

Real-Time Emotion Recognition using Deep Learning: A Comprehensive Approach

Khushi Jhunjhunwala
Computer Engineering
Department
Mukesh Patel School of
Technology Management and
Engineering

Devansh Banka
Computer Engineering
Department
Mukesh Patel School of
Technology Management and
Engineering

Haider Kachwalla
Computer Engineering
Department
Mukesh Patel School of
Technology Management and
Engineering

Dhirendra Mishra
Computer Engineering Department
Mukesh Patel School of Technology Management and Engineering

ABSTRACT

The work integrates a real-time facial emotion recognition system using the RAF-DB dataset, which contains 15,350 images annotated for seven basic emotions. Different aspects of dataset preparation and data augmentation, architectures to be used for the model, and performance evaluation on emotion classification are discussed within this study. A CNN is used in the emotion detection that can be drawn from live video streams in very high efficiency. Even though it is varied and contains images of various age ranges and demographic differences, the dataset nonetheless presents challenges like class imbalance and issues regarding the privacy and bias of data. Significant improvement in a model's generalization ability can be seen after using data augmentation techniques like rescaling, shearing, zooming, and horizontal flipping. Best accuracy was obtained at epoch 25, which was 78.32% for validation with three hidden layers and a filter of 3x3. This model therefore maintains an equilibrium between accuracy and real-time performance. The performance of the proposed model was tested across various demographic groups, and also with variations in the accuracy presented in most cases of the detection of fear and disgust. The research has compared its proposed model with existing models such as VGGNet and ResNet to give prominence to the computational efficiency, making it fit for real-time applications in situations where a dataset would be relatively small and limits the amount of computational resources available. The findings are important to put emotions recognition into practice in the real world, bringing to limelight the need for balancing model performance with ethical applicability in the field of AI.

Keywords

Real-time Emotion Detection, CNN, Deep Learning, Imbalanced Dataset, Facial Expression Recognition, Accuracy

1. INTRODUCTION

Facial emotion recognition is an area of affective computing, where it enables systems to know the human emotion from facial expressions. A very intelligent and interactive world will lead to an emotion detection through facial cues with a wide range of applications in health care, education, marketing, virtual assistance, and human computer interaction. The challenge here is in building models that are accurate enough

and fast enough for use in real-time applications, like video conferencing or surveillance systems and interactive systems.

Such variability and complexity in human facial expressions-person-to-person differ, varying from age to gender, ethnicity, light conditions etc-make FER the most challenging task. In fact, with large diversified datasets, such a kind of training is only possible that may achieve robust models capable of detecting emotions across variations. One major feature that distinguishes many datasets is that RAF-DB (Real-World Affective Faces Database) covers wide demographics in its compositions. It contains more than 15,000 labeled facial images, seven of the basic emotions: Anger, Disgust, Fear, Happiness, Sadness, Surprise, and Neutral. There is an age, gender, and ethnicity variation in faces as well. This dataset, specifically designed for emotion recognition, has come out to be one of the strongest resources in training systems for the detection of emotions.

The set of issues it may raise in the form of imbalance and ethical dilemmas with regard to privacy, however can cause potential interference in the model's performance in less-represented groups. Thus, addressing such problems while maintaining the real-time efficiency of the model will form the heart of the research.

A robust emotion recognition system developed using the RAF-DB dataset which emphasizes exploitation of deep convolutional neural networks, very much suitable for building an efficient model towards real-time emotion detection, capitalizes further on efficiency in the model with incorporation of preprocessing, data augmentation, face detection, and performance evaluation metrics to ensure diverse conditions of real-world scenarios. More implications of demographic imbalances and biases in the dataset are discussed in the article with challenges in overcoming them for accuracy and fairness in emotion classification. The results obtained from the final deployed model show effectiveness across various age, gender, and emotion classes in real-time applications for emotion recognition.

2. LITERATURE SURVEY

Recent advances in facial emotion recognition (FER) have demonstrated significant progress in real-time applications and deep learning approaches. Duncan et al. [1] developed a real-time FER system using convolutional neural networks (CNNs)

with transfer learning, achieving promising results through data augmentation and pre-trained models, though facing challenges with real-world accuracy drops and frame rate limitations. A comprehensive review by Mellouka and Handouzia [2] analyzed various deep learning architectures and databases for FER tasks, highlighting the superiority of deep learning approaches while noting constraints in emotion range and cross-database performance. Kaur and Kulkarni [3] expanded the scope to multimodal emotion recognition systems, demonstrating the effectiveness of CNN models while identifying challenges in person-independent recognition and real-time processing.

