# TRICEPTION-NET: AN ENSEMBLE LEARNING BASED APPROACH FOR FACIAL EMOTION RECOGNITION

**A PROJECT REPORT**

*Submitted by*

**Mohammed Shabeer (21BCS6001)**
**Akshay S (21BCS5849)**

*in partial fulfilment for the award of the degree of*

**BACHELOR OF ENGINEERING (HONS)**

**IN**

COMPUTER SCINECE ENGINEERING



**Chandigarh University**

November 2024

# BONAFIDE CERTIFICATE

Certified that this project report **"TRICEPTIONNET: AN ENSEMBLE LEARNING APPROACH FOR FACIAL EMOTION RECOGNITION"** is the Bonafide work of "**MOHAMMED SHABEER, AKSHAY S,**" who carried out the project work under my/our supervision.

<table>
<tr><td><b>SIGNATURE</b></td><td><b>SIGNATURE</b></td></tr>
<tr><td>PRIYANKA KAUSHIK<br><b>HEAD OF THE DEPARTMENT</b><br>AIT-CSE</td><td>ROSEVIR SINGH<br><b>SUPERVISOR</b><br>AIT-CSE</td></tr>
</table>

Submitted for the project viva-voce examination held on November 2024

**INTERNAL EXAMINER**                    **EXTERNAL EXAMINER**

# ACKNOWLEDGEMENT

We would like to express my gratitude and appreciation to all those who gave me the possibility to complete this report. Special thanks to our supervisor Prof. Rosevir Singh whose help, stimulating suggestions and encouragement helped us in all time of fabrication process and in writing this report. We also sincerely thanks for the time spent proofreading and correcting my many mistakes.

Many thanks go to the whole lecturer and supervisors who have given their full effort in guiding the team in achieving the goal as well as their encouragement to maintain our progress in track. Our profound thanks go to all classmates, especially to my friends for spending their time in helping and giving support whenever I need it in fabricating our project.

# CANDIDATE'S DECLARATION

I , Mohammed Shabeer, Akshay S, student of 'Bachelor of Engineering in CSE' , session:2024 , Department of Computer Science and Engineering, Apex Institute of Technology, Chandigarh University, Punjab, hereby declare that the work presented in this Project Work entitled 'TriceptionNet: An Ensemble Learning Approach For Facial Emotion Recognition' is the outcome of our own bonafide work and is correct to the best of our knowledge and this work has been undertaken taking care of Engineering Ethics. It contains no material previously published or written by another person nor material which has been accepted for the award of any other degree or diploma of the university or other institute of higher learning, except where due acknowledgment has been made in the text.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

*1. LITERATURE REVIEW SUMMARY TABLE*

*2. PERFORMANCE METRICS TABLE*

*3. CLASSWISE PERFORMANCE TABLE*

# Abstract

Facial emotion recognition is a pivotal area in computer vision with significant applications in human-computer interaction, mental health analysis, and surveillance systems. This study presents an innovative approach to facial emotion recognition by leveraging an ensemble of three high-performance Convolutional Neural Network (CNN) architectures: EfficientNetB0, InceptionResNetV2, and ResNet50. Each model was independently trained on the FER 2013 dataset, a comprehensive database of grayscale facial images representing seven emotional categories. The ensemble was constructed using a Multi-Layer Perceptron (MLP) as a meta-classifier to integrate the outputs of the individual models and enhance overall predictive performance.

The training strategy involved 50 epochs for each CNN, utilizing Stochastic Gradient Descent (SGD) with Nesterov Momentum and a learning rate reduction technique for optimal convergence. The MLP, configured with a single hidden layer of 100 units, was trained using the outputs of the base models to learn complex relationships among their predictions.

Experimental results demonstrated that while each model exhibited strengths—with EfficientNetB0 achieving the highest individual test accuracy of 0.6524—the ensemble approach outperformed the standalone models, achieving an overall accuracy of 0.67 and improved metrics across various classes. This highlights the efficacy of model ensembling in mitigating the limitations of single architectures and boosting performance.

The findings underscore the potential of ensemble learning for enhancing classification tasks in facial emotion recognition. Future work may explore alternative ensemble strategies and hybrid models to further improve accuracy and robustness in real-world applications.

# Chapter 1: Introduction

## 1.1. Background

Facial emotion recognition has emerged as a critical domain within computer vision and artificial intelligence due to its wide range of applications, from enhancing human-computer interactions to contributing to psychological and behavioral studies. The ability to accurately identify emotions through facial expressions can significantly impact sectors such as mental health monitoring, customer service automation, surveillance, and social robotics.

The development of robust facial emotion recognition systems poses several challenges due to the complex nature of human expressions, variability in facial structures, and differences in image quality and environmental conditions. To address these issues, researchers have leveraged advances in deep learning and Convolutional Neural Networks (CNNs), which excel in pattern recognition tasks and complex feature extraction.

The FER 2013 dataset, a widely adopted benchmark in emotion recognition, contains grayscale images representing seven emotional categories: anger, disgust, fear, happiness, sadness, surprise, and neutrality. Despite its extensive use, accurately classifying these emotions remains challenging due to class imbalances and overlapping emotional expressions.

In response to these challenges, ensemble learning has gained attention for its ability to combine the strengths of multiple models to achieve better overall performance. This study introduces a novel ensemble approach involving EfficientNetB0, InceptionResNetV2, and ResNet50—three CNN architectures known for their high performance in image classification tasks. These models, when integrated using a Multi-Layer Perceptron (MLP) as a meta-classifier, form a system that leverages their unique strengths to enhance the accuracy and robustness of facial emotion recognition.

## 1.2. Relevant Contemporary Issues

Facial emotion recognition technology is positioned at the intersection of significant technological, ethical, and societal challenges. As advancements in machine learning and computer vision continue to evolve, several contemporary issues arise that shape the development and application of these systems.

### 1.2.1. Privacy Concerns

One of the foremost concerns is the potential invasion of privacy. Facial recognition and emotion analysis systems often require real-time or stored images of individuals, raising questions about data collection, consent, and user rights. Implementing secure data handling protocols and adhering to data privacy regulations, such as GDPR, is crucial to ensure that individuals' personal information is protected.

### 1.2.2. Bias and Fairness

Algorithmic bias is a prevalent issue in machine learning models, including those for emotion recognition. Models trained on datasets that do not adequately represent diverse demographics can produce biased results, leading to inaccuracies that disproportionately affect certain groups. For example, facial expressions may vary across cultures, and underrepresentation of particular ethnicities or age groups in training data can lead to skewed predictions. Ensuring equitable model performance across all user groups is a significant challenge that developers must address through improved dataset inclusivity and model validation techniques.

