# GLaRE: A Graph-based Landmark Region Embedding Network for Emotion Recognition

**Debasis Maji**
Department of Computer & System Sciences
Visva-Bharati
India, 731235
youdebasis@gmail.com

**Debaditya Barman**
Department of Computer & System Sciences
Visva-Bharati
India, 731235
debadityabarman@gmail.com

August 29, 2025

## Abstract

Facial expression recognition (FER) is a crucial task in computer vision with wide range of applications including human–computer interaction, surveillance, and assistive technologies. However, challenges such as occlusion, expression variability, and lack of interpretability hinder the performance of traditional FER systems. Graph Neural Networks (GNNs) offer a powerful alternative by modeling relational dependencies between facial landmarks, enabling structured and interpretable learning. In this paper, we propose GLaRE, a novel Graph-based Landmark Region Embedding network for emotion recognition. Facial landmarks are extracted using 3D facial alignment, and a quotient graph is constructed via hierarchical coarsening to preserve spatial structure while reducing complexity. Our method achieves $\sim$64.89% accuracy on AffectNet[1] and $\sim$94.24% on FERG[2], outperforming several existing baselines. Additionally, ablation studies have demonstrated that region-level embeddings from quotient graphs have contributed to improved prediction performance.

***Keywords*** Facial Expression Recognition · Graph Neural Network · Deep learning · Human Computer Interaction.

## 1 Introduction

Emotion recognition through facial expressions is one of the essential components of artificial intelligence (AI), particularly in the area of computer vision and human-computer interaction (HCI). Among the various modalities for emotion recognition, such as audio, video, and bio-signals, facial gestures offer a direct and intuitive cue to human emotions. Facial landmarks, which abstract away identifiable facial details, offer an efficient and privacy-preserving alternative to the processing of the entire image of the face by capturing only key facial positions. Their subtle movements are strongly associated with different affective states, making them well-suited for real-time applications in smart devices [25, 26], driver monitoring systems [9, 14], e-learning [7], and immersive gaming environments [22]. This has resulted in development of Facial Landmark based Emotion Recognition, which holds promise for both accuracy and resource efficiency.

Although facial landmarks [4] provide semantically rich and privacy preserving indicators for Facial Expression Recognition (FER), they have remained underutilized due to the lack of effective methods to model their spatial and temporal dependencies. While early studies have demonstrated the effectiveness of combining landmarks with appearance features or integrating them in multimodal pipelines (i.e., with EEG or video data) [17, 12, 29] , standalone landmark based methods require advanced geometric reasoning to reach full potential. The advancement of geometric deep learning and graph neural networks (GNNs) provides further opportunities to effectively utilize facial landmarks for FER.

---

[1]https://www.mohammadmahoor.com/pages/databases/affectnet/
[2]https://grail.cs.washington.edu/projects/deepexpr/ferg-2d-db.html
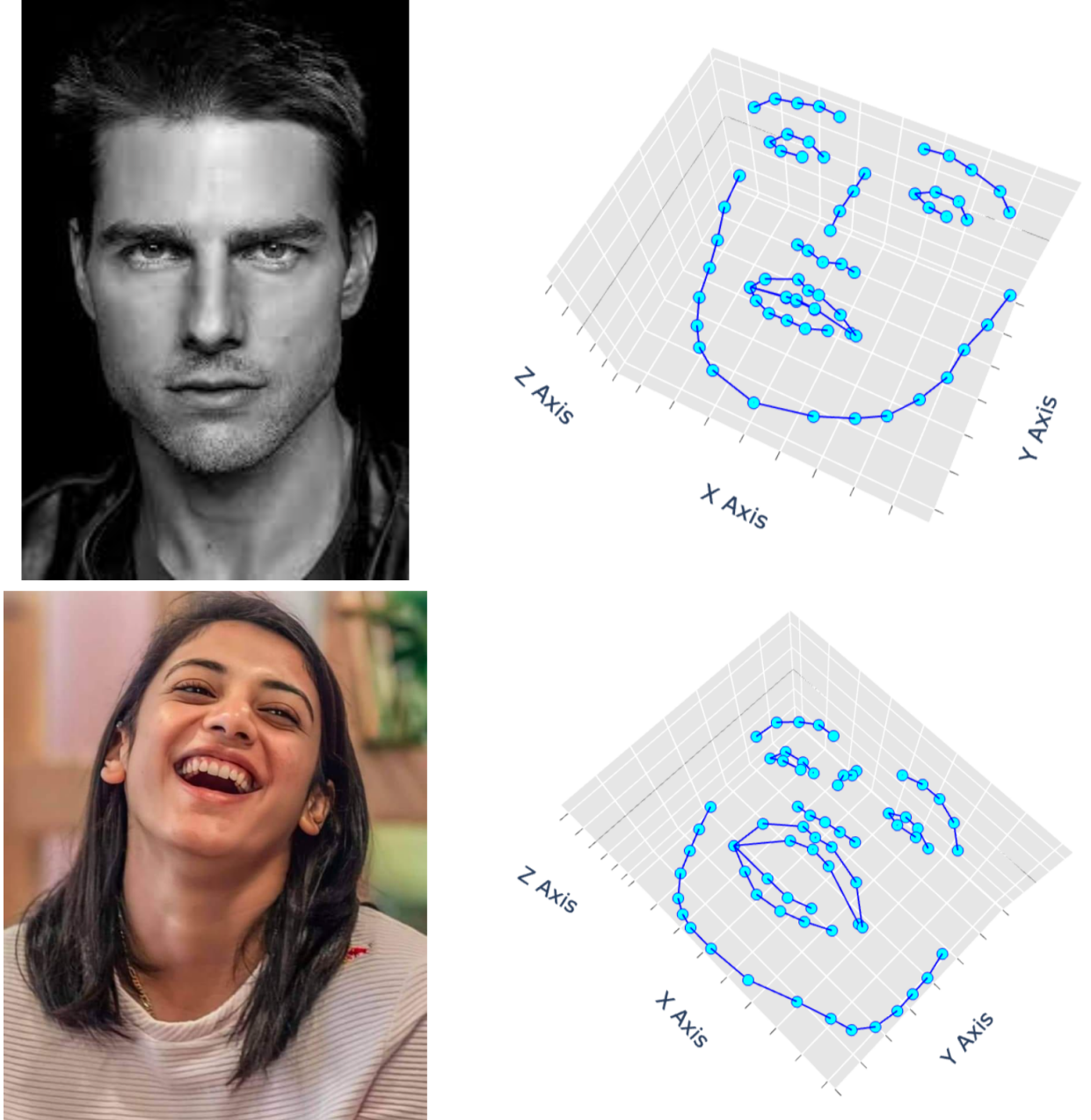
Figure 1: Visual comparison of facial expression images and their corresponding 3D facial landmarks.

