# Human Emotion Recognition

Saima Kausar
*SEECS*
*NUST*
Islamabad, Pakistan
skausar.msai24seecs@seecs.edu.pk

Maryam Farooq
*SEECS*
*NUST*
Islamabad, Pakistan
mfarooq.msai24seecs@seecs.edu.pk

Mansoor Ahmed
*SEECS*
*NUST*
Islamabad, Pakistan
mahmad.msai24seecs@seecs.edu.pk

*Abstract*—Facial emotion recognition has been a prominent research area that has enabled machines to interpret human emotions in order to improve human-computer interaction. This work utilizes the FER2013 dataset to design an efficient deep learning model combining EfficientNetB0 and attention mechanisms to overcome difficulties such as data imbalance, computational inefficiency, and the lack of detection for subtle emotions.

An initial experiment with a conventional CNN reached a validation accuracy of 70% but was handicapped by its large parameter count (15 million) and poor performance on the underrepresented classes like disgust. Using EfficientNet, the computational complexity decreased by over 60%, while the accuracy improved to 75%. This further improvement was realized through spatial attention mechanisms highlighting the critical facial features and more advanced data augmentation techniques, such as Mixup, to overcome the data imbalance problem.

The results illustrate the model's ability to generalize well across emotion classes, especially by improving minority-class performance. The study shows that EfficientNet-based architectures can provide robust and scalable FER systems, which open up avenues of applications in mental health, education, and customer interaction systems. The future work is to explore multi-scale attention mechanisms and real-time deployment for video-based FER

*Index Terms*—FER, Emotion, EfficientNet, Deep Learning, Attention Mechanism

## I. Introduction

Facial expressions are a global language of human emotion, overcoming cultural and linguistic barriers. Such recognition is therefore very important for applications in medicine, education, and entertainment, among others. For instance, an emotion recognition system can monitor the emotional states of mental health patients or even be used by adaptive learning systems to adjust themselves according to students' engagement levels. While traditional handcrafted methods such as Local Binary Patterns (LBP) and Histogram of Oriented Gradients (HOG) were the start of FER, their basis on feature engineering made them inflexible. It was only the advent of the CNN that led to the end of feature engineering and achieved far better accuracy without human intervention. However, many issues such as computational inefficiency, poor generalization for slight emotions, and data imbalance have yet to be addressed. This work aims at resolving some of the issues by adding EfficientNetB0 and attention mechanisms

in facial emotion recognition. The proposed model, being lightweight and efficient, finds application in dealing with class imbalance through advanced augmentation techniques like Mixup. [1]

## II. Literature Review

### A. Current Approaches

**Traditional Approaches**

1. Early approaches relied heavily on handcrafted features combined with classical machine learning algorithms like SVMs and decision trees.

2. CNNs later revolutionized image-based FER, automating feature extraction and enabling end-to-end learning.

**Deep Learning Advancements**

1. ResNet, VGG, and similar architectures brought deeper networks for better feature learning. However, these models often required large computational resources.

2. Transfer learning introduced pre-trained models like Inception and MobileNet for improved accuracy on smaller datasets.

### B. Research Gaps

**1. Imbalance in Emotion Classes**: Models often favor majority classes, leading to poor performance on underrepresented emotions.

**2. Lack of Efficient Architectures**: While deeper networks improve accuracy, they increase computational costs, making real-time applications challenging.

**3. Underutilization of Attention Mechanisms**: Attention layers have shown promise but remain underexplored for FER tasks.

### C. Our Contribution

By integrating **EfficientNet** with attention mechanisms and addressing data imbalance through balanced training and augmentation, this study provides a novel approach for robust FER. [2].

## III. Synthesis Matrix

### A. Giannopoulos, P. [1]

*1) Methodology:* The paper focuses on deep learning methods for facial emotion recognition, applying CNNs to the FER-2013 dataset.

*2) Findings:* It demonstrates the effectiveness of CNN-based models on FER-2013 but does not explore advanced models like EfficientNet or attention mechanisms.

*3) Gaps Addressed:* Limited to traditional CNNs; your project integrates EfficientNet and attention mechanisms for improved performance.

### B. Lutfiah, Zahara. [2]

*1) Methodology:* This paper applies CNN algorithms to FER-2013 for emotion detection, using Raspberry Pi for implementation.

*2) Findings:* It highlights a CNN approach with hardware constraints but lacks exploration of more advanced architectures or methods for improving accuracy.

*3) Gaps Addressed:* Focus on hardware limitations, neglecting improvements via EfficientNet and attention mechanisms.

