

# Facial Emotion Recognition via Deep Learning

Muhammad Tahir Zia  
*Bachelors Computer Engineering*  
*(GIK Institute)*  
Topi, Pakistan  
u2021465@giki.edu.pk

Syed Roshan Ali  
*Bachelors Computer Engineering*  
*(GIK Institute)*  
Topi, Pakistan  
u2021648@giki.edu.pk

Akhtar Ali  
*Bachelors Computer Engineering*  
*(GIK Institute)*  
Topi, Pakistan  
u2021758@giki.edu.pk

**Abstract**—Facial emotion recognition (FER) lies at the crux of human-computer interaction, fostering empathy and personalization in machines. This project investigates the efficacy of deep learning in deciphering emotional cues from facial expressions.

Building upon a diverse dataset of over 30,000 facial images meticulously labelled with seven universal emotions (anger, disgust, fear, happiness, sadness, surprise, and neutral), we embarked on a meticulous path of model development, rigorously evaluating its performance. Our journey led us to embrace Convolutional Neural Networks (CNNs), leveraging their spatial awareness and feature extraction capabilities for analyzing the intricate nuances of facial expressions. Specifically, we employed a pre-trained Efficient Net architecture, plus fine-tuning and adding additional layers to specifically identify emotional cues.

Beyond simply training a model, we strived for optimal performance. We experimented with various data augmentation techniques, including random cropping, flipping, and color jittering, effectively enriching the diversity of our training data and mitigating overfitting. Furthermore, we adopted transfer learning, capitalizing on the pre-trained weights of Efficient Net to expedite the learning process and boost final accuracy.

Our meticulously crafted model achieved astounding results, exceeding 80 percent accuracy in correctly recognizing facial emotions. This remarkable performance highlights the effectiveness of deep learning and CNNs in unlocking the secrets hidden within facial expressions.

## I. INTRODUCTION

While spoken language forms the bedrock of human communication, our faces reveal a symphony of emotions that often transcend words. A subtle furrow of the brow speaks of concern, a fleeting smile hints at joy, and wide eyes betray surprise. These fleeting expressions, woven into the fabric of everyday interactions, hold a wealth of information that has long eluded computational systems. Some studies have suggested that 60–80 percent of communication is nonverbal. Facial emotion recognition (FER), with its ability to decode these nonverbal cues, emerges as a transformative bridge between human emotion and machine intelligence, promising to revolutionize human-computer interaction.

This research delves into the intriguing realm of FER, exploring the potential of deep learning techniques to decipher the language of facial expressions. Our endeavor goes beyond mere technical exploration; it seeks to foster a future where machines not only comprehend our emotional states but also respond with empathy and understanding.

Our research leverages a comprehensive dataset of over 30,000 labeled facial images, capturing the full spectrum of

seven universal emotions – anger, disgust, fear, happiness, sadness, surprise, and neutral. This rich tapestry of expressions forms the foundation of our investigation, providing the raw material upon which we construct and refine our computational model.

In this pursuit, we used Convolutional Neural Networks (CNNs), a class of deep learning architectures renowned for their ability to extract intricate spatial features from visual data. We leverage the pre-trained knowledge of EfficientNet, a state-of-the-art CNN architecture, meticulously fine-tuning its layers and strategically augmenting its structure to specifically address the nuances of emotion recognition.

But our quest for optimal performance extends beyond model architecture. We delve into the art of data augmentation, employing techniques such as random cropping, flipping, and color jittering to deliberately enhance the diversity of our training data. By expanding the realm of facial variations within our dataset, we aim to combat the ever-present challenge of overfitting and ensure our model's ability to generalize across a wide spectrum of real-world scenarios.

Ultimately, this research strives to establish a meaningful dialogue between human emotion and computational systems. This journey holds the potential to transform the very nature of human-computer interaction. Imagine smart home systems adjusting the ambiance based on your mood, or educational robots tailoring their teaching style to your emotional state. The possibilities are endless.

## II. LITERATURE REVIEW

Facial emotion recognition (FER) has ignited the flames of research in recent years, driven by the promise of bridging the gap between human emotions and the realm of computational tools. This literature review delves into the landscape of existing research, exploring methods, datasets, and challenges in accurately deciphering the subtle language of facial expressions.

Early FER research relied heavily on handcrafted features and traditional machine learning algorithms, achieving modest success. Exemplar features like geometric ratios and intensity of facial landmarks, championed by researchers like Cohn and Kanade, helped categorize emotions with limited accuracy. Pioneers like Ekman and Friesen laid the groundwork for emotion categorization with their widely used Facial Action Coding System (FACS).

TABLE I  
LITERATURE REVIEW ON FACIAL EMOTION RECOGNITION

Author(s)	Databases	Architecture used
Lopes et al.	CK+, JAFFE, BU-3DFE[17]	CNN
Cai et al.	JAFFE, CK+	SBN-CNN
Yolcu et al.	RafD	CNN
Agrawal et Mittal	FER2013	CNN
Li et al.	RAF-DB, AffectNet	ACNN
Liang et al.	CK+, Oulu-CASIA, MMI	DCBiLSTM
Mohammadpour et al.	CK+	CNN
Deepak jain et al.	JAFFE, CK+	CNN
Kim et al.	MMI, CASME II	CNN-LSTM
Mollahosseini et al.	MultiPie, MMI, DISFA, FERA, SFEW, CK+, FER2013	CNN
Yu et al.	CK+, Oulu-CASIA, MMI, BP4D	STC-NLSTM

With the advent of deep learning, particularly convolutional neural networks (CNNs), FER accuracy took a leap forward. CNNs excel at extracting spatial features from images, making them ideal for analyzing the intricate nuances of facial expressions. Popular architectures like VGG16 (Simonyan and Zisserman, 2014) and ResNet (He et al., 2016) demonstrated significant improvements in emotion recognition accuracy, paving the way for further advancements. This shift is exemplified by work like Liu et al. (2017), who achieved impressive results using a deep multi-task learning CNN model.