Further developments in the field include Ballesteros et al.'s [4] work on real-time emotion recognition tools utilizing multi-task CNNs and facial landmark recognition, though their approach struggled with nuanced emotions and contextual information processing. Abdullah et al. [5] proposed a multimodal approach combining different data types, achieving high accuracy (>95%) but facing challenges with computational complexity and dataset dependencies. In educational contexts, Vanneste et al. [7] explored computer vision applications for student engagement monitoring, revealing limitations in behavior analysis and cultural variations in expression interpretation.

Recent applications of emotion recognition have expanded into various domains. James et al. [8] developed an emotion-based music recommendation system using support vector machines (SVM), achieving 90-95% recognition rates under controlled conditions but struggling with environmental constraints. Onyema et al. [9] addressed computational efficiency in patient facial recognition through a modified ResNet architecture, while Khopkar and Saxena [10] implemented a straightforward CNN model using Keras, achieving moderate accuracy with a simple architecture. Agrawal and Mittal [11] investigated the effects of kernel size and filter numbers on classification accuracy, demonstrating the potential for optimizing model architecture while maintaining reasonable performance.

This comprehensive review reveals several critical gaps in current research. While existing approaches have made significant strides in emotion recognition, they consistently face challenges in real-world accuracy, computational efficiency, and generalization across diverse environments. Most notably, the reliance on pre-trained architectures and general datasets has led to suboptimal performance in specific emotion recognition tasks. Current models struggle with accuracy drops between controlled and real-world environments [1], limited emotion range detection [2], and person-independent recognition [3]. Additionally, the field lacks a balanced approach that combines high accuracy with computational efficiency, as evidenced by the trade-offs seen in recent studies [9, 10, 11].

These identified gaps directly support the proposed problem statement, which aims to address the limitations in accuracy and dataset imbalances through a specifically tailored CNN-based model. The focus on optimizing the architecture for facial emotion recognition, rather than relying on general pre-trained models, aligns with the observed need for more specialized and robust solutions in real-world applications.

3. DATASET — RAF-DB

3.1 Overview

This real-world database, placed at the user's disposal, is of a high quality, having been especially designed for facial emotion recognition, with 15,350 images of the face. It represents a wide age range, including newborns, the elderly, and various kinds

of captured emotions (as seen in Figure 1) that allow the study across demographic groups regarding affective facial analysis. Each image is tagged with one of the seven basic emotions and compound emotions, making the dataset a richer, granular one with real-world applicability.

The RAF-DB images are obtained from the internet and have different settings, lighting conditions, and poses. Therefore, the test split of the dataset comprises 3,068 images, with an 80:20 train-test split ratio, which is a common routine used in machine learning for generalization on unseen data. Therefore, the RAF-DB is perfect for any task needing high variability, and it can be chosen highly by researchers in the emotion recognition field.



Figure 1: Sample Dataset

3.2 Data Collection and Limitation

This dataset is created through a careful selection of images and, with professional annotators, labeling them. They aim to comprise images from diverse online sources that reflect real-world variations in lighting, facial orientation, and image resolution. Diversity comes at the cost of the approach-it introduces certain biases. For instance, images retrieved through the internet would represent the youth and celebrities over proportionally because they have excellent Online media. Further, annotations rely on human interpretation hence affect subjective biases and labeling inconsistencies.

The dataset is missing metadata, such as conditions under which every image is taken and every subject's actual age. Also, the data set consists of ambiguous gender classification unsure category, about 5%, indicating some uncertainty about gender classification with vision alone.

3.3 Dataset Imbalance

As reported, the RAF-DB dataset has an imbalance across a number of demographics and categories. For gender, the data set is 52% female, 43% male, which falls under unsure 5%. This could potentially affect the performance of the model when it comes to faces against male vs. female unless it is corrected using balancing techniques or bias mitigation.

77% Caucasian, 8% African-American, 15% Asian. Such an unbalanced racial profile may lead to poor performances when testing the model with less represented groups because emotional cue recognition varies culturally. Thus, the model will be less accurate on demographic groups other than Caucasians, which should then be addressed by balancing or targeted augmentation in training.

For example, resampling techniques or even synthetically augmenting might be useful in developing model robustness. More extensive datasets may serve as the source for transfer learning. Then there are the comparisons based on performance metrics from each demographic group, which could uncover what weaknesses are caused by imbalances to the point where the model fails to deliver.