Convolutional neural networks (CNNs) [20, 23, 28] have achieved notable success in FER by learning local patterns from facial images. However, they often fall short in capturing the geometric relationships among facial landmarks. Most CNN-based methods focus on extracting deep appearance features, overlooking the underlying spatial structure of the face, which limits generalization and interpretability.

To overcome these limitations, GNNs [13, 21, 1] have emerged as a powerful alternative. Unlike CNNs, GNNs naturally operate on non-Euclidean data such as facial landmarks, modeling both the local and global structural relationships among keypoints. In the context of FER, each facial landmark is modeled as a node in a facial landmark graph. Edges in the graph encode anatomical or spatial connections, enabling the model to capture physiologically meaningful interactions. This facial landmark graph formulation enhances representation and generalization while improving robustness and interpretability by aligning with human cognitive mechanisms. Fig. 1 illustrates two example of faces

2

and their corresponding facial landmark graphs, highlighting how expression information can be structurally encoded. The key contributions of this work are as follows:

- We have proposed a novel graph-based model, GLaRE (Graph-based Landmark Region Embedding Network), which effectively captures both local and global geometric relationships among facial landmarks for emotion recognition.
- We have introduced a hierarchical quotient graph coarsening mechanism that has significantly reduced computational complexity while preserving essential facial structural information, thereby rendering the proposed model lightweight and well-suited for real-time facial emotion recognition tasks.

## 2  Problem Statement

A facial image can be represented as a landmark graph $G = (V, E)$, where $V = \{v_1, v_2, \cdots, v_n\}$ and $E \subseteq V \times V$ is the set of undirected edges connecting anatomically adjacent landmarks, and each node $v_i \in V$ is associated with a spatial coordinate feature vector $x_i \in \mathbb{R}^d$, derived using a landmark detector $\Phi$. Thus, the graph $G = (V, E, X)$ captures the facial structure, where $X \in \mathbb{R}^d$ is the matrix of node features. To incorporate anatomical hierarchy, the facial graph is partitioned into $R$ predefined facial regions.

$$V = V^{(1)} \cup V^{(2)} \cup \cdots \cup V^{(r)}, \text{ where } V^{(r)} \subset V \tag{1}$$

Each facial graph $G_i = (V_i, E_i, X_i)$ in the dataset is associated with a ground-truth emotion label $y_i \in \mathcal{Y}$ where $\mathcal{Y} = \{1, 2, \cdots, C\}$ denotes the set of $C$ discrete expression classes. The objective is to learn the following mapping that can accurately predict the corresponding expression label from the structural and spatial properties of the facial graph.

$$f : G_i \rightarrow y_i \tag{2}$$

This defines a graph classification problem, where each input graph encodes a face and its configuration, and the output is a categorical expression.

## 3  Proposed Method

### 3.1  Graph Construction from Facial Images

To prepare the input image for graph-based learning, a three-stage preprocessing pipeline has been adopted.

First, 3D facial landmarks have been extracted from images using the Facial Alignment Network (FAN) [3]. Each image has been processed to identify a fixed set of landmark points, thereby capturing the geometric structure of key facial regions such as the eyes, nose, mouth, and jawline. These landmark coordinates $\{\mathbf{p}_i \in \mathbb{R}^3\}_{i=1}^N$ have served as the geometric basis for node initialization.

Second, appearance features have been extracted using a pretrained ResNet-18 model [11], where the global average pooling and fully connected layers have been excluded to retain spatial feature maps for bilinear interpolation at landmark positions. The sampled descriptors have then been reduced to $f$ dimensions using principal component analysis (PCA). These appearance features have been concatenated with the normalized 3D landmark coordinates, yielding node features.

$$\mathbf{x}_i = \left[ \hat{\mathbf{p}}_i \parallel \mathbf{a}_i \right] \in \mathbb{R}^{f+3}, \quad i = 1, \ldots, N \tag{3}$$

where $\hat{\mathbf{p}}_i \in \mathbb{R}^3$ denotes normalized coordinates and $\mathbf{a}_i \in \mathbb{R}^f$ denotes the reduced appearance descriptor.

Finally, a fine-level graph $\mathcal{G}_f = (V_f, E_f)$ has been constructed, where each node $v_i \in V_f$ corresponds to a landmark with feature vector $\mathbf{x}_i$. Edges have been established using a $k$-nearest neighbor (kNN) criterion in Euclidean space:

$$E_f = \big\{ (v_i, v_j) \mid v_j \in \text{kNN}(v_i; \mathbf{p}_i) \big\}. \tag{4}$$

This procedure has ensured that each landmark is connected to its local neighborhood, thereby encoding fine-grained spatial relationships among facial components. Together, the node features and adjacency structure have provided a rich representation of facial topology.