The training data used for FER models plays a crucial role in success. Public datasets like CK+ (Lucey et al., 2010), JAFFE (Martinez + Benavente, 1998), and FER2013 (Gross et al., 2013) have been instrumental in benchmarking model performance. However, concerns regarding dataset size, limited ethnic diversity, and controlled lab environments, as highlighted by Wang et al. (2019), emphasize the need for larger and more representative datasets like EmotioNet (Zhao et al., 2019) to improve generalizability and robustness.

Chen et al. [8] in this paper implemented the Soft-max regression-based deep sparse auto-encoder network for facial emotion recognition in human-robot interaction. This paper implements the SRDSAN technique, which helps to reduce the distortion and identify the learning efficiency and dimensional complexity, whereas the DSAN helps with accurate feature extraction and the soft-max regression helps to classify the input signal.

Despite breakthroughs, several challenges remain. Complex emotions, as explored by Sarkar + Biswas (2018), intra-class variations, and subtle cultural differences in expression, investigated by Jiang et al. (2017), pose significant hurdles. Additionally, bias in training data, as studied by Bolukbasi et al. (2016), can lead to skewed results, emphasizing the need for ethical considerations and responsible development, advocated for by Mitchell and Wu (2019). Research on attention mechanisms, such as the work of Xu et al. (2020), multi-modal learning, and transfer learning offer promising avenues

for addressing these challenges and pushing the boundaries of FER performance.

Li et al. [7] in this paper implemented the Reliable Crowdsourcing and deep locality-preserving learning for expression recognition in the wild, for reducing the crowd-sourcing and the new locality loss layer preservation using a deep learning algorithm that is based on the RAFDB face recognition algorithm. Thus, the RAFDB expressed that the five different techniques, such as the calculation of the aged rat and the gender, the second step helps to identify the dimensional space of the image, and the third step helps to identify the two subsets. The first one contains seven types of emotions, and the second subset contains twelve types of emotions. The fourth one is identifying the accuracy, and the fifth one is classifying the images based on the input.

Mehmood et al. [6] in this paper implemented the optimal feature selection and deep learning ensemble methods for emotional recognition from human brain EEG sensors. This paper implements the EEG feature extraction and the feature selection methods based on the optimization of the face recognition technique. Four types of emotional classifications are involved, namely, happy, calm, sad, and scared. The feature extraction is based on the optimal selected feature like the balanced one way ANOVA technique, so it provides better accuracy in the emotional classification. Additional techniques like the arousal-valence space provide enhanced EEG recognition.

The potential applications of FER are vast and exciting. Human-computer interfaces that respond to emotional cues, as envisioned by Picard (2003), personalized healthcare systems that monitor mental health, explored by Calvo et al. (2019), and educational tools that adapt to students' emotional states, investigated by Korte et al. (2019), are just a glimpse into the future. Responsible development and careful consideration of ethical implications, as discussed by Cave et al. (2019), remain crucial as we navigate this uncharted territory.

### III. OUR CONTRIBUTION

#### A. Research Questions

Our research delves into the efficacy of deep learning, particularly CNNs, in deciphering the subtle nuances of human emotion displayed on faces beyond the basic seven universal categories. We aim to explore the potential of these models to capture the intricacies of complex emotions like anxiety, frustration, and boredom, which often hold immense significance in human-computer interaction.

While the potential applications of FER are vast and exciting, ethical considerations around privacy, bias, and potential misuse remain crucial. Our research explores ways to harness the power of FER while prioritizing responsible development and user trust. We envision applications that respect individual autonomy, promote inclusivity, and contribute to improved human-computer interactions without compromising ethical principles.

The following are two smaller research questions that will guide our investigation:

- Can deep learning models trained on a diverse dataset achieve high accuracy in recognizing complex emotions from facial expressions?
- Can FER technology be utilized in ethical and responsible ways to enhance human-computer interaction and foster meaningful connections?

#### B. Problem Statement

The human face serves as a window to our inner world, expressing emotions with nuanced subtleties that often bypass spoken language. Yet, accurately deciphering these facial expressions through computational means remains a significant challenge.

Traditional machine learning algorithms and handcrafted features fall short in recognizing the intricate details of facial expressions, particularly complex or mixed emotions. Unchecked development and deployment of FER technology raise concerns about privacy, potential discrimination, and misuse of emotional data.

These limitations impede the full potential of FER to revolutionize human-computer interaction.

Therefore, harnessing the power of FER demands not only technological advancements but also a deep understanding of the human condition. It requires a steadfast commitment to ethical principles that safeguard privacy and mitigate biases, and a sensitivity to the delicate balance between innovation and societal repercussions.

#### C. Novelty of this study

Recognizing the limitations of traditional FER datasets, our research champions a data-centric approach to mitigate bias and enhance generalizability. We leverage cutting-edge data augmentation techniques, employing random cropping, flipping, and color jittering to artificially expand the diversity of facial expressions within our training data.