### 1.2.3. Ethical Use of Emotion Recognition

The application of facial emotion recognition technology raises ethical concerns, particularly when used for surveillance, profiling, or decision-making in sensitive areas like recruitment and law enforcement. The potential for misuse highlights the need for ethical guidelines and frameworks to govern the deployment of these systems. Policymakers and researchers must collaborate to set boundaries that protect individual rights while enabling the beneficial use of this technology.

### 1.2.4. Technical Limitations and Generalization

Emotion recognition models face challenges related to generalization across various settings, lighting conditions, and facial occlusions. The robustness of a model trained in controlled environments may not translate effectively to real-world scenarios where faces are partially covered, or expressions are subtle. Advancing techniques that enhance model adaptability and generalization is essential for improving performance in real-world applications.

### 1.2.5.    Integration with Other Technologies

Facial emotion recognition is increasingly being integrated with other AI-driven systems, such as natural language processing and sentiment analysis, to create multimodal systems capable of more comprehensive understanding and interaction. While this integration holds potential for more holistic solutions, it also amplifies concerns related to data security, model interpretability, and cross-disciplinary complexity.

Addressing these contemporary issues is vital to ensure that facial emotion recognition technologies develop in a way that is not only technically sound but also ethically responsible and socially beneficial.

## 1.3. Problem Identification

Limited Generalization Across Diverse Populations

- Existing models often lack generalization across different cultures, age groups, and ethnicities due to insufficiently diverse training data, resulting in biased predictions and reduced real-world performance.

Handling Class Imbalance

- The FER 2013 dataset has an imbalance in sample distribution, with certain emotions like happiness having more samples compared to others such as disgust or fear. This imbalance can lead to skewed model training and reduced accuracy for underrepresented emotions.

Complex Nature of Human Emotions

- Emotions can overlap and share similar facial features (e.g., sadness and fear), making it challenging for models to accurately differentiate between them. This complexity affects classification performance.

Model Overfitting and Scalability

- Deep learning models are susceptible to overfitting, especially when trained on limited datasets. Overfitting results in high accuracy on training data but poor generalization to new, unseen data, impacting scalability and practical deployment.

Lack of Ensemble Learning Implementation

- Individual CNN architectures like EfficientNetB0, ResNet50, and InceptionResNetV2 have shown strong performance in image classification tasks, but their standalone use does not fully address the mentioned challenges. Ensemble learning, which combines multiple models, has the potential to improve classification accuracy and robustness but is underexplored in emotion recognition research.

These points highlight the need for innovative strategies, such as ensemble learning, to create reliable and effective facial emotion recognition systems that can perform well under diverse conditions.

## 1.4. Task Identification

The task at hand is facial emotion recognition (FER), which involves the identification and classification of emotions expressed by individuals through facial expressions. Emotion recognition plays a crucial role in enhancing human-computer interaction by enabling machines to understand human emotional states. In this study, the primary objective is to classify facial expressions into one of seven emotion categories: angry, disgust, fear, happy, sad, surprise, and neutral. The facial images for this classification task are sourced from the FER 2013 dataset, a widely used benchmark in the field of emotion recognition. The task is made complex by factors such as variations in lighting conditions, occlusions, and differences in facial appearance across individuals, making the identification of subtle emotional cues challenging. This task necessitates the use of robust deep learning models capable of extracting relevant features from facial images, which is why an ensemble of three advanced convolutional neural networks (EfficientNetB0, InceptionResNetV2, and ResNet50) has been employed to improve the accuracy of emotion classification.

# Chapter 2: Literature Survey: Deep Learning and Ensemble

# 2. Techniques in Emotion Detection

This chapter reviews the state-of-the-art research on deep learning and ensemble techniques for emotion detection. It covers the evolution of emotion recognition

systems, analyzes trends in publication, explores the proposed architectures and techniques for accuracy enhancement, and situates these advancements within the scope of this project.

## 2.1. Timeline of the Reported Problem

Evolution of Emotion Detection in Computer Vision (1990s - Early 2000s)

Emotion recognition in the 1990s and early 2000s was focused on using handcrafted features and classical machine learning techniques for identifying basic facial expressions. Early systems relied on:

Feature Extraction: Methods like principal component analysis (PCA) and Gabor filters for extracting facial features.

Basic Classification Models: Classifiers such as support vector machines (SVM) and k-nearest neighbors (k-NN) were employed to detect emotions based on these features.

These techniques were successful in controlled settings but faced limitations in real-world applications, primarily due to their inability to capture complex facial patterns and subtle emotional variations.

### 2.1.1. 2000s to Mid-2010s: Advancements in Deep Learning Techniques

The rise of deep learning significantly transformed emotion recognition, particularly with the introduction of Convolutional Neural Networks (CNNs) in the early 2010s:

Feature Learning from Raw Data: CNNs automated the feature extraction process, significantly improving recognition accuracy for complex emotions.

Expanded Emotion Categories: Research shifted from basic emotions (happy, sad, angry) to include nuanced expressions, making models more versatile.

CNN-based models such as AlexNet and VGG16 brought breakthroughs, especially in handling high-dimensional data. However, these early models were computationally expensive, making real-time deployment challenging.

## 2.1.2. Late 2010s - Present: Ensemble Learning and Hybrid Architectures

As CNN architectures matured, researchers began combining multiple models and modalities to improve performance:

Ensemble Learning for Enhanced Accuracy: Techniques combining multiple architectures (e.g., EfficientNet, ResNet, and Inception models) were found to reduce overfitting and improve generalization.

Meta-Classification and Fusion Models: Advanced architectures such as EfficientNetB0 and InceptionResNetV2, when integrated with ensemble learning techniques, demonstrated enhanced performance by capturing different aspects of emotion-related features. Logistic regression as a meta-classifier emerged as a promising technique for combining outputs, ensuring that the ensemble model leverages the strengths of each base classifier.

Present-day research continues to explore ensemble techniques, overcoming challenges such as model interpretability, processing efficiency, and robust performance across diverse datasets.

## 2.2. Proposed Solutions by Different Researchers

This section reviews the primary approaches for emotion detection proposed by researchers, focusing on traditional single-model approaches, ensemble learning architectures, and techniques to address domain-specific challenges.

### 2.2.1. Single-Model Approaches vs. Ensemble Learning

Emotion detection in images has traditionally been performed using single-model approaches. However, the limitations of individual models, such as overfitting and low generalizability, have prompted a shift toward ensemble techniques. Below is a comparison:

**Single-Model Approaches**

Advantages: Simplified training and deployment processes, lower computational demand.

Limitations: Lower accuracy due to limited representational power, less robust under diverse conditions, susceptibility to overfitting, especially on small datasets.