As illustrated in Fig. 2, the proposed GLaRE model has followed a two-level graph-based architecture that incorporates equivariant message passing. The architecture has been organized into three key stages: (a) fine-grained node embeddings have been computed from landmark graphs using equivariant GNN layers, (b) a quotient graph has been constructed to compress landmark-level information into region-level nodes, and (c) region embeddings have been further processed through equivariant GNN layers to produce the final global embedding. The subsequent subsections present each stage of the architecture in detail.
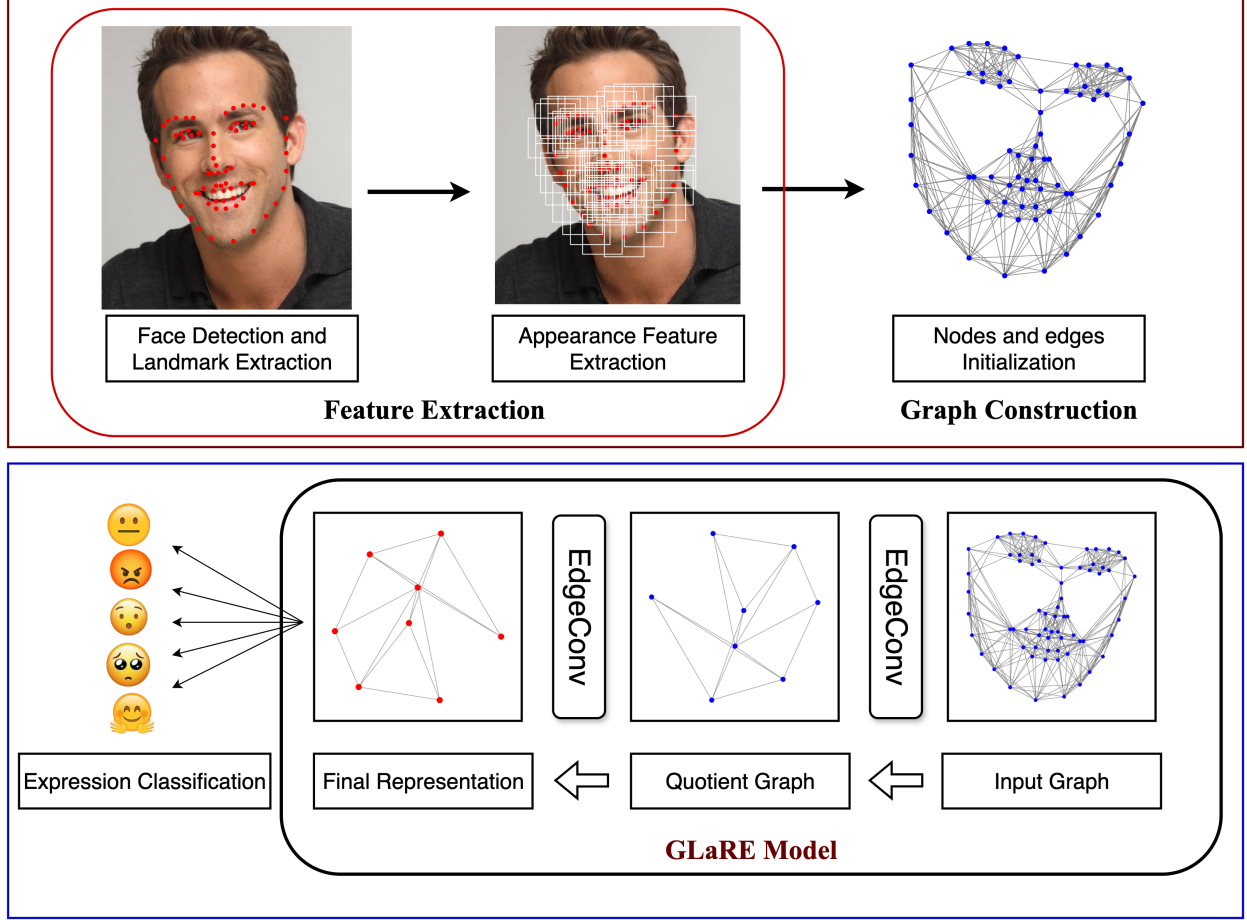
Figure 2: Schematic diagram of GLaRE model for facial expression recognition.

## 3.2 Fine-Level Node Embedding using EdgeConv

Given the fine-grained landmark graph, where each node corresponds to a facial landmark and edges encode spatial proximity, node embeddings are computed using an EdgeConv-based message passing scheme [27]. In this scheme, the embedding of each node $i$ is updated by aggregating messages from its neighbors $j \in \mathcal{N}(i)$ using max aggregation, as defined in (5):

$$\mathbf{h}_i^{(l+1)} = \max_{j \in \mathcal{N}(i)} \phi\left(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} - \mathbf{h}_i^{(l)}\right) \tag{5}$$

Here, $\mathbf{h}_i^{(l)}$ and $\mathbf{h}_j^{(l)}$ represent the feature vectors of the target and neighboring nodes at layer $l$, respectively. The function $\phi$, implemented as a two-layer multi-layer perceptron (MLP) with ReLU activation, transforms the concatenated source and relative features ($[\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)} - \mathbf{h}_i^{(l)}]$) into informative edge-level messages. This formulation captures local differences in landmark configurations and appearance. Moreover, it enhances the representation of facial expression patterns for subsequent processing.