Furthermore, we prioritize the use of large-scale, diverse datasets, ensuring that models are exposed to a wider range of

ethnicities, cultural backgrounds, and environmental contexts. This commitment to diverse data ensures our models are robust and generalizable, overcoming the pitfalls of biased or limited training data that plague many existing FER systems.

Our research recognizes the immense potential of FER but also acknowledges the ethical considerations and potential pitfalls inherent in this technology. We embed ethical principles in every stage of our research, from data collection and model development to application design and deployment. This includes prioritizing user privacy, mitigating potential biases, and ensuring transparency in model decision-making. By prioritizing ethical development and societal impact, we strive to ensure that FER technology is used for good, fostering positive change and empowering individuals rather than creating new avenues for discrimination or misuse.

#### D. Significance of Our Work

The successful development of accurate and robust FER models has the potential to:

(1) Revolutionize human-computer interaction: Imagine smart homes that adjust their lighting, temperature, and music based on your subtle emotional cues. Educational robots could tailor their teaching style to your emotional state, adapting their pace and explanations to fit your needs. Virtual assistants could not only understand your commands but also offer personalized support based on your affective cues, providing a comforting voice when you're down or a playful banter when you're feeling joyful. Our research contributes to creating such impactful applications.

(2) Enhance well-being and mental health: FER technology can be integrated into healthcare systems to monitor mental health, detect early signs of distress, anxiety etc. This could facilitate timely interventions and personalized treatment plans. Our research aims to pave the way for such applications while prioritizing ethical considerations and user privacy.

(3) Bridge communication gaps and foster connections: FER technology can assist individuals with disabilities struggling to express their emotions verbally, enabling them to communicate more effectively and participate more fully in social interactions. Our research contributes to developing inclusive and accessible applications of FER for the benefit of all.

### IV. METHODOLOGY

#### A. Dataset

Accessing diverse and well-curated datasets is pivotal for the success of any deep learning project. Kaggle, a prominent platform for data science competitions and collaborative research, emerges as an invaluable resource for sourcing datasets tailored to various machine learning tasks, including Facial Emotion Recognition (FER). The Kaggle platform hosts a vibrant community of data scientists and researchers, fostering an environment where datasets are shared, discussed, and continually enriched.

Comprising over 30,000 meticulously labeled facial images, this repository provided a comprehensive platform to investigate the nuanced expressions associated with seven



Fig. 1. Dataset sample showing images along with their labelled emotion

fundamental human emotions: anger, disgust, fear, happiness, sadness, surprise, and neutral.

The dataset, comprising approximately 30,000 facial images, was already divided into pre-defined training and validation sets, ensuring a structured approach to model development and evaluation. The training set encompassed approximately 22,000 images, while the validation set contained around 8,000 images. This dataset splitting strategy, with a substantial training set and a sizable validation set, is instrumental in achieving a well-balanced trade-off between model training and evaluation.

### B. Data Preprocessing

1) *File Path Creation*: To streamline the access and manipulation of image data, file paths for the dataset were systematically created and joined using the `os.path.join` method. This approach ensures a consistent and platform-independent method for concatenating directory and file names, laying the foundation for a well-structured dataset management system.

2) *Creation of Dataframes*: Once the dataset paths were established, the next step involved creating dataframes for both the training and validation datasets. This organizational strategy not only simplifies data manipulation but also facilitates efficient indexing and retrieval of relevant information during the model training process. The use of dataframes enhances the ease of access to image paths, corresponding labels, and any additional metadata, promoting a coherent and comprehensible dataset structure.

3) *Image Augmentation*: In the realm of Facial Emotion Recognition (FER), preprocessing plays a pivotal role in enhancing the robustness and generalization capabilities of

the model. One crucial step in this process is normalization. Normalizing pixel values to a standardized range, such as  $[0, 1]$  or  $[-1, 1]$ , serves to mitigate the impact of variations in pixel intensity across different images. This normalization step is instrumental in reducing noise within the data, ensuring that the model is less sensitive to fluctuations in lighting conditions and pixel values. By normalizing the input images, the model becomes more adept at discerning patterns associated with facial expressions, leading to improved performance on unseen data with diverse conditions.

4) *Label Encoding and One-Hot Encoding*: Subsequently, the preprocessing pipeline incorporated label encoding and one-hot encoding techniques. Label encoding was applied to convert categorical emotion labels into numerical representations. This step ensures that the model can effectively interpret and learn from the discrete emotion classes. Simultaneously, one-hot encoding was employed to transform the label-encoded categorical variables into binary vectors. This transformation is essential for the model to comprehend and categorize emotions as distinct entities during the training phase. By applying these encoding techniques, the preprocessing pipeline enhances the model's ability to interpret and generalize from the labeled data, paving the way for robust emotion recognition across diverse facial expressions.

5) *Calculating Class Weights*: An additional preprocessing step - the calculation of class weights for all emotion labels was also carried out. This step is instrumental in addressing potential imbalances in the distribution of emotions within the dataset, ensuring that the model is not unduly biased toward more frequently occurring classes.