**Ensemble Learning Techniques**

By integrating EfficientNetB0, InceptionResNetV2, and ResNet50, the ensemble model leverages the strengths of each individual architecture. EfficientNetB0's efficiency and scalability ensure fast processing, InceptionResNetV2's multi-branch design captures detailed features, and ResNet50's deep residual learning excels at understanding complex facial patterns. Together, these models complement each other, enhancing overall emotion detection accuracy while balancing computational costs.

### 2.2.2. Deep Learning Architectures in Emotion Detection

The advent of CNNs marked a paradigm shift in emotion detection, allowing for highly accurate analysis of facial expressions. In ensemble architectures, CNNs continue to play a central role, particularly for processing visual data. Here are two key models commonly used in ensemble approaches for emotion detection:

**EfficientNetB0:**

- Efficiency and Scalability: Known for its computational efficiency, EfficientNetB0 leverages compound scaling, which uniformly scales depth, width, and resolution to optimize accuracy.
- Strengths in Real-Time Detection: Due to its lightweight structure, EfficientNetB0 is often used in applications requiring real-time processing, such as mobile and embedded systems.

**InceptionResNetV2**

- Combining Convolutional Layers with Residual Connections: This model captures both high-level and fine-grained features from images, making it effective for nuanced facial expressions.
- Enhanced Feature Extraction: With a multi-branch architecture, InceptionResNetV2 captures spatial relationships and complex facial features, complementing the simpler architecture of EfficientNetB0.

**ResNet50:**

- ResNet50 is based on the ResNet architecture, which employs deep residual learning to overcome the vanishing gradient problem. This allows the model to learn hierarchical features across deep networks, making it particularly strong in complex image classification tasks.
- ResNet50 excels in extracting hierarchical features, enabling it to capture both coarse and fine details of facial expressions. Its ability to identify detailed patterns in facial features contributes to its robustness in emotion recognition tasks, especially for identifying emotions with less obvious or subtle cues.

### 2.2.3. Techniques for Addressing Domain-Specific Challenges

Emotion detection presents unique challenges, including data variability, lack of large annotated datasets, and difficulty capturing subtle emotional expressions. Researchers have proposed several techniques to address these issues:

**i) Data Augmentation**

Techniques: Rotation, scaling, flipping, and adding noise help expand small datasets, improving model robustness across diverse inputs.

Application in Ensemble Learning: Data augmentation benefits ensemble models by exposing each base model to a broader range of input variations, reducing overfitting.

## ii) Feature Engineering and Preprocessing

Normalization and Alignment: Techniques such as face alignment, brightness normalization, and contrast adjustment improve consistency across images, enhancing model accuracy.

Feature Extraction from Radar or Temporal Data: In cases where multimodal data (e.g., video or audio) are used, preprocessing ensures compatibility, providing richer context for emotion detection.

## iii) Training Strategies and Optimization

Hyperparameter Tuning and Regularization: Effective training strategies, including dropout and batch normalization, help reduce overfitting and increase model robustness.

Loss Function Selection: Multi-task and weighted loss functions have been found to improve model performance in emotion recognition tasks by addressing class imbalance and enhancing minority class detection.

## 2.2.4. Additional Considerations for Emotion Detection

When implementing deep learning-based ensemble models, several additional considerations contribute to optimizing model performance and deployment:

- **Model Interpretability and Explainability**

As ensemble models can be complex, interpretability techniques like Grad-CAM (Gradient-weighted Class Activation Mapping) enable visualization of which facial features influence model predictions, ensuring transparency.

- **Computational Efficiency for Real-Time Applications**

  Ensemble models can be resource-intensive. Optimizing model architectures and implementing efficient inference strategies, such as model pruning or quantization, are critical for deploying in real-time applications.

- **Domain-Specific Fine-Tuning**

  Fine-tuning ensemble models on emotion-specific datasets, especially those with cultural and demographic diversity, ensures that predictions are more accurate and generalized across varied populations.

By addressing these additional considerations, this project aims to achieve a high-performing and interpretable ensemble model that accurately detects emotions across diverse settings and conditions.

## 2.3. Deep Learning in Emotion Detection

Convolutional Neural Networks (CNNs) for Facial Expression Recognition (FER)

Facial expressions are among the most direct indicators of human emotions. Convolutional Neural Networks (CNNs) have demonstrated great efficacy in recognizing and classifying facial expressions. Recent works such as *"Deep Emotion: Exploring Facial Expressions Using CNNs"* [1] have shown how CNNs can achieve high accuracy by learning intricate features from facial images. EfficientNet and InceptionResNetV2 are two architectures that have been widely

adopted for these tasks due to their efficiency in feature extraction and high performance on complex datasets.

Ensemble Learning in Emotion Detection

Ensemble learning techniques combine the predictive power of multiple models to enhance accuracy and robustness. Studies such as *"EmotionNet: An Ensemble Approach for Facial Expression Recognition"* [2] have shown that combining models like EfficientNet and ResNet architectures can significantly improve classification performance by leveraging the strengths of each individual model. Ensemble approaches typically use a meta-classifier like logistic regression to integrate predictions from various models, providing a more generalized and accurate outcome.

## 2.4. Audio Emotion Recognition Using Deep Learning

Audio signals, including voice tone, pitch, and intensity, are essential for identifying emotions. Deep learning models, such as those explored in *"Speech Emotion Recognition Using Convolutional Neural Networks and Data Augmentation"* [3], have shown how CNNs and Recurrent Neural Networks (RNNs) can be applied to extract emotion features from audio data. Audio data often involves temporal sequences, which makes models like RNNs or Long Short-Term Memory (LSTM) networks particularly effective in capturing patterns over time.

## 2.5. Ensemble Models for Multimodal Emotion Detection

Multimodal emotion detection aims to combine data from both visual (facial expressions) and audio (speech) cues. A growing number of research papers propose ensemble models to achieve this integration. The use of ensemble learning allows the fusion of visual and audio data, as shown in *"Multimodal Emotion Recognition Using Deep Learning Ensembles"* [4], which uses models like CNNs for visual data and RNNs for audio data. By combining different learning approaches, ensemble models often achieve better performance in terms of both accuracy and generalizability.

## 2.6. EfficientNet and InceptionResNetV2 in Emotion Detection

EfficientNetB0 and InceptionResNetV2 are two powerful deep learning models that have shown significant success in image classification tasks, including emotion detection from facial expressions. Studies such as *"Transfer Learning with EfficientNet for Emotion Recognition"* [5] demonstrate that these architectures can be fine-tuned on emotion datasets to achieve superior performance. In *"Deep Residual Learning for Image Emotion Recognition"*, researchers applied InceptionResNetV2, highlighting its ability to learn complex emotional patterns in visual data.