3.4. Ethical Perspectives

There is an ethical issue raised by the utilization of RAF-DB, primarily concerning privacy and bias. Since the images were obtained from the internet, the dataset will contain faces of

individuals who have not granted consent to the use of photographs for facial recognition research. While RAF-DB is a publically accessible and commonly used dataset in research, it is important that privacy and possible autonomy of the person represented in this dataset are respected. A desirable research that utilizes RAF-DB would involve the conduct of transparency and privacy-preserving practices.

Further, this demographic imbalance potentially allows for the continued or even escalated perpetuation of bias. For example, deployment of a model trained on RAF-DB to real-world settings can result in performance biases against minority racial groups when the model was oversampled with Caucasian faces. Addressing this bias extends beyond technical balance. One needs to critically consider the point at which the constraints of the dataset begin to affect the deployment of the model in sensitive or high-stakes applications, such as surveillance applications or emotion analysis in clinical settings. The authors are encouraged to provide transparent reporting of bias and discussion of mitigation strategies toward fairness.

4. METHODOLOGY

4.1 Dataset Preparation

Real-time emotion recognition is required to be in a robust and structured dataset. The dataset for this project was arranged in two directories: training set and testing set. Images were standardized to the color mode of grayscale resolution, with a fixed resolution of 100x100 pixels, according to the specifications of the CNN model used. It obviates the computational overhead by not allowing the model to consider color, thereby giving more prominence to essential facial features rather than color, which are relatively less relevant for a classification of emotions.

Preprocessing involves resizing and normalizing pixel values between 0 and 1. This allows the input to be uniform across the training and testing datasets, thereby enhancing the model's ability to learn patterns well as well as generalize them properly with increased stability towards improvement in real-time predictions. The model is exposed to uniform data distributions regarding both training and testing, such that consistent and reliable results are assured through preprocessing images alike for every task..

4.2 Data Augmentation

For augmenting the ability of generalization of the model, a set of data augmentation techniques was applied. Data augmentation artificially increases the size of the training dataset by making variations to the images simulating real-world differences. The various augmentation strategies that were used in this project are tested in a trial-and-error process and observed for what effect they have on accuracy. The augmentation methods include rescaling, shearing, zooming, and horizontal flipping (as seen in Table 1):

TABLE I. Data Augmentation

Augmentation Techniques	Accuracy	Val Accuracy
Rescale	0.9495	0.7360
Rescale, Shear	0.9667	0.7746
Rescale, Shear, Zoom	0.8542	0.7746
Rescale, Shear, Zoom, Horizontal Flip	0.8172	0.7832

Rescale Only: The pixel values were normalized by a fixed factor, and the results were the base training at 94.95% and validation at 73.6%.

Rescaling and Shearing: The addition of the shear transformation increased the accuracy to 77.46%, demonstrating that slight geometric distortion does actually help in validation.

Rescale, Shear and Zoom: This combination maintained the classification accuracy at 77.46%. That is to say, zooming transformations did not have a significant effect on the model.

Rescale, Shear, Zoom, Horizontal Flip: Horizontal flipping improved validation accuracy up to 78.32%, which implied that mirror-image variations were useful in making the model better generalize facial asymmetries in emotional expressions.

The best combination of the final model applies rescale, shear, zoom, and a horizontal flip to yield maximum performance on the test set with minimal overfitting. These transformations expose the model to more facial orientations and lighting conditions and thus enhance its ability to classify emotions in real-world varied environments.

4.3 Face Detection

In this real-time emotion recognition project, face detection is of great importance because it identifies faces from live video frames captured by the webcam. Isolating faces from each frame ensures that this model only deals with parts that matter in each image frame and reduces unnecessary computation. Extracted detected faces are then converted to grayscale, resized to 100 by 100, and sent to the CNN model trained earlier.

This face detection system doesn't just improve the precision of detection by isolating the face within images but it also adds to model's efficiency in real time—a prime feature that this needs for interactive applications. Once a face is detected, a bounding box is drawn around it and the predicted emotion label is written above the box. The face detection, combined with the real-time feedback, ensures that the system is responsive and user-friendly for practical deployments.

4.4 Evaluation Metrics

Both the training and validation accuracy have been analyzed, which is the usual metric to consider in classification problems. From accuracy on the training data set, we can consider the ability of the model to learn, whereas its validation accuracy can be a kind of indication about whether it generalizes well or not. The best validation accuracy obtained was 78.32%, in which the operating result on new data with different facial expressions was reliable.