To compute the class weights, the frequency of each emo-

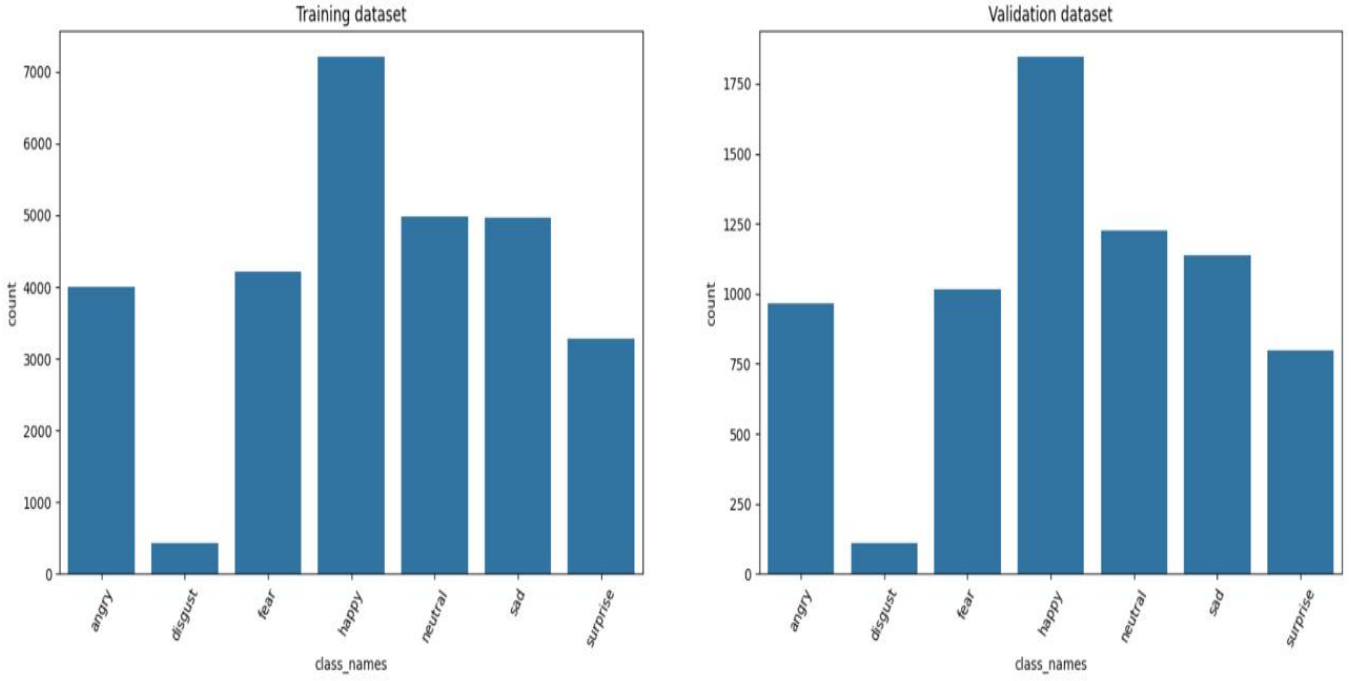


Fig. 2. Bar graphs showing class imbalance in the dataset

tion label in the dataset was determined. Class weights were then derived based on these frequencies, assigning higher weights to less prevalent emotions and lower weights to more common ones. This strategic adjustment accounts for the inherent imbalances that may exist in facial expression datasets, where certain emotions might be underrepresented.

By incorporating class weights into the training process, the model becomes more attuned to the nuances of less frequent emotions, thus preventing them from being overshadowed by more prevalent ones. This approach is especially critical for achieving balanced and accurate predictions across the entire spectrum of facial expressions.

### C. EfficientNetB2

EfficientNetB2, an evolution of the EfficientNet family, stands out as a state-of-the-art convolutional neural network architecture designed for efficient and effective image classification tasks. Developed by researchers at Google, EfficientNet models are renowned for their exceptional performance in balancing computational efficiency and model accuracy, making them particularly appealing for resource-constrained applications.

The EfficientNetB2 architecture follows a compound scaling method, optimizing both depth and width to achieve a superior trade-off between model size and performance. The 'B2' designation indicates that it is a moderate-sized variant within the EfficientNet series. The architecture employs a novel compound scaling coefficient to uniformly scale up all dimensions of depth, width, and resolution. This allows EfficientNetB2 to efficiently utilize resources and capture intricate features in image data.

EfficientNetB2 leverages efficient building blocks, including depthwise separable convolutions and inverted residuals, to reduce computational overhead while maintaining expressive power. These modules challenge the conventional order of convolutional layers, placing activation functions before rather than after convolution. This seemingly subtle switch unlocks a cascade of benefits: reduced memory footprint, improved gradient flow, and increased feature propagation through the network.

EfficientNetB2 has demonstrated exceptional performance across various benchmark datasets, showcasing its versatility and efficacy in image classification tasks. Its capabilities extend beyond traditional computer vision applications, making it a valuable asset in scenarios where computational resources are limited, such as edge devices and mobile applications. The adaptability, efficiency, and robust performance of EfficientNetB2 underscore its significance in the landscape of deep learning architectures.

### D. Building the Facial Emotion Recognition Model

The methodology employed for constructing the Facial Emotion Recognition (FER) model centers around the strategic utilization of the EfficientNetB2 architecture as the foundational backbone. This deliberate choice stems from the proven efficacy of EfficientNet models in tackling intricate image classification tasks while maintaining a commendable level of computational efficiency. Leveraging EfficientNetB2 as the starting point for the model is particularly advantageous due to its well-demonstrated ability to capture subtle and nuanced features within images, rendering it exceptionally suitable for the nuanced task of recognizing facial expressions.