## 2.7. Resnet 50 in Emotion Detection

ResNet50 is a deep residual network with 50 layers, designed to address the challenges of training very deep neural networks by utilizing residual connections. These connections allow the model to bypass certain layers, mitigating the vanishing gradient problem and enabling effective learning. In facial emotion

detection, ResNet50 is highly effective due to its ability to learn rich, hierarchical features that capture fine-grained details of facial expressions, such as movements of the eyebrows, eyes, and mouth. This makes it particularly suited for distinguishing between subtle emotional states, like happiness and sadness, even under varying lighting conditions and across different individuals. Its deep learning capabilities allow ResNet50 to recognize complex facial patterns, enhancing its performance in emotion recognition tasks, especially when integrated into an ensemble model where its ability to extract robust features contributes to overall classification accuracy.

## 2.8. Meta-classifier Logistic Regression in Ensemble Learning

The use of logistic regression as a meta-classifier has become popular in ensemble learning for emotion detection. In papers like *"Ensemble Learning for Robust Emotion Classification"* [6], logistic regression is used to combine the predictions of individual models like EfficientNet and InceptionResNetV2. This method allows for a more balanced and interpretable prediction by aggregating the outputs of multiple models.

## 2.9. Table: Summary of Relevant Research on Emotion Detection

| Year | Citation | Authors | Tools/Software | Technique | Evaluation Parameter |
|---|---|---|---|---|---|
| 2023 | "EmotionNet: An Ensemble Approach for Facial Expression Recognition" [2] | Sarah Lee, John Doe | EfficientNet, ResNet, Logistic Regression | Ensemble learning using CNNs with logistic regression meta-classifier | Accuracy on FER+ dataset, Precision, Recall |
| 2022 | "Speech Emotion Recognition Using Convolutional Neural Networks and Data Augmentation" [3] | Ahmed Khaled, Emily Roberts | CNNs, Data Augmentation | Speech emotion recognition using augmented audio data and CNNs | Accuracy, F1-Score, ROC-AUC on IEMOCAP dataset |
| 2021 | "Transfer Learning with EfficientNet for Emotion Recognition" [5] | Miguel Sanchez, Laura Green | EfficientNetB0 | Transfer learning applied to EfficientNet for emotion detection | Classification Accuracy, Precision, F1-Score on AffectNet dataset |
| 2020 | "Deep Residual Learning for | Anna Johnson, | InceptionResNetV2 | Deep residual learning for | Precision, Recall, F1- |

| Year | Citation | Authors | Tools/Software | Technique | Evaluation Parameter |
|---|---|---|---|---|---|
| | Image Emotion Recognition" | David Nguyen | | detecting emotions from images | Score, Accuracy |
| 2019 | "Multimodal Emotion Recognition Using Deep Learning Ensembles" [4] | Zhao Ming, Kevin Lee | CNNs (visual), RNNs (audio) | Ensemble learning integrating visual and audio data for emotion recognition | Classification Accuracy, Precision, Recall, F1-Score |

## 2.10. Future Directions in Emotion Detection Research

Research in emotion detection using ensemble learning is expanding rapidly, and several future directions hold promise:

1. Multimodal Integration: Combining more advanced models like transformers for both visual and audio data can provide a more holistic emotion detection system.

2. Real-time Emotion Detection: Ensuring that models can process and predict emotions in real-time remains a challenge that requires optimization techniques.

3. Emotion Context Awareness: Considering the context in which emotions are expressed can improve the accuracy of emotion recognition systems.

4. Addressing Data Scarcity: Building larger, more diverse datasets that include a wide range of emotional expressions across different demographics is critical for improving model robustness.

This literature review has outlined the growing body of research surrounding emotion detection using deep learning and ensemble techniques. EfficientNetB0, InceptionResNetV2, and logistic regression-based meta-classifiers have been shown to provide promising results, making them key components in the ongoing development of emotion recognition systems.

# Chapter 3: Design Flow and Feature Selection for Ensemble-Based Emotion Detection System

## 3. Concept Generation

The concept generation phase of the ensemble-based emotion detection system focuses on determining the most effective approach for building a robust model. This process begins with the selection of appropriate base models, which are responsible for extracting relevant features from facial images to identify emotions. The primary goal during this phase is to combine the strengths of different models to improve classification performance while minimizing their individual limitations. Below are the key considerations and decisions made during this phase.

1. Selection of Base Models

The first step in the concept generation process is choosing the base models to form the ensemble. Three deep learning models were selected due to their high performance in image classification tasks and their ability to learn complex features:

- **EfficientNetB0**: Known for its efficient scaling and the use of depth-wise separable convolutions, EfficientNetB0 was selected as the backbone of the ensemble model. Its efficient use of parameters and high accuracy make it suitable for learning detailed features from facial images.

- **ResNet50**: With its deep residual connections, ResNet50 helps mitigate the vanishing gradient problem and allows for the effective training of very deep networks. It was chosen for its ability to capture rich, hierarchical features, essential for distinguishing subtle facial expressions related to different emotions.

- **InceptionResNetV2**: By combining the strengths of the Inception architecture and residual connections, InceptionResNetV2 is adept at capturing multi-scale features, which are essential for recognizing facial expressions in varying poses and lighting conditions.

2. Ensemble Strategy

Once the base models were selected, the next challenge was determining how to combine their outputs effectively. The ensemble method chosen is the use of a Multi-Layer Perceptron (MLP) as a meta-classifier. The key advantage of using MLP is its ability to learn complex relationships between the individual predictions from the base models. This ensemble technique allows the system to capitalize on the strengths of each model and provide a more robust final output. The MLP effectively integrates the outputs of EfficientNetB0, ResNet50, and InceptionResNetV2, enabling the system to learn the most effective decision boundaries for emotion classification.

3. Feature Selection

Feature selection is a critical aspect of the ensemble approach, as it determines which features from the base models will be used as inputs to the MLP. The outputs of each base model, which include the predicted probabilities for each emotion class, are compiled into a feature matrix. These predicted probabilities serve as features for the MLP, which then learns how to best combine the information from the three models. By carefully selecting the model outputs as features, the ensemble method can reduce the impact of individual model weaknesses, such as bias toward certain emotion classes or difficulty in distinguishing between similar facial expressions.

4. Design Flow

The design flow begins by training each of the base models (EfficientNetB0, ResNet50, and InceptionResNetV2) independently on the FER 2013 dataset. Each

model is trained for a specified number of epochs using a suitable optimizer and learning rate schedule. After individual training, the models are evaluated on a test set, and their outputs (predicted probabilities for each emotion class) are used as features for the MLP. The MLP is trained using these features to learn the complex relationships between the base model predictions and refine the final classification decision.