## 3.3 Region-Level Embedding via Quotient Graph

To capture higher-level facial structures critical for FER, a quotient graph is introduced to model regional patterns by aggregating landmark nodes into coarse-grained regions. This approach reduces computational complexity, enhances robustness to local landmark variations, and enables hierarchical processing of local and global facial features, improving the model's ability to distinguish expressions.
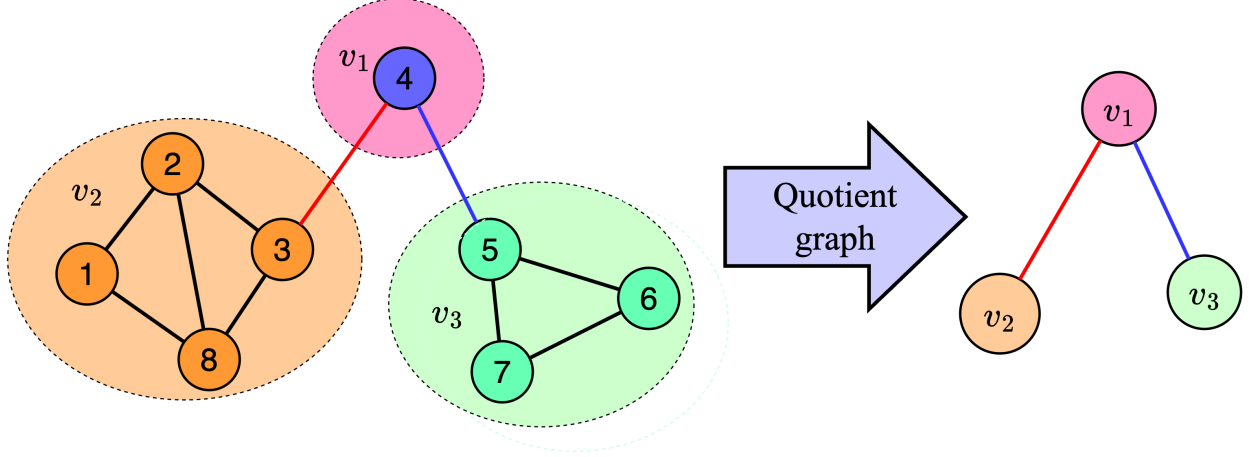
4

Figure 3: The right hand graph is the quotient graph under the equivalence relation specified by the partition set $\mathcal{G}_f/P = \{v_1, v_2, v_3\}$, where $V(G_f) = \{1, 2, ..., 8\}$, $v_1 = \{4\}$, $v_2 = \{1, 2, 3, 8\}$, and $v_3 = \{5, 6, 7\}$.

Let $\mathcal{G}_f = (V_f, E_f)$ be a graph with node set $V_f$ and edge set $E_f$, and let $P = \{V_1, V_2, \ldots, V_k\}$ be a partition set of $V_f$. The quotient graph [10] $\mathcal{G}_q = \mathcal{G}_f/P = (V_q, E_q)$ has node set $V_q = \{V_1, V_2, \ldots, V_k\}$, where each node $V_i$ represents a cluster of nodes from $V_f$. Each node $v \in V_f$ belongs to exactly one subset $V_i$ in the partition, where $i \in \{1, 2, \ldots, k\}$. This one-to-one assignment naturally induces an *equivalence relation* $\sim$ on $V_f$, defined as following.

$$v \sim w \iff \exists V_i \in P \text{ such that } v \in V_i \text{ and } w \in V_i \tag{6}$$

Fig. 3 presents an example of how a quotient graph can be formed.

Let the landmark-level graph be defined as $G = (V, E)$, where $V$ is the set of landmark nodes and $E$ is the set of landmark edges. The landmark nodes are partitioned into $k$ disjoint regions using $k$-Means clustering based on their coordinates.

$$\mathcal{P} = \{R_1, R_2, \ldots, R_k\}, \quad R_k \subseteq V, \quad R_i \cap R_j = \varnothing \text{ for } i \neq j \tag{7}$$

The quotient graph is defined as $G_q = (V_q, E_q)$, with $|V_q| = k$, where each node $u_k \in V_q$ corresponds to a region $R_k$. The embedding of a region node is computed by mean pooling over the embeddings of its constituent landmarks:

$$\mathbf{h}_{R_k} = \frac{1}{|R_k|} \sum_{v_i \in R_k} \mathbf{h}_i \tag{8}$$

Similarly, the coordinates of the region node are aggregated:

$$\mathbf{p}_{R_k} = \frac{1}{|R_k|} \sum_{v_i \in R_k} \mathbf{p}_i \tag{9}$$

Edges in the quotient graph are defined using $k$-nearest neighbors ($k$-NN), based on the Euclidean distances between region coordinates $\mathbf{p}_{R_k}$, ensuring connectivity reflects the spatial proximity of facial regions. In practice, this process is applied independently to each graph, producing a quotient graph with $k$ nodes per input graph.

The quotient graph is processed with two EdgeConv layers, as described in (5). These layers enrich the region embeddings by capturing higher-order dependencies, such as correlations between the eyes and mouth or between the eyebrows and cheeks. The refined embeddings improve the effectiveness of the downstream classification.

### 3.4 Graph-Level Prediction from Region Embeddings

To produce a unified representation of the facial expression for classification, the region embeddings of the quotient graph are processed to capture inter-region interactions and aggregated into a permutation-invariant graph-level embedding. This phase ensures that the model captures unique facial patterns, such as coordinated movements of the eyes and mouth, critical for distinguishing expressions like happiness or sadness in FER.

The refined region embeddings are aggregated via global additive pooling [16] to obtain a permutation-invariant graph-level representation:

$$\mathbf{h}_G = \sum_{R \in V_q} \mathbf{h}_R \tag{10}$$

Additive pooling ensures robust aggregation across regions, accommodating variations in region contributions across samples.

A linear layer maps $\mathbf{h}_G$ to a probability distribution over 8 expression classes (neutral, happy, sad, surprise, fear, disgust, anger, contempt), enabling the model to classify facial expressions effectively.

### 3.5 Theoretical Achievements

The hierarchical architecture of the proposed GLaRE model achieves both computational efficiency and improved generalization for facial expression recognition by incorporating a quotient graph to represent regional facial structures. This approach reduces the computational complexity of message passing while maintaining the essential expressive information contained in facial landmarks.

By coarsening the fine-level facial landmark graph into a higher-level quotient graph, the number of nodes and edges involved in computation is significantly reduced. Given a fine-level graph with $N$ nodes and $E$ edges, GNN-based message passing typically incurs $\mathcal{O}(E)$ time complexity per layer. After coarsening, the coarse graph with $N' \ll N$ nodes and $E' \ll E$ edges reduces the message passing complexity to $\mathcal{O}(E')$, which leads to significant speedup.

Beyond efficiency, the quotient graph groups landmarks into semantically meaningful facial regions (such as eyes, mouth, or eyebrows), thereby capturing higher-order interactions that are crucial for distinguishing expressions. This regional abstraction improves robustness to variations in landmark positions and local noise, enhancing the generalization ability of the model across diverse facial configurations.