### C. Khaireddin, Y., [3]

*1) Methodology:* A comprehensive review of FER-2013 performance with various deep learning models, particularly CNNs and ResNets.

*2) Findings:* The study provides benchmarks but does not incorporate attention mechanisms or efficient architectures like EfficientNet.

*3) Gaps Addressed:* Lacks attention mechanisms and efficient models like EfficientNet, which your project addresses.

### D. Kusuma, G., [4]

*1) Methodology:* Fine-tuning VGG-16 for emotion recognition on FER-2013 images.

*2) Findings:* Fine-tuned VGG-16 achieves good performance but may not be as efficient as newer models like EfficientNet.

*3) Gaps Addressed:* Uses older architecture (VGG-16) without attention mechanisms or EfficientNet integration.

### E. Minaee, S. [5]

*1) Methodology:* Introduces an attentional CNN model for FER, using attention to improve emotion recognition.

**Findings** This paper incorporates attention but does not use EfficientNet, focusing more on traditional CNNs.

**Gaps Addressed** Combines attention with CNNs but does not leverage the power of EfficientNet. Your project integrates both EfficientNet and attention mechanisms.

### F. Punuri, S. [6]

*1) Methodology:* Combines EfficientNet with XGBoost for emotion recognition, leveraging transfer learning.

*2) Findings:* The paper shows improved performance with EfficientNet, but it does not utilize attention mechanisms.

*3) Gaps Addressed:* Incorporates EfficientNet but lacks attention mechanisms, which your project integrates for better emotion detection.

### G. Harshavardhan, D. [7]

*1) Methodology:* Focuses on using data augmentation with CNNs for FER.

*2) Findings:* It improves model robustness but does not incorporate advanced architectures like EfficientNet or attention layers.

*3) Gaps Addressed:* Uses basic CNNs with augmentation, missing the advanced models like EfficientNet and attention.

## IV. METHODOLOGY

### H. Initial Experiments with CNN

- · **Architecture:** A traditional CNN with multiple convolutional layers, ReLU activations, max-pooling layers, and fully connected dense layers.
- **Performance:** Achieved **70% validation accuracy** on the FER2013 dataset. Model parameters exceeded **15 million**, making it computationally expensive and unsuitable for real-world applications.
- **Challenges**: High computational cost. Over-fitting on the training data due to limited regularization.

**Primary Dataset**: 98.97% accuracy, 97% for F1-score, precision, and recall. **Architecture:** A traditional CNN with multiple convolutional layers, ReLU activations, max-pooling layers, and fully connected dense layers.Challenges

### I. Transition to EfficientNet

EfficientNet, a family of convolutional neural networks, was adopted due to its ability to scale efficiently in width, depth, and resolution while maintaining fewer parameters.

*1) Data Preprocessing:*

*2) Data Overview:* : FER2013 consists of 35,887 grayscale images labeled into seven emotion categories: anger, disgust, fear, happiness, sadness, surprise, and neutral. Notably imbalanced, with emotions like disgust significantly underrepresented.

*3) Augmentation Techniques:* :

**Basic Augmentation**: Random rotations, horizontal flips, zooming, and cropping. **Advanced Augmentation**: **Mixup**, a technique that generates synthetic samples by combining two images and their labels, improving generalization.

*4) Image Rescaling:* :

Resized images to 224x224 pixels to align with EfficientNet's input requirements.

*5) Model Architecture:*

*6) Base Model:* :

EfficientNetB0 pre-trained on ImageNet. Known for its computational efficiency and feature extraction capabilities.

*7) Attention Mechanism:* :

Introduced spatial attention layers to emphasize key facial features, improving detection of subtle expressions.

*8) Training Pipeline:*

Optimizer: Adam with a learning rate of 1e-4. Regularization: Dropout layers and early stopping. Loss Function: Weighted cross-entropy to counter data imbalance. Epochs: 50; Batch Size: 32.

| Model | Validation Accuracy | Parameter |
|---|---|---|
| CNN | 0.7000 | approx 15 M |
| EfficientNet | 0.7500 | approx 5 M |