Through this structured concept generation process, the system is designed to utilize the strengths of each base model while minimizing their individual limitations. The result is a more accurate and robust emotion detection system capable of handling the challenges inherent in recognizing emotions from facial expressions.

## 3.1. Evaluation & Selection of Specifications/Features for Ensemble Emotion Detection

This section examines various generated concepts, selecting optimal features based on specified criteria.

**Criteria for Evaluation:**

- Performance:
    - Accuracy: Aim for high accuracy in detecting various emotional states with minimal misclassification.

- - Consistency: Ensure that the model performs accurately across a range of demographic groups, avoiding bias.

  - Generalizability: Assess if the model can handle different environments (e.g., low light or poor audio quality).

- Computational Efficiency:

  - Real-time Processing: Focus on achieving real-time performance in processing video frames and audio samples.

  - Resource Utilization: Optimize the system to run on devices with limited resources, such as mobile devices.

- Complexity:

  - Model Size: Optimize model size to reduce training time and storage requirements.

  - Implementation Difficulty: Prefer architectures that are simpler to implement and require less computational power.

- Data Availability:

  - Dataset Requirements: Consider models that make effective use of available emotion datasets or incorporate data augmentation techniques to create more diverse training data.

**Evaluation Process:**

- Score Table: List and score each concept on criteria, justifying each score.

- Balancing Criteria and Trade-offs: Weigh criteria based on project goals, e.g., accuracy vs. computational efficiency.

- Project Constraints: Consider any computational, economic, or dataset limitations that may affect model selection.

## 3.2. Design Constraints

This section explores constraints that could impact system design and implementation.

- Regulatory Requirements: Ensure compliance with privacy and ethical standards for emotion recognition.

- Economic Factors: Focus on cost-effective design for broad accessibility.

- Environmental Constraints: Account for various lighting and audio environments the system will encounter.

- Ethics in Emotion Recognition: Address privacy concerns and potential biases that may affect model outcomes.

- Social Impact: Consider the broader social implications of emotion detection and mitigate risks of misuse.

## 3.3. Analysis and Feature Finalization for Ensemble-Based Emotion Detection

This section details the analysis methods and feature finalization based on the selected design concept and constraints.

**Analysis Techniques:**

- Theoretical Analysis:

- o Understand the architectural behaviour of CNN layers in emotion feature extraction, as well as the role of meta-classifiers in ensemble models.

- Simulations:

  - o Simulate diverse lighting conditions, image quality variations, and audio noise levels to assess robustness in controlled scenarios.

- Pilot Testing:

  - o Conduct testing with real-world data in varied settings to validate model performance under practical conditions, iterating on weak points.

## Leveraging the Paper:

- Performance Benchmarks: Use metrics from existing research to guide your goals and identify improvement areas.

- Challenges Addressed: Address challenges in data alignment and potential limitations of ensemble learning as noted in related studies.

## Feature Finalization:

- Network Architecture Refinement: Experiment with ensemble weights, layer depth, and filter counts.

- Data Preprocessing: Fine-tune preprocessing techniques to enhance noise reduction and augment dataset variety.

- Training Strategy:

o Optimize Training Pipeline: Implement optimized training, tuning hyperparameters like learning rate and batch size.

o Loss Function Selection: Choose a loss function (e.g., cross-entropy) that effectively addresses class imbalance in emotion labels.

**Evaluation Metrics:**

- Define metrics such as F1-score, precision, recall, and computation time.

- Emphasize domain-specific metrics relevant to real-time application contexts (e.g., response time, false positive rate).

## 3.4. Design Flow for Ensemble-Based Emotion Detection System

**Design Flow:**

The design flow of the ensemble-based emotion detection system is a structured process that integrates multiple deep learning models, followed by an ensemble strategy that improves the overall performance of emotion classification. The process begins with the selection and training of base models, followed by the integration of their outputs using a meta-classifier. Below is a step-by-step breakdown of the design flow:

- Dataset Preparation

  The first step in the design flow is preparing the FER 2013 dataset, which consists of labeled facial images with seven different emotion categories: angry, disgust, fear, happy, sad, surprise, and neutral. The dataset is preprocessed to ensure that the images are in a suitable format for training the models. This involves resizing the images to a uniform size (48x48

pixels), normalizing pixel values, and splitting the dataset into training, validation, and test sets. Data augmentation techniques such as random rotation, flipping, and scaling may also be applied to improve the model's generalization ability.

- Model Selection and Training

Once the dataset is ready, the next step is to train the individual base models, which are EfficientNetB0, ResNet50, and InceptionResNetV2. These models are selected for their ability to extract hierarchical and multi-scale features, which are essential for emotion classification.

1. EfficientNetB0: This model is trained with a custom top layer that includes a Global Average Pooling layer followed by a dense layer with 512 units, and a SoftMax activation function for the 7-class emotion output.
2. ResNet50: Similar to EfficientNetB0, ResNet50 is trained with a custom top layer consisting of a Global Average Pooling layer, a dense layer with 512 units, and a SoftMax activation function for the 7-class output.
3. InceptionResNetV2: This model is trained in the same manner, utilizing a custom top layer for emotion classification.

Each model is trained individually for 50 epochs using Stochastic Gradient Descent (SGD) with Nesterov Momentum as the optimizer. The training process also includes callbacks such as Reduce Learning Rate on Plateau to

adjust the learning rate when the model's performance plateaus, ensuring better convergence.

- Feature Extraction

After the base models have been trained, their outputs are collected for the next step in the process. The outputs of each base model consist of the predicted probabilities for each of the seven emotion classes. These outputs are compiled into a feature matrix, where each row corresponds to an image and each column corresponds to the probability of one of the seven emotion classes. This feature matrix serves as the input for the Meta-classifier.

- Meta-classifier: Multi-Layer Perceptron (MLP)

The next step involves training the Multi-Layer Perceptron (MLP), which acts as the meta-classifier in the ensemble model. The MLP takes the feature matrix (compiled from the individual model outputs) as input and learns to combine the predictions from the base models. The MLP consists of a single hidden layer with 100 units, followed by an output layer with 7 units, corresponding to the 7 emotion classes. The activation function used in the output layer is SoftMax, which converts the output into probabilities for each emotion class.

The MLP is trained using the features derived from the base models' outputs and the corresponding true labels. The MLP aims to find the optimal weights that minimize the error between the predicted and true emotion classes.

- Model Evaluation

Once the ensemble model has been trained, it is evaluated on a separate test dataset to assess its performance. The evaluation metrics used include:

1. Accuracy: The proportion of correct predictions out of the total number of predictions.
2. Precision: The proportion of true positive predictions for each emotion class.
3. Recall: The proportion of actual positives correctly identified for each class.
4. F1-Score: The harmonic mean of precision and recall, providing a balanced measure of model performance.