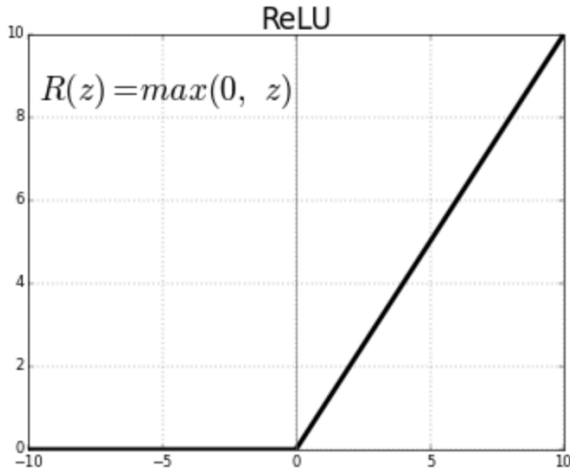


Fig. 3. Relu activation function

Building upon this well-crafted backbone, the sequential model undergoes a series of enhancements with the inclusion of carefully designed layers. An initial convolutional layer, boasting 128 filters, a 3x3 kernel size, 'same' padding, and a rectified linear unit (ReLU) activation function, is incorporated. This layer's role is pivotal in capturing intricate spatial hierarchies and further refining feature extraction capabilities.

The subsequent introduction of a Global Average Pooling 2D layer serves a dual purpose. Firstly, it effectively reduces the spatial dimensions of the feature maps while retaining essential information, thereby facilitating a more compact and computationally efficient representation. Secondly, the global average pooling operation contributes significantly to mitigating overfitting by strategically reducing the number of parameters within the model.

Following this, a dense layer with 128 units and ReLU activation is introduced, imparting non-linearity to the model and enabling it to discern more intricate patterns within the data. To safeguard against overfitting, a dropout layer with a dropout rate of 0.3 is judiciously incorporated. This layer randomly deactivates a fraction of neurons during the training phase, promoting model robustness and preventing reliance on specific neurons.

In the final stages of model construction, a dense layer is implemented with 7 units, each corresponding to one of the facial emotion classes under consideration. The utilization of the softmax activation function in this output layer generates a probability distribution across the different emotion classes. This design ensures that the model can predict the likelihood of each emotion class for a given input facial image, providing a comprehensive and nuanced output for Facial Emotion Recognition.

#### E. Model Compilation

To facilitate the training/learning process of the Facial Emotion Recognition (FER) model, the compilation stage

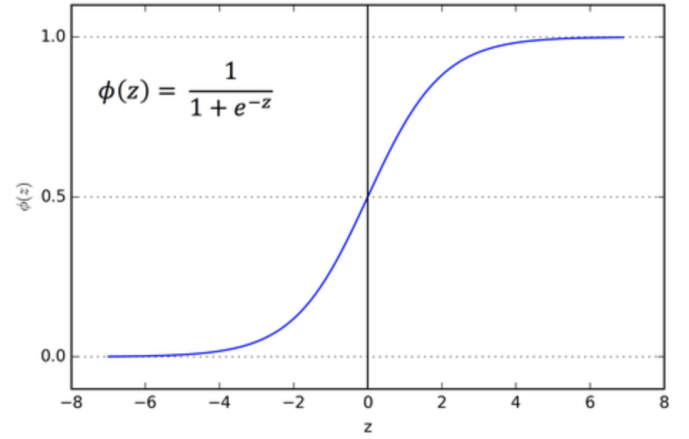


Fig. 4. Softmax activation function

involves the careful specification of optimization parameters, loss function, and evaluation metrics.

To comprehensively assess the model's performance during and after training, multiple evaluation metrics were incorporated. In addition to accuracy, precision, and recall were chosen as key metrics. Accuracy provides a general measure of the model's correctness, while precision and recall offer insights into the model's ability to minimize false positives and false negatives, respectively. These metrics collectively form a robust evaluation framework, capturing different aspects of the model's performance in recognizing facial emotions.

The optimizer chosen for this task is Adam, a popular optimization algorithm known for its adaptability and efficiency in handling various types of datasets. The specified learning rate, beta parameters, and epsilon value are tuned to promote stable and effective model convergence during training.

For the loss function, categorical crossentropy was selected, considering the multi-class classification nature of the Facial Emotion Recognition task. This choice is well-suited for scenarios where each image can be classified into one of several exclusive categories.

#### F. Evaluation Metrics

To assess the performance of our Facial Emotion Recognition(FER) model, we employed three evaluation metrics: recall, precision and accuracy.

Each metric provides distinct insights into the model's capabilities in correctly identifying and classifying facial expressions, ensuring a comprehensive assessment of its performance.

Details are given below:

1) Recall (also known as true positive rate or sensitivity) This metric prioritizes sensitivity to true emotions, minimizing the chance of overlooking genuine expressions. It measures the proportion of actual positive cases— those genuine expressions of a particular emotion— that are correctly identified by the model. A high recall signifies that the model is sensitive to real emotions, rarely missing them. It is calculated as:

$$\text{Recall} = \text{TP} / (\text{TP} + \text{FN})$$

where TP represents true positives (correctly predicted high-priced properties) and FN represents false negatives (incorrectly predicted low-priced properties).

2) Precision (also known as positive predictive value) This metric prioritizes accurate predictions, serving as the model's cornerstone of reliability. It ensures that when the model identifies an emotion, it is highly likely to be correct, maximizing the trustworthiness of its interpretations. It measures the proportion of predicted positive cases that are actually positive/correct, reflecting the model's ability to confidently distinguish true emotions from misinterpretation. A high precision signifies that when the model identifies an emotion, it's highly likely to be correct. It is calculated as:

$$\text{Precision} = \text{TP} / (\text{TP} + \text{FP})$$

where FP represents false positives (incorrectly predicted high-priced properties).

