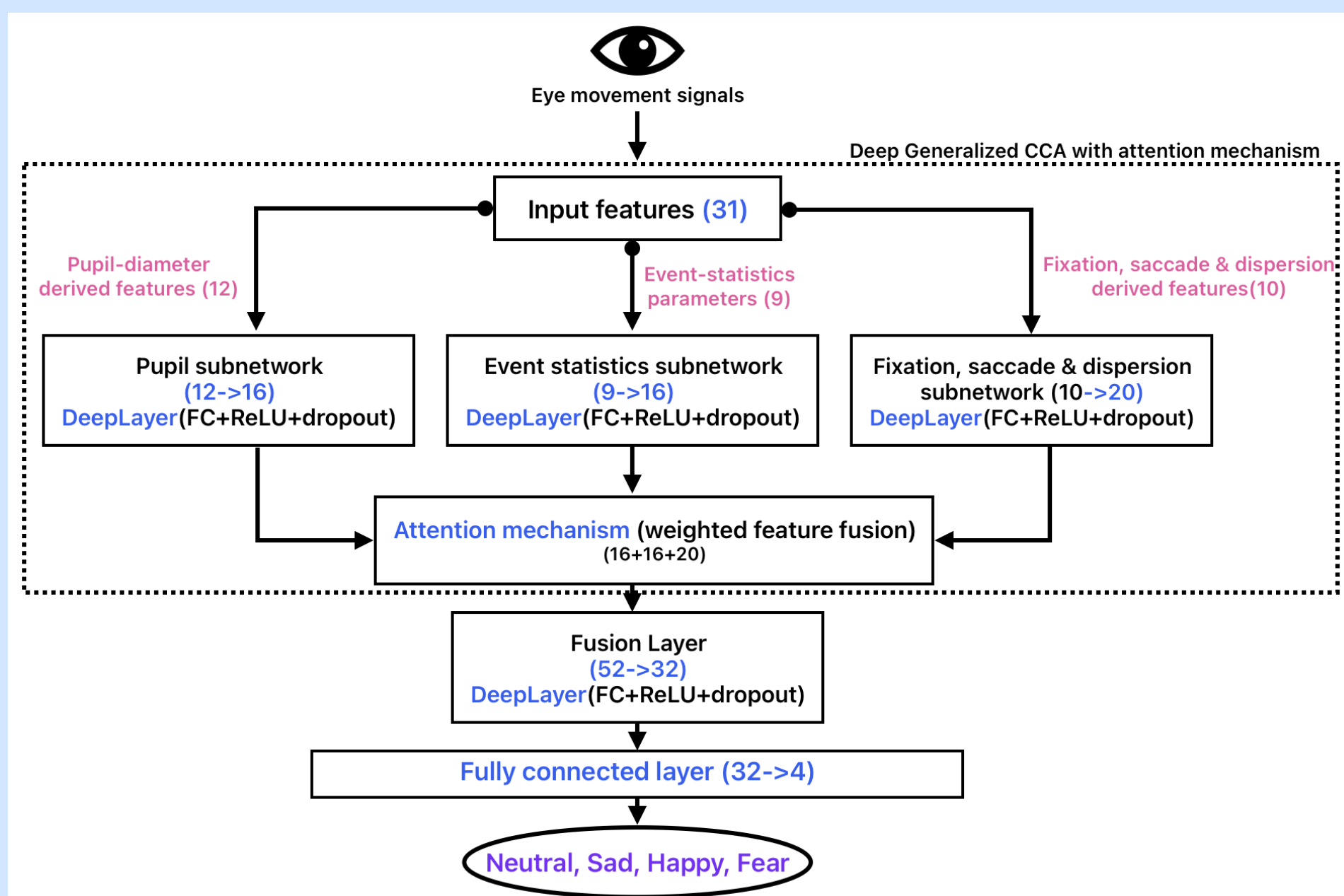
DATASETS

Class distribution in training, validation, and test sets for Session 3 of the SEED-IV dataset.

Emotion	Training	Validation	Test
Neutral	2473	275	687
Sad	2462	274	677
Fear	1760	196	479
Happy	2182	242	602

PROPOSED METHOD

To leverage semantic structure in eye movement data, we divided features into three semantic groups (31 features total):



- Pupil-Based Features (12):** Capture emotional arousal via pupil diameter stats (mean, std. for both eyes) and Differential Entropy (DE) across 4 frequency bands: 0–0.2, 0.2–0.4, 0.4–0.6, and 0.6–1.0 Hz.
- Fixation, Saccade Dispersion Features (10):** Model gaze control and attention using saccade amplitude/duration (mean, std.), fixation duration, and gaze dispersion along X and Y axes.
- Event-Based Features (9):** Quantify blink rate, fixation/saccade counts and durations, average saccade amplitude, and latency to capture temporal gaze event patterns.

Deep Generalized Canonical Correlation Analysis with Attention Mechanism (DGCCA-AM) model structure.