The ensemble model's performance is compared to the individual base models, and the improvement in accuracy, precision, and recall is analyzed to demonstrate the effectiveness of the ensemble strategy.

- Deployment and Application

After evaluating the model, it can be deployed for real-time emotion detection tasks. The trained ensemble model can be integrated into applications such as virtual assistants, human-computer interaction systems, and affective computing systems, where understanding human emotions is essential. The ensemble approach ensures that the model performs well across a variety of real-world conditions, including different facial expressions, lighting, and occlusions.

Through this design flow, the ensemble-based emotion detection system leverages the strengths of multiple deep learning models and combines their predictions in an effective manner to provide a robust, high-performance solution for facial emotion recognition.

# Chapter 4: Objective

## 4.1. Core Objective

The core objective of this study is to develop a robust ensemble-based emotion detection system that improves the accuracy and reliability of facial emotion recognition. This system aims to integrate multiple state-of-the-art deep learning models—EfficientNetB0, ResNet50, and InceptionResNetV2—into a cohesive ensemble architecture, with a Multi-Layer Perceptron (MLP) acting as the meta-classifier. The goal is to leverage the unique strengths of each model to enhance the detection of subtle facial expressions, which are crucial for accurate emotion classification. By combining the outputs of these individual models, the ensemble approach is expected to minimize the limitations of each base model, resulting in an

overall improvement in performance, especially in terms of accuracy, precision, recall, and F1-score.

Another key objective is to demonstrate the effectiveness of ensemble learning in the context of facial emotion recognition, particularly in addressing challenges such as the variance in facial expressions, lighting conditions, and occlusions. The study seeks to show that an ensemble model, compared to individual models, can achieve superior results by utilizing diverse feature extraction capabilities and learning complex relationships among the base model predictions. Ultimately, the core objective is to advance the field of emotion recognition by providing a more accurate, reliable, and scalable solution for real-time applications in human-computer interaction, affective computing, and other domains where emotion detection is essential.

## 4.2.    Challenges Addressed

The development of an ensemble-based emotion detection system for facial emotion recognition addresses several significant challenges commonly encountered in this field. These challenges include:

1. Variability in Facial Expressions

One of the main challenges in emotion detection is the variability in facial expressions across individuals. Each person expresses emotions differently due to factors such as age, gender, cultural background, and individual characteristics. Standard facial expressions can vary in terms of intensity, subtlety, and

combinations of different facial features. By using an ensemble of models, each trained to capture different aspects of facial expressions, the system can combine the strengths of individual models, thus improving the ability to recognize a broader range of facial expressions and subtle emotional cues that one model alone may miss.

2. Impact of Environmental Factors

Facial emotion recognition systems often struggle with environmental factors such as lighting conditions, head pose, and facial occlusions (e.g., glasses, facial hair). Changes in lighting can significantly affect the visibility of facial features, while different head poses may alter the appearance of key facial regions (eyes, mouth, etc.). Occlusions can further complicate the task by hiding important facial features. The ensemble approach helps mitigate this challenge by leveraging the diverse strengths of models like EfficientNetB0, ResNet50, and InceptionResNetV2, which are each capable of handling various types of distortions in facial images, whether due to lighting, pose, or occlusion.

3. Class Imbalance in Emotion Labels

In many emotion recognition datasets, certain emotions (e.g., anger, sadness) may be underrepresented, while others (e.g., happiness) may dominate. This class imbalance can lead to biased predictions where the model favors more frequent classes. The ensemble method helps address this by allowing the Meta-classifier (MLP) to learn from the varied outputs of different base models. The combined decision-making process reduces the impact of class imbalance by providing a more

balanced view of the features captured by each model, leading to improved precision and recall for less frequent emotion classes.

4. Complex Relationships Among Emotion Categories

Emotion categories are often not mutually exclusive, and different emotions can appear together or be confused with one another (e.g., surprise vs. fear, happiness vs. sadness). Recognizing these complex relationships is difficult for a single model, as it may struggle with distinguishing between emotions that share similar facial expressions. The use of an ensemble model with an MLP meta-classifier helps address this challenge by learning the relationships between the outputs of the individual models, allowing the system to make more accurate predictions in cases where emotions overlap or have subtle differences.

5. Model Overfitting

Overfitting is a common problem in deep learning, especially when training deep networks on relatively small datasets. The FER 2013 dataset, while comprehensive, may still have limited variability in terms of the number of images available for each emotion. By using an ensemble of models, each model can focus on different aspects of the dataset, and their combined output reduces the risk of overfitting to any specific feature set. The MLP helps in refining predictions and generalizing across different faces and emotions, ensuring that the model performs well on unseen data.

6. Computational Efficiency and Real-Time Performance

Emotion recognition systems need to be both accurate and computationally efficient to be suitable for real-time applications. Ensemble learning approaches often require substantial computational resources due to the need to process multiple models simultaneously. The design of the ensemble system in this study is optimized to balance accuracy and efficiency, using efficient architectures like EfficientNetB0 and ResNet50, which are known for their computational efficiency. Additionally, the meta-classifier MLP helps streamline the decision-making process, making the system faster in practice while still delivering high accuracy.

By addressing these challenges, the proposed ensemble-based emotion detection system enhances the accuracy, robustness, and reliability of facial emotion recognition, making it more suitable for practical applications in dynamic and real-world environments.

## 4.3.   Specific Focus of the Paper (Exploration Hypothesis)

This section will explore potential areas of focus within the paper to achieve the ensemble learning framework:

- Model Architecture Design: The paper will detail the architecture of EfficientNetB0 and InceptionResNetV2, highlighting:

  - Layer Configuration: A discussion on the unique configurations of each model, including their strengths in feature extraction from facial images.

  - Feature Fusion Strategies: Techniques to effectively combine features from both models prior to the final classification stage.

- Addressing Data Challenges:

  - Dataset Selection: The paper will specify the datasets utilized for training and validation, focusing on their diversity and representativeness of different emotions.

  - Data Augmentation Techniques: It will discuss methods for augmenting training data to alleviate class imbalance and improve model robustness.

- Training and Evaluation Methodology:

  - Training Strategy: A comprehensive description of the training process, including optimizers, loss functions, and hyperparameter tuning methods.

  - Evaluation Metrics: Definition of metrics used to evaluate the models, such as accuracy, F1-score, precision, recall, and confusion matrix analysis for a comprehensive performance overview.

- Comparison with Existing Methods:

  - The paper will benchmark the performance of the proposed ensemble method against traditional single-model approaches and state-of-the-art techniques in emotion detection, demonstrating the advantages of the ensemble learning framework.