In such circumstances, cross-entropy loss is taken as the measure of optimization of the model during training. Cross-entropy loss is appropriate for problems in categorical classification and enables the quantification of the mismatch between the predicted probabilities and the actual labels. The ability to make more accurate predictions is reflected in low cross-entropy loss; therefore, model optimization is aligned during training.

Given the nature of the application as real-time, latency is thus another metric that is very informal but very important. The architecture is optimized for low latency with video frame processing being done as quickly as possible to enable instant feedback. Early stopping has been implemented in the model, with the model watching its validation loss so as not to overfit and save the best checkpoint for deployment.

5. IMPLEMENTATION

5.1 Preprocessing and Data Loading

Standardized Input for Uniformity: The preparation and data loading portion of the implementation involves standardizing inputs to the CNN. All images are resized to 100x100 pixels before conversion into grayscale, ensuring that the model receives uniform input dimensions along with color channels, simplifying computational requirements while retaining facial features necessary for emotion classification.

Various data augmentation techniques applied on the training set help in enhancing generalization to the model, like rescaling, shearing, zooming, and horizontal flipping; it simulates diverse conditions of real life and makes the model resistant to variations in terms of lighting, angle, and orientation of faces. Images also maintained in batches where each batch consists of a group of images for optimal memory usage and faster training.

5.2 Model Architectures

Now, multiple configurations such as layer counts, hidden layers, filter number, and filter size needed to be tested and have the proper architecture of CNN selected for testing. After multiple setups, the configuration with three hidden layers containing filters that were gradually increasing in complexity from 32 to 192 in each hidden layer was found optimal. This design achieved an appropriate balance between both accuracy and generalization since it presented good training accuracy as well as validation accuracy.

TABLE II. No. of Layers

No. of Hidden Layers	Accuracy	Val Accuracy
1	0.7695	0.7759
2	0.7922	0.7816
3	0.8172	0.7832
4	0.8202	0.7779
5	0.8291	0.7693

As extracted from Table II, the best possible performance with a three-layer configuration used a training accuracy of 81.72%, and validation accuracy of 78.32%. Increased beyond three layers did not improve performance significantly and sometimes overfit, where the model performed great on the training data but poor in unseen data. This choice reflects a structured and efficient approach, focusing on layers that help in extracting the essential spatial and hierarchical features.

TABLE III. No of Filters

No of Filters for each Layer	Accuracy	Val Accuracy
16, 32, 64, 96	0.7772	0.7307
32, 64, 128, 192	0.8172	0.7832
64, 128, 256, 384	0.8143	0.7410

This is as shown in Table III where from 32 to 192 filters are used that progressively increased across layers were helpful. This forced the model to be deeper in knowledge with the increase in depth. For instance, basic features concerning the face might be captured by an initial layer with 32 filters whilst

more subtle muscle movements concerned to recognize emotional input might be detected by the deeper layers and 192 filters. This configuration provided a stable validation accuracy of 78.32% and reflected on how the model was strong in its ability to handle different input while avoiding over-complexity.

Table IV. Filter Sizes

Filter Sizes	Accuracy	Val Accuracy
2x2	0.7938	0.7822
3x3	0.8172	0.7832
5x5	0.8150	0.7909

As shown in Table IV, the best filter size selected for the experiments conducted in all configurations was 3x3. The designed filter was able to capture sufficient details in images without letting the associated computational cost get out of hand, leading to 81.72% accuracy and 78.32% validation accuracy as opposed to accuracy of 2x2 and 5x5 filters respectively. This exactly maintained a balance between extracting detail and processing efficiency with the requirements of the real-time applications that focus on accuracy and speed.

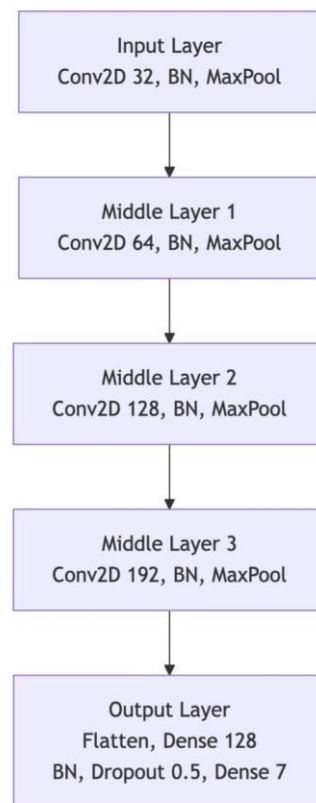


Figure 2: Model Architecture Deployed

Figure 2 shows the final model architecture that was deployed after running all the tests and trials.