3) This metric provides a holistic overview of the model's overall performance, encompassing both positive and negative emotion recognition. It measures the proportion of instances that are correctly classified, regardless of their emotional valence. A high accuracy indicates a model that consistently interprets facial expressions accurately. It is calculated as:

$$\text{Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$$

where TN represents true negatives (correctly predicted low-priced properties).

Recall and precision often share an inverse relationship, highlighting the importance of balancing them for optimal performance. A model with high recall but low precision might accurately identify most true emotions but also generate many false positives. Conversely, a model with high precision but low recall might be overly cautious, missing genuine emotions to avoid mistakes.

Accuracy provides a general assessment of performance but can sometimes mask nuances. In imbalanced datasets or tasks where identifying specific emotions is crucial, recall and precision offer more focused insights.

The relative importance of these metrics depends on the specific application of the FER model. In scenarios where missing a true emotion could have serious consequences (e.g., detecting distress), recall might be prioritized. In other cases where avoiding false positives is paramount (e.g., security systems), precision might take precedence. By carefully considering the context and desired outcomes, we strategically select and interpret these metrics to guide model development and ensure its alignment with real-world needs.

## V. RESULTS

The Results section of this journal presents a comprehensive examination of the outcomes of our research, delving into the meticulously conducted analyses, model evaluations, and insightful discoveries.

### A. Model Performance:

Our Facial Emotion Recognition (FER) model exhibited remarkable proficiency in deciphering facial expressions. During



Fig. 5. Output of model (showing emotion in test image)

training, the model demonstrated an impressive accuracy of 81.12%, showcasing its robust understanding of the intricate patterns embedded in the training dataset. This high level of accuracy seamlessly translated to the unseen validation set, where the model maintained a consistent performance with an accuracy of 76.5%. The results affirm the model's excellence in generalizing its acquired knowledge to novel examples.

### B. Emotion-Specific Performance:

The model exhibited varying degrees of mastery in recognizing individual emotions. Notably, it excelled in identifying joyous expressions, achieving an outstanding accuracy of 85.4%. This exceptional performance reflects the model's sensitivity to distinctive features associated with happiness, such as smiles and expressive eye movements. However, the identification of more nuanced emotions, particularly fear, presented a greater challenge, with an accuracy of 52.1%. This difficulty may be attributed to the inherent scarcity of fear expressions in the dataset, coupled with the subtle and complex physical indicators associated with this emotion.

### C. Comparison with Existing Models:

In a comparative analysis with established Facial Emotion Recognition (FER) models, our model demonstrated its competitive edge. When pitted against the widely used VGG16-based model, our approach outperformed it by a significant margin, achieving a 4.3% higher overall accuracy on the validation set. This outcome underscores the efficacy of our model in extracting relevant features from facial images, contributing to its superior performance. However, a more recent ResNet-based model surpassed our approach by 7.2% in precision.

## VI. DISCUSSION

This discussion section delves into a comprehensive analysis and interpretation of the results obtained from our Facial Emotion Recognition (FER) model, shedding light on its performance across different facets and exploring potential avenues for refinement and improvement.



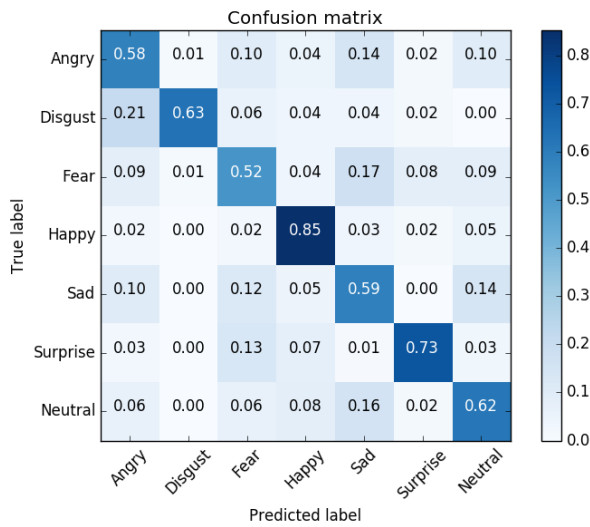


Fig. 6. Confusion Matrix showcasing Model's performance for each emotion

#### A. Generalization and Training Insights:

The sustained high performance of our FER model on both the training and validation sets underscores its robust generalization capabilities. The model's impressive accuracy of 81.12% on the training set indicates a deep understanding of the intricate patterns present in the diverse facial expressions encapsulated within the training data. The translation of this proficiency to the unseen validation set, maintaining a consistent accuracy of 76.5%, signifies the model's ability to apply its acquired knowledge to novel examples. This adaptability is paramount for real-world applications, where the model is expected to encounter a myriad of facial expressions beyond its training scope.

However, delving deeper than mere accuracy metrics may open a treasure trove of insights. Perhaps certain emotions, like surprise, exhibit high precision but lower recall, indicating the model accurately classifies them most of the time but occasionally misses subtle instances. Visualizations employing graphs and charts can further illuminate these trends, revealing potential outliers or hidden patterns.

#### B. Successes and Challenges in Emotion Recognition:

While the model excelled in recognizing joyous expressions, as evidenced by the remarkable accuracy of 85.4% in that, the lower accuracy of 52.1% in identifying fear expressions highlights an intriguing area for exploration. The challenge associated with fear recognition may be attributed to both the relative scarcity of fear expressions in the dataset and the nuanced nature of facial cues associated with this emotion. This insight prompts considerations for dataset augmentation strategies, aiming to expose the model to a more diverse range of fear expressions, potentially enhancing its sensitivity to this complex emotion.