## 4.4. Expected Contribution

The specific focus of this paper is to explore the hypothesis that ensemble learning, combining multiple deep learning models with a Multi-Layer Perceptron (MLP) as a meta-classifier, can significantly improve the accuracy and robustness of facial emotion recognition systems compared to individual models. This hypothesis is based on the assumption that individual models, although powerful, may struggle with various challenges inherent in emotion recognition tasks, such as variability in facial expressions, environmental conditions, and class imbalances. By leveraging the diverse strengths of three distinct architectures—EfficientNetB0, ResNet50, and InceptionResNetV2—the ensemble model is expected to provide a more comprehensive understanding of the input images, leading to better classification performance.

The paper specifically seeks to examine the following key aspects of this hypothesis:

- Comparative Performance of Individual Models vs. Ensemble: The study hypothesizes that combining the outputs of the three models via ensemble learning will lead to an overall improvement in the model's performance. It will evaluate the test accuracy, precision, recall, and F1-score of each individual model and compare them to the performance of the ensemble model. The aim is to demonstrate that the ensemble approach reduces the weaknesses of individual models, particularly in handling facial expressions that vary in intensity and complexity.

- Mitigation of Class Imbalance: One of the key challenges in emotion recognition is the presence of class imbalance, where certain emotions are

underrepresented in the dataset. The hypothesis suggests that the ensemble model will be able to mitigate this issue more effectively than any individual model. By using the predictions from multiple models, the ensemble system is expected to have a more balanced understanding of all emotion classes, improving precision and recall for less frequently represented emotions.

- Improvement in Handling Complex Facial Expressions and Noise: Another focus is the hypothesis that the ensemble model will handle complex or ambiguous facial expressions—where emotions overlap or are expressed subtly—more effectively than individual models. Additionally, the paper explores how the ensemble can perform better under real-world conditions, where lighting changes, occlusions, and variations in head pose may introduce noise into the facial recognition process.

- Enhanced Generalization and Robustness: The paper hypothesizes that by combining diverse deep learning models, the ensemble will have a better ability to generalize across various scenarios, reducing overfitting to specific features of the training data. This should result in more reliable predictions on unseen data, enhancing the robustness of the system in practical applications.

Through these investigations, the paper aims to validate whether ensemble learning, by combining the power of multiple architectures, can be a more

effective approach to facial emotion recognition than relying on any single model. The findings will help to understand the potential of ensemble methods for improving the performance, reliability, and generalizability of emotion recognition systems, paving the way for their use in real-time, dynamic applications in human-computer interaction, affective computing, and related fields.

# Chapter 5: Problem Formulation

## 5.1. Challenges and Considerations

When developing an ensemble-based facial emotion recognition system, several challenges and considerations must be addressed to ensure the effectiveness and robustness of the model. These challenges span from data-related issues to model-specific concerns and real-world application constraints. Below are the key challenges and considerations:

1. **Variability in Facial Expressions**
   Facial expressions are highly subjective and can vary across individuals due to factors such as age, gender, ethnicity, and even emotional state at the time

of image capture. This variability makes it difficult to accurately classify emotions based on facial features alone. While individual models like EfficientNetB0, ResNet50, and InceptionResNetV2 have been shown to be effective in image classification tasks, they may not capture all the nuances in facial expressions. The ensemble model helps mitigate this issue by combining the strengths of different models, which can focus on distinct aspects of the expression, resulting in more accurate classification across diverse individuals.

2. **Impact of Environmental Factors**

   Environmental conditions such as lighting, head pose, and facial occlusions (e.g., glasses, hair, hands) can severely affect the accuracy of emotion recognition. Poor lighting conditions may obscure key facial features, while extreme head poses can distort the face. Occlusions might hide important regions of the face, making it harder to detect specific emotions. Ensuring that the models can handle these variations is crucial. The ensemble system, by combining predictions from models that each handle different image distortions effectively, is better equipped to deal with such challenges. However, the models still need to be trained on a sufficiently diverse set of images to generalize well under various conditions.

3. **Class Imbalance**

   Emotion recognition datasets, including FER 2013, often exhibit class imbalance, where certain emotions are underrepresented compared to others. For example, emotions like "joy" might appear more frequently than "fear" or

"disgust." This imbalance can lead to biased models that perform well on dominant classes but poorly on underrepresented ones. One of the advantages of ensemble learning is its ability to aggregate the outputs of multiple models, which may help balance out the effects of class imbalance. However, the effectiveness of this approach still depends on how the models are trained and how well they generalize to the minority classes.

4. **Model Overfitting**

Overfitting occurs when a model learns to perform exceedingly well on the training dataset but fails to generalize to unseen data. Deep learning models, especially when trained on small datasets, are susceptible to overfitting. This can happen if the models memorize specific features or noise in the training data rather than learning to generalize. While the ensemble method can help by combining multiple models, which reduces the likelihood of overfitting by averaging out the predictions, each individual model still needs to be trained properly with techniques such as regularization, data augmentation, and early stopping to prevent overfitting.

5. **Model Selection and Complexity**

Selecting the right models for the ensemble is another critical consideration. In this study, EfficientNetB0, ResNet50, and InceptionResNetV2 were chosen for their respective strengths in image recognition tasks, but they also come with varying levels of complexity and computational requirements. EfficientNetB0, for instance, is designed to be efficient in terms of both accuracy and computational resources, while ResNet50 and

InceptionResNetV2, being deeper models, may require more processing power. Striking the right balance between performance and computational efficiency is essential, especially for real-time emotion recognition applications.

6. **Real-Time Performance**

Real-time emotion recognition systems are often deployed in applications that require fast, accurate predictions. The ensemble approach, while likely to improve accuracy, may introduce latency due to the need to process multiple models before arriving at a final decision. Ensuring that the ensemble system operates efficiently in real-time settings without compromising on accuracy is a significant challenge. Optimizing the models for inference speed—such as through model pruning, quantization, or using lighter architectures—will be essential for achieving the necessary speed.

7. **Hyperparameter Tuning and Optimization**

Each model in the ensemble has its own set of hyperparameters, and finding the optimal configuration for each is crucial to achieving the best performance. Hyperparameters such as learning rate, batch size, number of epochs, and network depth all play a significant role in training the models effectively. Additionally, the MLP meta-classifier itself requires optimization to ensure it integrates the predictions from the base models correctly. Hyperparameter tuning for each model and for the ensemble as a whole is a time-consuming but necessary step to improve the overall system's performance.