- Group-Specific Subnetworks:** Each feature group (Pupil, Fixation/Saccade, Events) is passed through a separate subnetwork with a fully connected (FC) layer and ReLU activation to learn group-specific embeddings. Each group X_i is passed through a dedicated subnetwork:

$$X_i \rightarrow FC_i \rightarrow \text{ReLU} \rightarrow h_i$$

producing latent vectors h_1, h_2, h_3 .

- DGCCA Alignment Layer:** Learns shared latent representations by maximizing correlation between the feature groups using Deep Generalized Canonical Correlation Analysis (DGCCA).

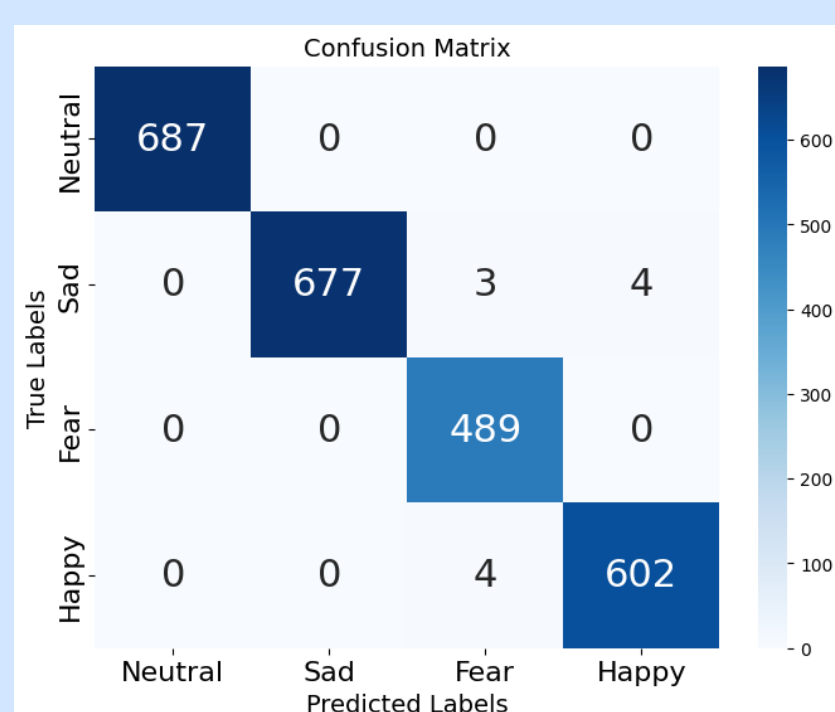
$$\max_{W_1, W_2, \dots, W_n} \sum_{i < j} \text{corr}(W_i X_i, W_j X_j)$$

- Attention Fusion Mechanism:** An attention layer computes weights for each aligned feature group, dynamically emphasizing the most informative signals per sample.

$$\mathbf{z} = \sum_{i=1}^3 \alpha_i h_i, \quad \text{where} \quad \alpha_i = \frac{\exp(e_i)}{\sum_j \exp(e_j)}$$

- Incorporating DGCCA (Deep Generalized CCA) encourages each subnetwork (modality) to learn representations that are maximally correlated across the different modalities.
- Weighted features from each modality are concatenated and passed through a fully-connected fusion layer that reduces the feature dimension from 52 to 32. The resulting fused vector is then fed into a Softmax classifier (32→4) to predict one of four emotion classes (Happy, Sad, Neutral, Fear).
- The model is trained end-to-end with the Adam optimizer and a categorical cross-entropy loss function.
- Performance is evaluated under both subject-dependent and subject-independent settings, using accuracy scores and confusion matrices to measure classification effectiveness.

RESULTS



- Session 3 (intra-subject) confusion matrix shows minimal misclassifications.
- Intra-subject accuracy peaked at 99.92
- Feature ablation: baseline (all features) accuracy 99.92%; without pupil features 93.71% (−6.2%); without event-statistics 97.77% (−2.1%); without fixation/saccade 97.93% (−2.0%).


CONCLUSION

- We demonstrate that eye movement features alone can achieve near-perfect emotion classification in a subject-dependent setting (99.92% accuracy)
- The DGCCA-AM model effectively fuses semantically grouped features to capture emotion-relevant signals.
- A key finding is the dominant role of pupil-based features and event statistics in recognition performance.
- Eye data alone can rival multimodal methods when models are tailored to them. However, cross-subject generalization remains limited (63% accuracy), underscoring the need for subject-specific or adaptive techniques.

LIMITATION AND FUTURE WORK

- Limited cross-subject generalization due to individual variability in eye movement patterns.
- Domain adaptation methods needed for robust, subject-invariant emotion recognition.
- Future work: evaluate in diverse, real-world settings for practical deployment.

REFERENCES

-  Lan, Y.-T., W. Liu, and B.-L. Lu (July 2020). "Multimodal emotion recognition using deep generalized canonical correlation analysis with an attention mechanism". In: *Proceedings of the 2020 International Joint Conference on Neural Networks (IJCNN)*. Glasgow, UK, pp. 1–8. DOI: 10.1109/IJCNN48605.2020.9207625. URL: <https://doi.org/10.1109/IJCNN48605>.