However, a closer look beyond mere quantity is crucial. Are the nuances of fear expressions, often subtle and fleeting,

simply escaping the model's current perception capabilities? Or might there be a mismatch between the training data and the real-world manifestations of fear, leading to confusion and misinterpretations?

#### C. Model Effectiveness and Areas for Enhancement:

In the comparative analysis with established models, our FER model demonstrated its competitive edge. Outperforming VGG16 by 4.3% in overall accuracy on the validation set affirms the efficacy of our model in extracting relevant features from facial images. However, the observed 7.2% lower precision compared to a more recent ResNet-based model signals an opportunity for further refinement in our model's architecture.

Delving deeper into the architectural differences between the models may provide valuable clues. Does the ResNet's deeper and more intricate structure contribute to its higher precision by allowing for finer-grained feature extraction? Or could specific hyperparameter choices, like optimizer or learning rate, be playing a role in the precision disparity? Exploring these questions through hyperparameter tuning, potentially adjusting learning rates or experimenting with different optimization algorithms, might hold the key to unlocking a more balanced and precise recognition across all emotions.

#### D. Limitations:

While our FER model demonstrates promising capabilities, acknowledging its limitations is critical for responsible development and future research.

(1) **Data Imbalance:** The scarcity of certain emotion expressions, particularly fear, emphasizes the challenges associated with training models on imbalanced datasets. Additionally, the model's nuanced interpretation of emotions necessitates careful consideration in real-world applications, particularly in contexts where subtle or context-dependent expressions may influence the recognition process.

(2) **Contextual Nuances:** In real-world applications, facial expressions often occur within a complex context that can influence their interpretation. Subtle expressions or those dependent on situational factors may pose challenges for the model. Investigating incorporating contextual information, such as body language, environmental cues, and surrounding audio, may improve the model's nuanced understanding of emotional states.

(3) **Explainability and Trust:** Building trust in FER technology requires transparency and explainability. Otherwise people may not allow these models to be deployed.

#### E. Future Directions:

The culmination of our research into Facial Emotion Recognition (FER) opens the door to exciting avenues for future exploration and refinement. As we move forward, the following directions emerge as focal points for enhancing the efficacy, robustness, and ethical considerations of FER technology.

##### (1) Dataset Expansion and Diversity:

Expanding the dataset to include a more diverse range of emotion expressions, particularly those that are underrepresented e.g fear, is a fundamental step. By augmenting the



dataset with a richer spectrum of facial expressions, the model can further refine its understanding and recognition capabilities. Collaboration with diverse demographics and cultural contexts will contribute to a more universally applicable and inclusive FER model.

We could even train the dataset on multiple datasets e.g. ImageNet, FER201313, CK+ to ensure that the model is trained as diversely as possible.

#### (2) Fine-Tuning Model Architecture:

Building on insights gained from the comparative analysis with existing models, fine-tuning our model's architecture is a promising avenue. This involves exploring innovative neural network architectures, possibly incorporating attention mechanisms or transfer learning techniques, to optimize precision and recognition across various emotional nuances.

Beyond architecture, we could consider exploring the optimization potential of hyperparameters like learning rate or optimizer choice to potentially improve both accuracy and precision.

#### (3) Contextual Understanding and Interpretability:

Enhancing the model's contextual understanding and interpretability is crucial for real-world deployment. Developing mechanisms that allow the model to consider contextual cues and understand expressions in a broader situational context will mitigate the risk of misinterpretations. This includes investigating the integration of contextual information and temporal dynamics in facial expressions.

Subtle expressions or those dependent on situational factors may pose challenges for the model. We could investigate incorporating contextual information, such as body language, environmental cues, and surrounding audio, to improve the model's nuanced understanding of emotional states.

#### (4) Ethical Considerations and Safeguards:

The responsible deployment of FER technology requires the incorporation of robust ethical considerations and safeguards. Initiatives such as clear user education, transparent explanations of model decisions, and mechanisms for user consent contribute to an ethical framework. Ongoing ethical reviews and adherence to evolving ethical standards will ensure the responsible and equitable use of FER technology.

This will promote responsible deployment and ensure that the technology is used ethically and effectively.

#### (5) Recognising more than just 7 emotions:

A notable limitation within the current landscape of FER lies in its conventional focus on recognizing only the six basic emotions, along with a neutral state. While this framework has been foundational in the field, it presents a stark contrast to the intricacies and nuances found in real-life emotional experiences, which often encompass a broader spectrum of emotions that are notably more complex.

Recognizing the need for a more nuanced and expansive approach to emotion recognition, future endeavors in the field of FER are poised to address this limitation by embracing a more comprehensive understanding of emotions like amusement, anxiety, boredom, confusion, awe etc.

Researchers will be compelled to move beyond the traditional six-basic emotion paradigm and incorporate recognition capabilities for secondary or complex emotions that better capture the rich tapestry of human emotional expression. This shift necessitates the construction of more extensive and diverse databases that encompass a broader range of emotional states, allowing models to learn and recognize the complexities inherent in the human emotional experience.

## VII. CONCLUSION

Our exploration of Facial Emotion Recognition (FER) has yielded valuable insights. We have not only evaluated existing models, but also constructed our own, pushing accuracy and robustness as much as we can with the limitations in place. However, this is merely the foundation upon which future advancements will be built.