5.3 Performance Training Process

Optimization and Regularization: The training of the model is owed to a set of optimizations and regularization to prevent overfitting and stable convergence. Adam's Optimizer was used in compilation along with dropout regularization rate of

50%, thus preventing the dependency of the model on one specific neuron while showing improvement in generalization on the validation set.

Checkpointing and Monitoring: In training, checkpointing monitored validation loss, saving the best model only. This keeps the model state when validation loss is minimized and avoids the scenario of model degradation especially when training for a long period in terms of epochs. The training is performed over 25 epochs with real-time validation to make sure the parameters are getting correctly adjusted for the model to ensure accuracy and generalization.

Real-Time Inference

Deployment in a Flask Web Application: The implemented model is deployed in a real-time emotion recognizer in a Flask web application. This configuration enables the system to process video feeds, detect faces, and classify emotions from live-streamed frames. Face detection functionality is achieved using the OpenCV library. During user access to the application, a video stream from a webcam will be shown, and the model will continuously process each frame that will do face detection and emotion classification on the detected faces.

Emotion Classification with Immediate Feedback: Within real time, the application detects multiple faces in every frame with bounding boxes around detected faces. For each detected face, the model predicts the emotion, which is then shown up as a label over the bounding box as shown in Figure 3. This interface proved very useful for analyzing group dynamics or personal interactions, which qualifies it for practical applications such as virtual meetings or assessment in person.

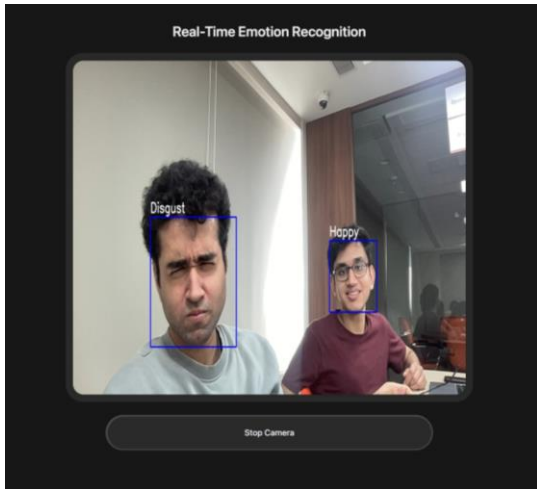


Figure 3: Real Time Implementation

6. RESULTS AND DISCUSSION

6.1 Performance Analysis

To evaluate the model's efficiency across demographics, we further tested its performance by emotion class, by age and sex across epoch settings of 15, 25, and 50 epochs. From this evaluation, accuracy maximization and minimization of overfitting were observed at an epoch setting of 25 epochs.

Table V. Neutral Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.4	0.8	0.8

Age Range (years)	Sex	No. of Epochs		
		15	25	50
25 - 39	Male	0.8	1.0	0.8
	Female	0.8	0.8	0.6
40 - 54	Male	1.0	1.0	0.8
	Female	0.8	1.0	0.6
55 - 70	Male	0.8	1.0	0.8
	Female	0.8	0.8	1.0

Table 5 illustrates that, for all but older subjects, the male group achieved accuracy close to 100% by the 25th epoch, and even increased to or stabilized at 50 epochs in some cases. Young males and females (10–24 years) improved dramatically from 15 to 25 epochs to around 0.8 accuracy, then stabilized or decreased at 50 epochs. Older age groups, 40–54 years, had the highest accuracies, meaning that the ability of the model to recognize neutral expressions may be more effective among adults than among the younger or the older age groups.

Table VI. Happy Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.8	1.0	1.0
	Female	0.8	1.0	0.8
25 - 39	Male	1.0	1.0	1.0
	Female	0.8	0.8	1.0
40 - 54	Male	1.0	1.0	1.0
	Female	0.8	1.0	1.0
55 - 70	Male	1.0	1.0	1.0
	Female	0.6	0.8	0.8

As seen in Table 6, Happiness, as an expressive emotion, achieved strong model accuracy across all age groups, with most categories stabilizing at or near 1.0 accuracy by 25 epochs. This finding supports the model's strength in detecting universally recognizable expressions like happiness, which tend to be visually consistent across age and gender.