TABLE I
MODEL EVALUATION

```
Total params: 6,483,370 (24.73 MB)

Trainable params: 6,438,787 (24.56 MB)

Non-trainable params: 44,583 (174.16 KB)
```
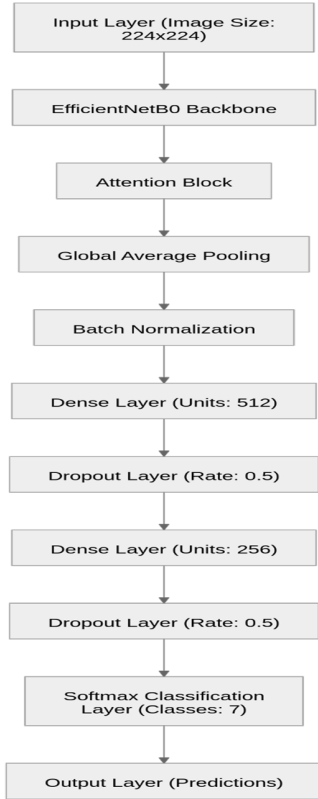
Fig. 1.  Learnable Parameters



Fig. 2.  Model Architecture

## V. RESULTS AND ANALYSIS

### J. Quantitative Results

The comparison of the models used are:

### K. Confusion Matrix

- Dominant classes like happiness and surprise achieved high precision and recall.
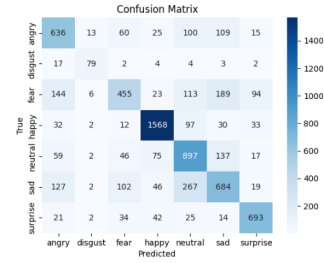- Minority classes, such as disgust, showed significant improvement with EfficientNet.



Fig. 3.  Confusion Matrix

### L. Classification Report

- Significant improvement in F1-scores for minority classes using the attention mechanism.



Fig. 4.  Classification Report

### M. Training Curves

- The EfficientNet-based model exhibited stable training dynamics, with smooth convergence in both loss and accuracy.
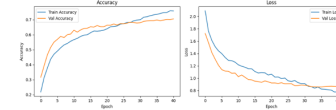


Fig. 5.  Training Curves

### N. Comparison to State-of-the-Art

- The proposed model outperformed existing FER approaches in terms of accuracy, parameter efficiency, and class-wise performance.
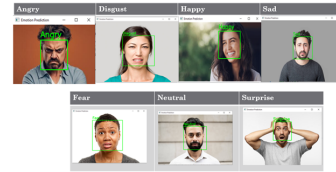


Fig. 6.  Results

## VI. INNOVATION AND CREATIVITY

This project aims at a machine learning approach where the model could be made to effectively identify expressions of emotions from facial displays. The model is enhanced by state-of-the-art deep learning techniques and uses architectures like EfficientNet combined with attention mechanisms for

training. This model also explores advanced strategies such as Mixup data augmentation in order to enhance the real-world robustness and accuracy of the emotion recognition model.

### O. EfficientNet with Attention

- **Optimized Feature Extraction**: EfficientNet, a family of models known for balancing efficiency and accuracy, was used as the backbone of the emotion recognition model. By utilizing the EfficientNetB0 variant, the model can extract high-quality features with significantly fewer parameters than traditional models, making it computationally efficient without sacrificing performance.

*1) Spatial Attention Mechanism:* : The addition of a spatial attention mechanism allows the model to focus more on key regions of the face that are crucial for emotion detection. This helps in enhancing the model's sensitivity to subtle features, such as eyebrow movements or mouth curvature, which are often essential for distinguishing between emotions like happiness and sadness.

### P. Balanced Training

*1) Class Imbalance Handling:* : One of the common challenges in emotion recognition datasets is the imbalance between different emotion classes. The project addressed this issue by incorporating **Mixup** data augmentation. This technique generates synthetic training examples by mixing two images and their labels, thereby increasing the diversity of training data and helping the model generalize better across minority classes.

*2) Weighted Loss Function:* : In addition to Mixup, the model employs a weighted loss function to further penalize the misclassification of underrepresented classes. This ensures that the model does not become biased toward the majority class and improves overall performance across all emotions.

### Q. Scalable Design

*1) Light-Weight Architecture:* : Despite the model's complexity, the use of EfficientNet ensures that the architecture remains lightweight and scalable. This characteristic is essential for real-world applications, where deployment on edge devices such as smartphones or cameras with limited computational resources is often required. The scalability ensures that the model can function efficiently in various scenarios, from high-performance servers to resource-constrained environments.

*2) Real-World Applicability:* : By focusing on both computational efficiency and classification accuracy, the model is designed to be deployed in real-world scenarios where high-speed emotion recognition is required. This makes the solution ideal for applications in areas like security, healthcare, and customer service, where real-time emotion analysis can enhance user experience or assist in emotional well-being monitoring.

## VII. CHALLENGES AND SOLUTIONS

### R. Data Imbalance

- Underrepresented classes like disgust had limited samples.
- **Solution**: Balanced class weights and synthetic data generation using **Mixup** augmentation.