Finally, a permutation-invariant pooling mechanism aggregates region-level embeddings into a compact graph-level representation, ensuring consistent predictions irrespective of the input node ordering. This hierarchical approach allows the model to jointly exploit fine-grained landmark relations and coarse regional structures, resulting in a more efficient and reliable solution for facial expression recognition.

## 4 Result & Discussion

This section presents the experimental evaluation of the proposed model GLaRE, which employs a hierarchical quotient graph to achieve computational efficiency and robustness to variations in facial landmark configurations. We provide a detailed analysis and discussion of classification accuracy, loss function, and comparisons with state-of-the-art models, including CNN-based approaches and graph-based methods, in this section.

### 4.1 Dataset Description

We have conducted experiments on two benchmark facial expression datasets, AffectNet and FERG-DB, covering both real-world and stylized facial data.

**AffectNet** [24] has emerged as the largest facial expression dataset to date, containing over one million facial images collected from the internet using emotion-related keywords in six different languages. Out of these, approximately 450,000 images have been manually annotated. The dataset defines eleven categories, including six basic emotions (anger, disgust, fear, happiness, sadness, surprise), as well as neutral, contempt, none, uncertain, and non-face. For our experiments, we have randomly created a subset consisting of the six basic expressions along with neutral, yielding a total of 83,901 training images and 1500 validation images.

**FERG-DB** [2] (Facial Expression Research Group Database) consists of 555,767 images of six stylized characters exhibiting seven types of facial expressions: the six basic emotions and neutral. Since the facial figures in this dataset are cartoon-style and lack realistic facial texture and contour details, extracting accurate facial landmarks has been particularly challenging. This difficulty arises because many landmark detection models are trained on real human faces and often fail to generalize to synthetic domains.

### 4.2 Baseline Methods

The gACNN [18] proposes an occlusion-aware attention mechanism that combines local and global facial features to improve facial expression recognition under occluded conditions. LDL-ALSG [5] introduces soft label supervision and

models the correlations between samples to enhance the learning of facial expression features. OADN [6] incorporates a landmark-guided attention branch that directs the model to focus on non-occluded facial regions, thereby improving recognition performance under occlusions. FERGCN [19] employs a deep graph-based architecture that includes a feature extraction module, a graph convolution module, and a graph matching mechanism to recognize facial expressions.

### 4.3 Additional Details

The model has been trained using the PyTorch-Geometric [8] framework for 200 epochs. The Adam optimizer [15] has been employed with a learning rate of 0.0001 and default betas of $(0.9, 0.999)$. Cross-entropy loss has been used as the objective function, which is well-suited for multi-class classification tasks. It measures the discrepancy between the predicted class probabilities and the true labels, encouraging the network to assign higher confidence to the correct class. The loss for a single sample is defined as (11).

$$\lambda_{CE} = -\sum_{c=1}^{C} y_c log(\hat{y}_c) \tag{11}$$

where $C$ is the total number of classes, $y_c$ is the ground truth indicator, and $\hat{y}_c$ is the predicted probability for class $c$.

### 4.4 Performance Evaluation and Analysis

Table 1: Emotion-wise accuracy and precision (with standard deviation) of the GLaRE model on AffectNet Dataset subset

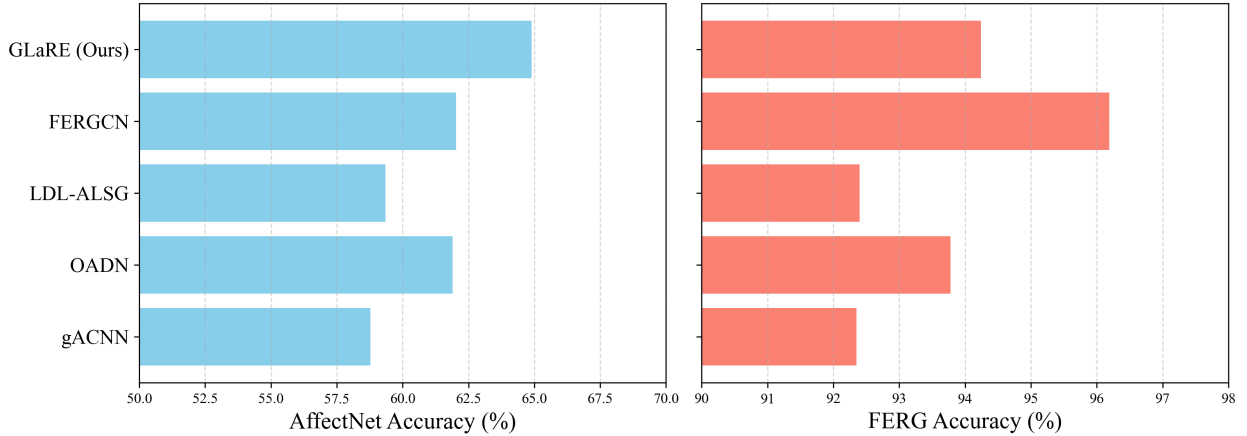| Emotion | Samples per Class | Accuracy (%) | Precision (%) |
|---------|------------------|--------------|---------------|
| Fear | 1190 | $65.29 \pm 0.82$ | $66.07 \pm 0.90$ |
| Anger | 1236 | $65.13 \pm 0.76$ | $63.79 \pm 0.85$ |
| Joy | 1157 | $65.25 \pm 0.79$ | $64.04 \pm 0.88$ |
| Sadness | 1255 | $65.50 \pm 0.84$ | $66.18 \pm 0.91$ |
| Surprise | 1228 | $66.12 \pm 0.87$ | $66.56 \pm 0.92$ |
| Neutral | 1218 | $64.94 \pm 0.80$ | $66.19 \pm 0.89$ |
| Disgust | 1139 | $63.62 \pm 0.85$ | $63.07 \pm 0.93$ |



Figure 4: Model performance comparison on AffectNet and FERG datasets.