TABLE VII. Sad Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.6	0.8	0.8
	Female	0.6	0.6	0.8
25 - 39	Male	1.0	1.0	0.8
	Female	0.8	0.8	0.8

Age Range (years)	Sex	No. of Epochs		
		15	25	50
40 - 54	Male	1.0	1.0	1.0
	Female	0.8	1.0	1.0
55 - 70	Male	0.6	0.8	0.6
	Female	0.6	0.6	0.6

As displayed in Table 7, the results for accuracy were age dependent, and with 25 epochs, only moderate levels of accuracy appeared for the younger subjects at 0.8. All epochs for adults falling into the category of between 25 and 54 years resulted in high accuracy, which implies that sadness was often recognized by the model, especially in middle-aged individuals. For older adults aged between 55 and 70, accuracy was varied, thus displaying some limitations of the model when the expressions given represent sadness for this age category.

TABLE VIII. Fear Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.4	0.8	0.6
	Female	0.4	0.6	0.4
25 - 39	Male	0.6	0.6	0.4
	Female	0.0	0.6	0/6
40 - 54	Male	0.4	0.8	0.6
	Female	0.2	0.4	0.2
55 - 70	Male	0.2	0.8	0.6
	Female	0.0	0.6	0.6

As seen in Table 8, Fear detection showed marked variability across demographics, especially among female subjects, with lower accuracies observed in younger females at all epoch settings. Although male accuracies were generally higher than females in detecting fear, accuracy remained below 1.0 across the board, highlighting a potential challenge in consistently recognizing fear expressions.

TABLE IX. Disgust Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.4	0.6	0.4
	Female	0.0	0.4	0.2
25 - 39	Male	0.6	0.6	0.8
	Female	0.2	0.4	0.4
40 - 54	Male	0.2	0.8	0.6
	Female	0.4	0.6	0.6
55 - 70	Male	0.6	0.6	0.6

Age Range (years)	Sex	No. of Epochs		
		15	25	50
	Female	0.4	0.6	0.2

It also varies with age and sex, as can be seen from Table 9. For middle-aged males at 40–54 years, the accuracy rises very high after 25 epochs, and the improvement was moderate for female subjects. These patterns indicate that the model finds it a bit challenging to detect disgust across different demographics.

TABLE X. Angry Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.2	0.4	0.4
	Female	0.0	0.4	0.2
25 - 39	Male	0.6	0.8	0.8
	Female	0.4	0.4	0.6
40 - 54	Male	0.8	1.0	1.0
	Female	0.8	1.0	1.0
55 - 70	Male	1.0	1.0	0.8
	Female	0.4	0.6	0.6

As it is evident from Table 10, significantly, the accuracy of Anger increased for older males, reaching 1.0 at 25 epochs for males in the age group 40 to 70; Females showed improvement with epochs but at slightly lower accuracies. Repeated recognition of anger in older males suggests that there may be a role of demographics which influence model performance like visual difference in anger expressions between the age groups.

TABLE XI. Surprise Emotion

Age Range (years)	Sex	No. of Epochs		
		15	25	50
10 - 24	Male	0.6	0.6	0.6
	Female	0.2	0.6	0.6
25 - 39	Male	0.8	1.0	1.0
	Female	0.8	0.8	0.8
40 - 54	Male	0.6	0.8	0.6
	Female	0.4	0.4	0.6
55 - 70	Male	0.6	0.6	0.6
	Female	0.6	0.6	0.4

As shown by Table 11, Balanced accuracy in detecting surprise was achieved over all demographics with the architecture showing higher accuracies predominantly reached at 25 epochs. That, again, would not indicate overfitting since the architectures managed to capture some relevant visual cues for

the notion of surprise adequately well in order to provide stability in prediction.

(Figure 4) Overall, the model works best at 25 epochs for optimal balance between accuracy and generalization across emotion classes and demographics. Epoch settings of more than 25 tend to show overfitting symptoms, particularly for a few emotions. This argument strengthens the basis for real-time applications to opt for the 25-epoch model.