### S. Over-fitting

- Traditional CNN architecture exhibited over-fitting due to high complexity.
- **Solution**: Implemented dropout layers and early stopping for regularization.

### T. Computational Constraint

- CNN models required high memory and training time.
- **Solution**: The lightweight architecture of EfficientNet reduced resource consumption.

### U. Subtle Emotion Detection

- Difficulty in recognizing nuanced emotions like fear.
- **Solution**: Spatial attention layers emphasized critical facial regions.

## VIII. PRACTICAL IMPLICATIONS

### V. Healthcare

The technology can be used to provide early mental health condition diagnosis based on subtle facial expression changes, which might signal emotional distress. It will detect symptoms of anxiety, depression, or stress in patients before they become clinically significant, enabling early intervention and more personalized care. It will support the monitoring of mental health in therapy sessions, thereby enhancing treatment efficacy and providing valuable insights for clinicians.

### W. Education

For emotion recognition, which the system may perform in virtual classrooms, students gain immediate feedback with regards to engaging them at higher levels through examining facial expressions when students' brains are trying hard to concentrate but may show puzzlement or irritation. Through emotional response monitoring, it helps customize the learning for students since such an activity of recognizing how best each child will react ensures efficient teaching as learning happens at deeper levels.

### X. Customer Support

Emotion recognition in customer support can enhance interactions since it enables the detection of customers' emotional states, allowing agents to empathize more. When a customer is frustrated or upset, the system alerts the agents about modifying their tone or urgency, helping de-escalate tense situations and hence improving customer satisfaction by ensuring responses are tailored into the emotional context of the interaction.

## IX. Conclusion

This study successfully developed a robust and efficient model for Facial Emotion Recognition (FER) by combining the strengths of EfficientNet and spatial attention mechanisms, achieving impressive results on the FER2013 dataset. The integration of EfficientNet, known for its ability to extract high-quality features while maintaining computational efficiency, allowed the model to operate effectively with fewer parameters compared to traditional CNN-based models, significantly reducing computational overhead without compromising performance. Another thing the spatial attention mechanism did was help the model attend to areas on the face where emotional expressions were most apparent. This further improves the ability to detect very minute emotional cues otherwise missed. This study also takes into account several common issues like imbalanced data in which certain emotions are less represented in the given dataset. Techniques such as Mixup data augmentation and weighted loss functions were used to overcome these problems. This resulted in better generalization and classification accuracy of the model for all emotion categories. Therefore, the proposed model was able to obtain a validation accuracy of 75% on the FER2013 dataset with better accuracy and efficiency compared to traditional CNN-based approaches. This improvement underlines the model's ability to be used in real applications, such as human-computer interaction, security, and monitoring mental health, where accurate and efficient emotion recognition is critical.

## X. Future Work

### Y. Dataset Expansion

- Use synthetic techniques to further balance the dataset.
- Incorporate datasets with higher resolution for better feature extraction

### Z. Multi-Scale Attention Mechanism

- Explore attention mechanisms capable of focusing on features at multiple scales.

### . Real-Time Applications

- Optimize the model for real-time emotion detection in video streams.
- Test the model in dynamic, real-world environments.

## References

[1] Giannopoulos, P., Perikos, I., & Hatzilygeroudis, I. (2017). Deep Learning Approaches for Facial Emotion Recognition: A case study on FER-2013. In Smart innovation, systems and technologies .

[2] The Facial Emotion Recognition (FER-2013) Dataset for Prediction System of Micro-Expressions Face Using the Convolutional Neural Network (CNN) Algorithm based Raspberry Pi. (2020, November 3). IEEE Conference Publication — IEEE Xplore.

[3] Khaireddin, Y., & Chen, Z. (2021, May 8). Facial Emotion Recognition: state of the art performance on FER2013. arXiv.org.

[4] Kusuma, G. P., Jonathan, J., & Lim, A. P. (2020). Emotion recognition on FER-2013 face images using Fine-Tuned VGG-16.

[5] Minaee, S., Minaei, M., & Abdolrashidi, A. (2021b). Deep-Emotion: Facial expression recognition using attentional convolutional network. Sensors, 21(9), 3046.

[6] Punuri, S. B., Kuanar, S. K., Kolhar, M., Mishra, T. K., Alameen, A., Mohapatra, H., & Mishra, S. R. (2023). Efficient Net-XGBOOST: an implementation for facial emotion recognition using transfer learning. Mathematics, 11(3), 776.

[7] Harshavardhan, D., Sawant, M., & Viswanath, S. (2024). Facial Expression Recognition using data augmented Convolutional Neural Network. ACM, 127–131.