Table 1 summarizes the number of samples per class and their corresponding accuracies, while Fig. 4 provides a visual representation of this performance comparison with baseline models. These results indicate that GLaRE is highly effective on real-world facial data, especially when landmark features are accurate and consistent.

As seen in Table 2, the proposed GLaRE model achieves the best accuracy of 64.89% on this reduced AffectNet subset, outperforming all other baselines, including gACNN [18], OADN [6], LDL-ALSG [5], and FERGCN [19]. To

Table 2: Accuracy (%) comparison of various models on AffectNet and FERG datasets. GLaRE outperforms baselines on AffectNet and demonstrates competitive performance on FERG.

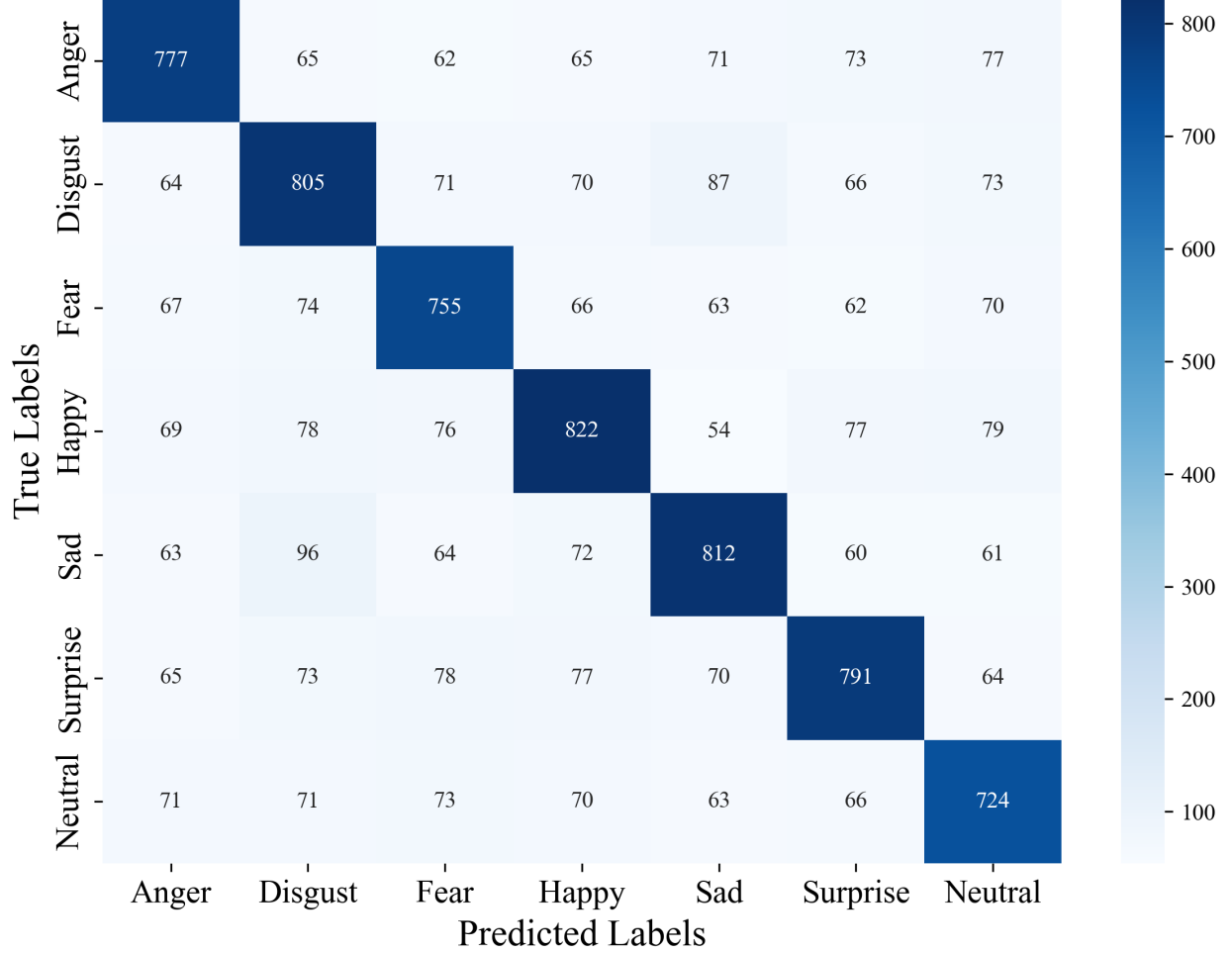| Method | AffectNet | FERG |
|---|---|---|
| gACNN | 58.78±0.31 | 92.35±0.27 |
| OADN | 61.89±0.28 | 93.78±0.34 |
| LDL-ALSG | 59.35±0.25 | 92.40±0.29 |
| FERGCN | 62.03±0.22 | **96.19±0.21** |
| **GLaRE** | **64.89±0.23** | 94.24±0.38 |



Figure 5: Confusion matrix of Emotion-wise accuracy on AffectNet Dataset.

assess stability, ten independent runs have been conducted, and a low standard deviation of ±0.27% has been observed, confirming the consistency of the model's performance. This demonstrates the model's ability to generalize well, even with limited training data. Table 1 further confirms that the performance remains consistent across different categories, with only minor variations. Performance comparisons have been conducted across all models on the same dataset, and statistical significance has been established through paired t-tests, yielding $p < 0.01$.

In contrast, while evaluating on the FERG dataset, which consists of synthetic or artificially-rendered characters, GLaRE achieves 94.24% accuracy which is slightly lower than FERGCN's 96.19%. This reduction is primarily due to the difficulty in extracting consistent facial landmarks from stylized characters. Unlike real human faces, cartoon faces often lack anatomically grounded features, making landmark-based graph construction noisy and less reliable. Since

GLaRE relies on high-quality 3D landmark features to build expressive graphs, its performance is naturally impacted under such conditions.