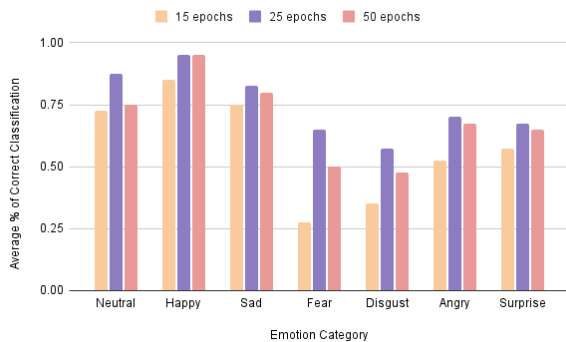


Figure 4: Comparative Analysis based on Epochs

6.2 Comparison with Existing Models

Comparison with popular FER models stresses the efficiency of this model with a specific dataset and architecture. Most models like VGGNet and ResNet that are extensively used for emotion detection are very accurate for tasks related to FER; however, these models are computationally intensive, usually requiring large datasets and much greater processing power in order to make sure that the performance is real-time.

Performance with Small Datasets: Compared to large scale datasets used by sophisticated models, for instance, AffectNet with 1 million images, our model reaches the levels of accuracy with much fewer images. This demonstrates the viability of our CNN design toward efficiency and its deployment in real-time. It thus is pretty well suited for applications with limited data and limited computational resources.

Most FER models, especially deep networks like ResNet, are very computationally expensive that limit their real-time applicability on standard hardware without the optimization. In contrast, the CNN model used here, with optimized parameters or 3 hidden layers and filter sizes, etc., shows practical implications through significant processing time and inference time, which provides real-time results applicable in applications such as live video analysis.

Emotion-Specific Accuracy Comparison: As more complex architectures sometimes break down when dealing with the visually subtle emotions of fear and disgust, especially in other demographics, such findings have already appeared in previous studies. These trends are set within this model as well, with similar trends set for error-prone consistency in accurately classifying fear and disgust. However, the model matched the performance of most other CNN-based systems in FER with universally recognizable emotions like happiness and anger.

Feasibility of Deployment: As the architecture of this model is pretty simple, it can be deployed easily on systems with a moderate grade of specification as opposed to big models that require much more computing capacity. This leads to possible applications in varied education, healthcare, and social research fields where real-time facial emotion recognition can add value by providing meaningful insights without high-end hardware setups.

7. FUTURE WORK AND LIMITATIONS

7.1 Challenges with RAF-DB

Though the RAF-DB dataset is highly popular and well-curated for facial emotion recognition, it involves many challenges that have affected model performance in this study. This inherent characteristic causes class imbalance within the dataset which results in fewer samples for emotions such as "disgust" and "fear" as compared to more universally recognizable emotions such as "happiness" and "neutral." This results in skewed learning patterns wherein the model overfits the majority classes, thus lowering predictive power for minority classes. Data augmentation techniques that were somewhat effective in alleviating this further can be built upon to carry forward the advanced methods like Synthetic Minority Over-sampling or even to explore other data generation strategies to further improve the class representation.

Another weakness of RAF-DB is poor demographic diversity in age, ethnicity, and cultural backgrounds, which might affect the ability of the model to be generalized in real application fields. Emotional expression varies with cultures and different demographics; therefore, any models based on RAF-DB would not represent all the nuances of such variations. Enlarging the dataset to include a broader demographic representation or merging it with another dataset, such as AffectNet or FER+, would have strong implications for improved model robustness and application to vastly diverse user bases.

7.2 Model Expansion

There is lots of scope for augmenting this model for the incorporation of additional functionality and better effectiveness. Future work might include adding extra network architectures, such as attention mechanisms, which have been demonstrated to produce effective good results in sharpening the attention in regions of the face that may contain the more important emotional cues. The system then will become more efficient at handling complex or even overlapping facial expressions.

Moreover, combining existing CNN with other architectures may improve the accuracy of models since different types and models from different categories will contribute toward the right model. Finally, the idea of the multi-modal data fusion for facial expressions might be adding vocal intonations or body language cues towards revealing a more holistic expression of emotional states. Such an approach for a multi-modal system would be really useful in real-time demanding applications where high accuracy is required, and it may be applied to human-computer interaction, mental health assessments, and even immersive virtual environments.

7.3 Privacy and Ethical Concerns

Real-time applications of facial emotion recognition systems also bring forth serious privacy and ethical concerns. The gathering and analysis of facial data demand strict measures to prevent their malicious use or unauthorized access to the data. In this work, we have paid attention to optimization of the model's performance; however, real-world deployment of such models would need compliance with data privacy regulation, such as GDPR. Data should be anonymized, encrypted, and stored in safe practices from privacy risks.