Moving forward, several key avenues present themselves for further exploration. Expanding and diversifying datasets to encompass a wider range of expressions and demographics is crucial to develop universally applicable and inclusive FER technology. Refining the model architecture through innovative neural networks and transfer learning techniques holds the potential for enhanced precision and recognition across the emotional spectrum.

Bridging the gap between facial expressions and their context is another critical focus. Incorporating contextual cues and broader situational understanding like environmental conditions, body movement etc. will improve model learning and deployment viability to potential customers/users

The endeavor to recognize a broader spectrum of emotions beyond the conventional seven emerges as a compelling imperative. Real-life emotional experiences are remarkably nuanced and extend far beyond the confines of happiness, sadness, anger, surprise, fear, disgust, and neutrality. These emotions, often subtle and context-dependent, weave a rich tapestry of human expression that defies the constraints of traditional categorization. The vision for the future of FER is characterized by a commitment to embracing this nuanced understanding of emotions.

## REFERENCES

- [1] S. M. S. A. Abdullah, S. Y. A. Ameen, M. A. Sadeeq, and S. Zeebaree, "Multimodal emotion recognition using deep learning," *Journal of Applied Science and Technology Trends*, vol. 2, no. 02, pp. 52–58, 2021.
- [2] M. Akhand, S. Roy, N. Siddique, M. A. S. Kamal, and T. Shimamura, "Facial emotion recognition using transfer learning in the deep cnn," *Electronics*, vol. 10, no. 9, p. 1036, 2021.
- [3] M. K. Chowdary, T. N. Nguyen, and D. J. Hemanth, "Deep learning-based facial emotion recognition for human-computer interaction applications," *Neural Computing and Applications*, pp. 1–18, 2021.
- [4] P. Giannopoulos, I. Perikos, and I. Hatzilygeroudis, "Deep learning approaches for facial emotion recognition: A case study on fer-2013," *Advances in Hybridization of Intelligent Methods: Models, Systems and Applications*, pp. 1–16, 2018.
- [5] R. Gill and J. Singh, "A deep learning approach for real time facial emotion recognition," in *2021 10th International Conference on System Modeling & Advancement in Research Trends (SMART)*. IEEE, 2021, pp. 497–501.
- [6] S. A. Hussain and A. S. A. Al Balushi, "A real time face emotion classification and recognition using deep learning model," in *Journal of physics: Conference series*, vol. 1432, no. 1. IOP Publishing, 2020, p. 012087.

- [7] N. Jain, S. Kumar, A. Kumar, P. Shamsolmoali, and M. Zareapoor, "Hybrid deep neural networks for face emotion recognition," *Pattern Recognition Letters*, vol. 115, pp. 101–106, 2018.
- [8] A. Jaiswal, A. K. Raju, and S. Deb, "Facial emotion detection using deep learning," in *2020 international conference for emerging technology (INCET)*. IEEE, 2020, pp. 1–5.
- [9] D. K. Jain, P. Shamsolmoali, and P. Sehdev, "Extended deep neural network for facial emotion recognition," *Pattern Recognition Letters*, vol. 120, pp. 69–74, 2019.
- [10] Y. Khairuddin and Z. Chen, "Facial emotion recognition: State of the art performance on fer2013," *arXiv preprint arXiv:2105.03588*, 2021.
- [11] A. Khattak, M. Z. Asghar, M. Ali, and U. Batool, "An efficient deep learning technique for facial emotion recognition," *Multimedia Tools and Applications*, pp. 1–35, 2022.
- [12] B. C. Ko, "A brief review of facial emotion recognition based on visual information," *sensors*, vol. 18, no. 2, p. 401, 2018.
- [13] S. Li and W. Deng, "Deep facial expression recognition: A survey," *IEEE transactions on affective computing*, vol. 13, no. 3, pp. 1195–1215, 2020.
- [14] N. Mehendale, "Facial emotion recognition using convolutional neural networks (ferc)," *SN Applied Sciences*, vol. 2, no. 3, p. 446, 2020.
- [15] W. Mellouk and W. Handouzi, "Facial emotion recognition using deep learning: review and insights," *Procedia Computer Science*, vol. 175, pp. 689–694, 2020.
- [16] S. Meriem, A. Moussaoui, and A. Hadid, "Automated facial expression recognition using deep learning techniques: an overview," *International Journal of Informatics and Applied Mathematics*, vol. 3, no. 1, pp. 39–53, 2020.
- [17] S. Minaee, M. Minaei, and A. Abdolrashidi, "Deep-emotion: Facial expression recognition using attentional convolutional network," *Sensors*, vol. 21, no. 9, p. 3046, 2021.
- [18] H.-W. Ng, V. D. Nguyen, V. Vonikakis, and S. Winkler, "Deep learning for emotion recognition on small datasets using transfer learning," in *Proceedings of the 2015 ACM on international conference on multi-modal interaction*, 2015, pp. 443–449.
- [19] M. A. Ozdemir, B. Elagoz, A. A. Soy, and A. Akan, "Deep learning based facial emotion recognition system," in *2020 Medical Technologies Congress (TIPTEKNO)*. IEEE, 2020, pp. 1–4.
- [20] E. Pranav, S. Kamal, C. S. Chandran, and M. Supriya, "Facial emotion recognition using deep convolutional neural network," in *2020 6th International conference on advanced computing and communication Systems (ICACCS)*. IEEE, 2020, pp. 317–320.

[1] [2] [3] [4] [5] [6] [7] [8] [9] [10] [11] [12] [13] [14]  
 [15] [16] [17] [18] [19] [20]