CNN-based models for FER often require millions of parameters, ranging from 2 to 18 million, which results in high computational demand. In contrast, graph-based approaches can achieve competitive results with far fewer parameters. The proposed GLaRE model contains only 44,424 parameters, highlights that GLaRE operates with significantly fewer parameters than established CNN-based methods while remaining more expressive than earlier graph-based baselines such as FERGCN [19]. Such a lightweight design facilitates faster training, reduced memory consumption, and greater suitability for deployment in real-world scenarios.

## 4.5 Ablation Studies

We have conducted ablation studies to evaluate the contribution of the quotient graph and the choice of region granularity in our model. In the first setting, we have removed the quotient graph and directly processed landmark-level graphs. This has resulted in a significant increase in computational time and a noticeable drop in recognition accuracy, highlighting the efficiency and effectiveness of the quotient graph representation. Furthermore, we have varied the number of facial regions when constructing the quotient graph. The results presented in the Table 3 and Fig. 6 have shown that both decreasing and increasing the number of regions beyond a certain point leads to performance degradation. The highest accuracy has been achieved when the number of regions is set to eight or nine, suggesting that this range provides an optimal balance between structural compression and information preservation.

Table 3: Ablation study on the number of regions in the quotient graph. Accuracy (%) and inference time per batch (ms) are reported.

| Number of Regions | Accuracy (%) | Time (ms) |
|:---:|:---:|:---:|
| 5 | 51.32±0.28 | 30.5 |
| 6 | 57.47±0.31 | 31.2 |
| 7 | 63.58±0.26 | 31.9 |
| 8 | **64.89±0.23** | 32.2 |
| 9 | **64.72±0.25** | 33.4 |
| 10 | 62.12±0.29 | 38.0 |

Additional ablation studies have been conducted to evaluate the contribution of different feature types. Two restricted variants of the input have been examined: (i) using only the 3D position vector of facial landmarks, and (ii) using only the appearance feature vector of size 16 extracted from local patches around the landmarks. As shown in Table 4, both variants have resulted in significantly lower performance compared to the joint featurization, demonstrating that the positional and appearance cues are complementary in nature. The position-only model has achieved $33\%$ accuracy, but it lacks the texture details necessary for reliable expression discrimination. Conversely, the appearance-only model has yielded only $56\%$ accuracy, highlighting that geometric information plays a crucial role in stabilizing the representation. In contrast, the joint featurization has achieved the best performance, confirming the necessity of integrating both modalities.

Table 4: Ablation study on different feature choices for facial expression recognition. The results indicate that joint featurization of position and appearance vectors has provided a substantial improvement over using either feature alone.

| Feature Type | Accuracy (%) |
|:---|:---:|
| Position vector only (3D) | 56 |
| Appearance vector only (16D) | 33 |
| Joint (Position + Appearance) | **64.89** |

## 5 Conclusion & Future Work

To conclude, GLaRE has been introduced as a hierarchical graph-based model that integrates fine-grained landmark information with high-level regional interactions for emotion recognition. The model has effectively captured both local geometry and global facial structure, while maintaining computational efficiency through graph coarsening. The approach has achieved performance surpassing existing baselines on a balanced subset of AffectNet and has demonstrated strong generalization on FERG-DB, despite the challenges posed by stylized facial textures.
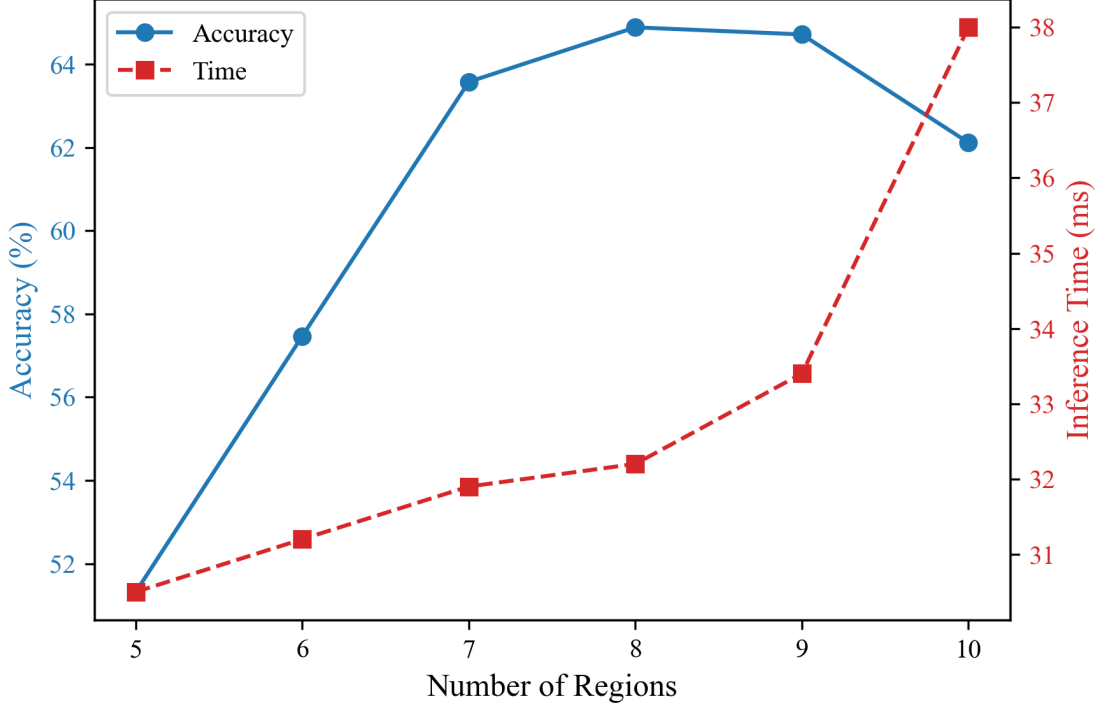
Figure 6: Ablation study on the number of quotient graph regions. The accuracy improves as the number of regions increases, peaking at 8–9 regions, beyond which performance slightly declines and the inference time increases.

Future work has been directed toward improving landmark detection on non-realistic faces, extending the framework to spatio-temporal modeling for video-based recognition, and adopting self-supervised strategies to reduce reliance on extensive annotated datasets.