Furthermore, the question of ethics in facial emotion recognition, about bias on grounds of race, age, or gender, should be handled. Models trained on biased datasets risk propagating this in real-world applications and are likely to raise discriminatory outcomes. Further work may involve using fairness constraints or bias mitigation techniques during

training of models. Furthermore, integration of strong evaluation metrics might ensure that the fairness of model decisions is preserved across groups of demographics. Communication in relation to the limitations and ethical considerations of the model in its deployment also spells indispensable capabilities towards the responsible use of AI.

8. CONCLUSION

This study presents a real-time emotion recognition model using a CNN trained on the RAF-DB dataset, achieving a validation accuracy of 78.32% with an optimized configuration of three hidden layers and a 3x3 filter size. The 3x3 filter was selected to capture essential facial details without excessive computational cost, balancing detail extraction and processing efficiency. Data augmentation techniques—rescaling, shearing, zooming, and horizontal flipping—were specifically chosen to improve generalization, with horizontal flipping increasing validation accuracy to 78.32%, as it introduced facial symmetry variations beneficial for diverse emotional expressions.

The model performed well in detecting happiness and neutrality, reaching up to 100% accuracy in certain age and gender groups by 25 epochs, while accuracy for fear and disgust remained lower, particularly in female and younger demographics. Despite these challenges, our CNN-based model achieved comparable accuracy to larger models like VGGNet and ResNet, while requiring fewer computational resources and a smaller dataset, making it suitable for deployment in limited-resource environments.

To further improve performance, future work could address class imbalances in the RAF-DB dataset, possibly incorporating fairness constraints and multi-modal data sources for enhanced demographic representation. This study demonstrates that careful design of model architecture and augmentation can effectively balance real-time efficiency with accurate emotion recognition across diverse applications, while also highlighting the importance of privacy and ethical safeguards for deployment.

9. REFERENCES

- [1] D. Duncan, G. Shine, and C. English, "Facial emotion recognition in real time," Stanford University, 2020.. Available: duncand@stanford.edu
- [2] W. Mellouka and W. Handouzia, "Facial emotion recognition using deep learning: Review and insights," in *Proc. 2nd Int. Workshop Future Internet Everything (FIoE)*, Leuven, Belgium, 2020, pp. 1-8.
- [3] S. Kaur and N. Kulkarni, "Emotion recognition – A review," CSE, MIT SOE, MIT ADT University, Pune, India, 2020.
- [4] J. A. Ballesteros, G. M. Ramírez, F. Moreira, A. Solano, and C. A. Pelaez, "Facial emotion recognition through artificial intelligence," University of Cauca, Colombia, 2020.
- [5] S. M. Saleem Abdullah, S. Y. Ameen, M. A. M.sadeeq, and S. R. M. Zeebaree, "Multimodal emotion recognition using deep learning," 2020.
- [6] S. J. Oh, J. Y. Lee, and D. K. Kim, "The design of CNN architectures for optimal six basic emotion classification using multiple physiological signals," 2020.
- [7] P. Vanneste, J. Oramas, T. Verelst, T. Tuytelaars, A. Raes, F. Depaepe, and W. Van den Noortgate, "Computer vision and human behaviour, emotion and cognition detection: A use case on student engagement," 2020.
- [8] H. I. James, J. J. Anto Arnold, J. M. Masilla Ruban, M. Tamilarasan, and R. Saranya, "Emotion-based music recommendation system," Valliammai Engineering College, Chennai, Tamil Nadu, 2020.
- [9] E. M. Onyema, P. K. Shukla, S. Dalal, M. N. Mathur, M. Zakariah, and B. Tiwari, "Enhancement of patient facial recognition through deep learning algorithm: ConvNet," 2020.
- [10] A. Khopkar and A. A. Saxena, "Facial expression recognition using CNN with Keras," DES's Navinchandra Mehta Institute of Technology and Development & Research Scholar Pacific University, Udaipur, India, 2020.
- [11] A. Agrawal and N. Mittal, "Using CNN for facial expression recognition: A study of the effects of kernel size and number of filters on accuracy," 2020.
- [12] I.-K. Choi, H.-E. Ahn, and J. Yoo, "Facial expression classification using deep convolutional neural network," 2020.
- [13] Mollahosseini, A., Hasani, B., & Mahoor, M. H. (2017). AffectNet: A database for facial expression, valence, and arousal computing in the wild. *IEEE Transactions on Affective Computing*, 10(1), 18-31.
- [14] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT press.
- [15] Viola, P., & Jones, M. (2001). Rapid object detection using a boosted cascade of simple features. *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 511-518.
- [16] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 770-778).
- [17] Simonyan, K., & Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.