# References

[1] Ai, W., Shou, Y., Meng, T., Li, K.: Der-gcn: Dialog and event relation-aware graph convolutional neural network for multimodal dialog emotion recognition. IEEE Transactions on Neural Networks and Learning Systems (2024)

[2] Aneja, D., Colburn, A., Faigin, G., Shapiro, L., Mones, B.: Modeling stylized character expressions via deep learning. In: Asian Conference on Computer Vision. pp. 136–153. Springer (2016)

[3] Bulat, A., Tzimiropoulos, G.: How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In: International Conference on Computer Vision (2017)

[4] Chen, B., Guan, W., Li, P., Ikeda, N., Hirasawa, K., Lu, H.: Residual multi-task learning for facial landmark localization and expression recognition. Pattern Recognition **115**, 107893 (2021)

[5] Chen, S., Wang, J., Chen, Y., Shi, Z., Geng, X., Rui, Y.: Label distribution learning on auxiliary label space graphs for facial expression recognition. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 13984–13993 (2020)

[6] Ding, H., Zhou, P., Chellappa, R.: Occlusion-adaptive deep network for robust facial expression recognition. In: 2020 IEEE International Joint Conference on Biometrics (IJCB). pp. 1–9. IEEE (2020)

[7] Fang, W., Love, P.E., Luo, H., Ding, L.: Computer vision for behaviour-based safety in construction: A review and future directions. Advanced Engineering Informatics **43**, 100980 (2020)

[8] Fey, M., Lenssen, J.E.: Fast graph representation learning with pytorch geometric. arXiv preprint arXiv:1903.02428 (2019)

[9] Gao, H., Yüce, A., Thiran, J.P.: Detecting emotional stress from facial expressions for driving safety. In: 2014 IEEE International Conference on Image Processing (ICIP). pp. 5961–5965. IEEE (2014)

[10] Hajiabolhassan, H., Taheri, Z., Hojatnia, A., Yeganeh, Y.T.: Funqg: Molecular representation learning via quotient graphs. Journal of chemical information and modeling **63**(11), 3275–3287 (2023)

[11] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

[12] Hossain, M.S., Muhammad, G.: Emotion recognition using deep learning approach from audio–visual emotional big data. Information Fusion **49**, 69–78 (2019)

[13] Huang, C., Jiang, F., Han, Z., Huang, X., Wang, S., Zhu, Y., Jiang, Y., Hu, B.: Modeling fine-grained relations in dynamic space-time graphs for video-based facial expression recognition. IEEE Transactions on Affective Computing (2025)

[14] Kaushik, P., Yadav, K., Gopika, B., Kumar, N., Kaushik, M.M., et al.: Deep learning for driver vigilance and road safety. In: 2024 4th International Conference on Advancement in Electronics & Communication Engineering (AECE). pp. 285–290. IEEE (2024)

[15] Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)

[16] Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks (2017), `https://arxiv.org/abs/1609.02907`

[17] Kuo, C.M., Lai, S.H., Sarkis, M.: A compact deep learning model for robust facial expression recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition workshops. pp. 2121–2129 (2018)

[18] Li, Y., Zeng, J., Shan, S., Chen, X.: Occlusion aware facial expression recognition using cnn with attention mechanism. IEEE transactions on image processing **28**(5), 2439–2450 (2018)

[19] Liao, L., Zhu, Y., Zheng, B., Jiang, X., Lin, J.: Fergcn: facial expression recognition based on graph convolution network. Machine Vision and Applications **33**(3), 40 (2022)

[20] Liu, M., Li, S., Shan, S., Wang, R., Chen, X.: Deeply learning deformable facial action parts model for dynamic expression analysis. In: Asian conference on computer vision. pp. 143–157. Springer (2014)

[21] Liu, S., Huang, S., Fu, W., Lin, J.C.W.: A descriptive human visual cognitive strategy using graph neural network for facial expression recognition. International Journal of Machine Learning and Cybernetics **15**(1), 19–35 (2024)

[22] Marín-Morales, J., Llinares, C., Guixeres, J., Alcañiz, M.: Emotion recognition in immersive virtual reality: From statistics to affective computing. Sensors **20**(18), 5163 (2020)

[23] Meng, Z., Liu, P., Cai, J., Han, S., Tong, Y.: Identity-aware convolutional neural network for facial expression recognition. In: 2017 12th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2017). pp. 558–565. IEEE (2017)

[24] Mollahosseini, A., Hasani, B., Mahoor, M.H.: Affectnet: A database for facial expression, valence, and arousal computing in the wild. IEEE Transactions on Affective Computing **10**(1), 18–31 (2017)

[25] Pampouchidou, A., Pediaditis, M., Kazantzaki, E., Sfakianakis, S., Apostolaki, I.A., Argyraki, K., Manousos, D., Meriaudeau, F., Marias, K., Yang, F., et al.: Automated facial video-based recognition of depression and anxiety symptom severity: cross-corpus validation. Machine Vision and Applications **31**(4), 30 (2020)

[26] Srinivasan, S., Raja, R., Jehan, C., Murugan, S., Srinivasan, C., Muthulekshmi, M.: Iot-enabled facial recognition for smart hospitality for contactless guest services and identity verification. In: 2024 11th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions)(ICRITO). pp. 1–6. IEEE (2024)

[27] Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M.: Dynamic graph cnn for learning on point clouds. ACM Transactions on Graphics (tog) **38**(5), 1–12 (2019)

[28] Yu, J., Zheng, Y., Wang, L., Wang, Y., Xu, S.: Cross-modal facial expression recognition with global channel-spatial attention: Modal enhancement and proportional criterion fusion. In: Proceedings of the Computer Vision and Pattern Recognition Conference. pp. 5707–5714 (2025)

[29] Zhang, T., Zheng, W., Cui, Z., Zong, Y., Li, Y.: Spatial–temporal recurrent neural network for emotion recognition. IEEE transactions on cybernetics **49**(3), 839–847 (2018)