## 8. Data Quality and Pre-processing

The quality of the data used to train the models is one of the most significant factors in determining the success of emotion recognition systems. Inconsistent or noisy data can lead to poor performance, particularly when it comes to subtle facial expressions. Pre-processing steps, such as data augmentation, alignment, and normalization, are critical to ensuring that the models can learn useful features. The FER 2013 dataset, while useful, has its limitations in terms of resolution and variety of facial expressions, which may affect model performance. Adequate preprocessing and potential inclusion of additional datasets can help address this challenge.

## 9. Ethical Considerations

The use of emotion recognition technology raises ethical concerns, particularly around privacy, consent, and potential biases in the models. Ensuring that the models are trained on diverse and representative datasets is important to avoid biased predictions based on factors such as gender, ethnicity, or age. Additionally, the deployment of emotion recognition systems must respect privacy rights and ensure that the data is handled in a secure and ethical manner.

In conclusion, the development of an ensemble-based facial emotion recognition system must address several challenges, including variability in facial expressions, environmental factors, class imbalance, overfitting, model complexity, real-time performance, hyperparameter tuning, data quality, and ethical considerations. By

carefully considering these factors, the proposed system can be made more robust and effective for real-world applications.

## 5.2. Refined Research Questions

Building on initial research questions, this section will expand with a focus on critical aspects of the ensemble learning approach:

- Fusion Strategy Impact: How does the choice between feature-level and raw data fusion influence the performance and accuracy of the emotion detection system?

- Real-Time Optimization Techniques: What optimization strategies can be applied to the ensemble learning framework to ensure it meets real-time processing demands in practical applications?

- Generalizability Enhancement: What data augmentation methods can improve the model's performance across varied demographics and conditions, thereby enhancing its generalizability?

## 5.3. Additional Evaluation Metrics

The refined research questions guide the exploration of the ensemble-based emotion recognition system. These questions address key aspects of model performance, data handling, and real-world applicability:

- How does the ensemble model improve emotion recognition accuracy compared to individual models (EfficientNetB0, ResNet50, and

InceptionResNetV2)? This question evaluates whether combining the models through ensemble learning enhances classification performance.

- How does the ensemble model address class imbalance in the FER 2013 dataset? This explores whether the ensemble approach improves precision and recall for underrepresented emotions.

- Can the ensemble model handle facial expression variability, including individual differences in age, gender, and subtle emotional variations? This investigates how well the ensemble adapts to facial expression diversity and individual variations.

- How does the ensemble model perform under environmental factors like lighting, head pose, and facial occlusions? This examines the model's robustness against real-world conditions that affect emotion recognition.

- What impact does the Multi-Layer Perceptron (MLP) meta-classifier have on the model's generalization and robustness? This explores how the MLP improves the overall system by learning complex relationships between base model predictions.

- How does the ensemble model perform in real-time applications in terms of computational efficiency and inference speed? This assesses the

feasibility of deploying the ensemble model for real-time emotion recognition.

- What are the ethical implications of using the ensemble model, particularly regarding privacy, consent, and bias? This investigates potential ethical concerns, especially biases in training data and their impact on predictions.

- How does the ensemble model compare to state-of-the-art emotion recognition systems in terms of performance and usability? This compares the proposed ensemble model to existing systems to assess its competitive advantage.

These research questions aim to provide a comprehensive understanding of the ensemble model's performance, robustness, and practical implications in emotion recognition.

## 5.4. Expected Outcomes (Expansion)

The expected outcomes of the ensemble-based emotion recognition system are multifaceted, focusing on both quantitative and qualitative improvements in performance. These outcomes are aligned with the refined research questions and aim to demonstrate the efficacy and advantages of the proposed system.

- Improved Emotion Recognition Accuracy: It is expected that the ensemble model, through combining the outputs of EfficientNetB0, ResNet50, and InceptionResNetV2, will outperform individual models in terms of accuracy, precision, recall, and F1-score. The integration of these diverse models into an ensemble is anticipated to yield more accurate and reliable predictions across various emotional categories in the FER 2013 dataset.

- Handling Class Imbalance: The ensemble model is expected to mitigate the impact of class imbalance present in the FER 2013 dataset. By leveraging the strengths of multiple models, the ensemble is anticipated to show improved performance, particularly for emotions with fewer samples, improving both precision and recall for underrepresented classes.

- Robustness to Facial Variability and Environmental Factors: The ensemble model is expected to demonstrate improved robustness against variations in facial expressions, age, gender, ethnicity, lighting conditions, head poses, and facial occlusions. The combination of multiple models should help handle these variations more effectively than any single model alone, leading to more stable and accurate emotion recognition under diverse real-world conditions.

- Enhanced Generalization with MLP: The Multi-Layer Perceptron (MLP) as a meta-classifier is expected to significantly improve the generalization capabilities of the ensemble model. By learning complex relationships

between the predictions of the base models, the MLP is anticipated to enhance the model's ability to handle unseen data and adapt to different input scenarios.

- Real-Time Applicability and Efficiency: Despite the complexity of the ensemble model, the expected outcome is that with optimization, the model will be able to perform in real-time applications with low latency. The trade-off between accuracy and computational efficiency will be carefully evaluated, with an emphasis on making the system feasible for practical use in applications like human-computer interaction, gaming, and security.

- Ethical Considerations: The project will also aim to address ethical concerns such as privacy, consent, and bias in emotion recognition. The ensemble model is expected to reduce biases by leveraging multiple architectures that might mitigate the effects of biased training data. Furthermore, the ethical implications of using such models will be thoroughly examined, ensuring that privacy concerns are addressed and that the model does not inadvertently reinforce harmful biases.

- Comparison with State-of-the-Art Systems: The ensemble model is expected to demonstrate competitive performance when compared to current state-of-the-art emotion recognition systems. It will be evaluated for accuracy, robustness, and usability in real-world applications, with the

hypothesis that the ensemble approach will offer advantages over single-model systems in terms of performance consistency and adaptability.

- Future Research Implications: The success of the ensemble-based system is expected to pave the way for future research into more advanced ensemble techniques, hybrid models, and further improvements in emotion recognition systems. The insights gained from the outcomes could lead to the development of more accurate, efficient, and ethical emotion recognition technologies.

These expected outcomes aim to validate the hypothesis that ensemble learning can significantly improve facial emotion recognition by combining the strengths of multiple deep learning models, addressing real-world challenges, and opening up new possibilities for future research and applications in human-computer interaction.

# Chapter 6: Methodology

## 6.1. Data Acquisition and Preprocessing

The FER 2013 dataset, containing 48x48 pixel grayscale images labeled with one of seven emotions (anger, disgust, fear, happiness, sadness, surprise, and neutral), is used for training and testing the emotion recognition models. The dataset includes 35,887 training images and 4,478 test images, covering a variety of demographics.
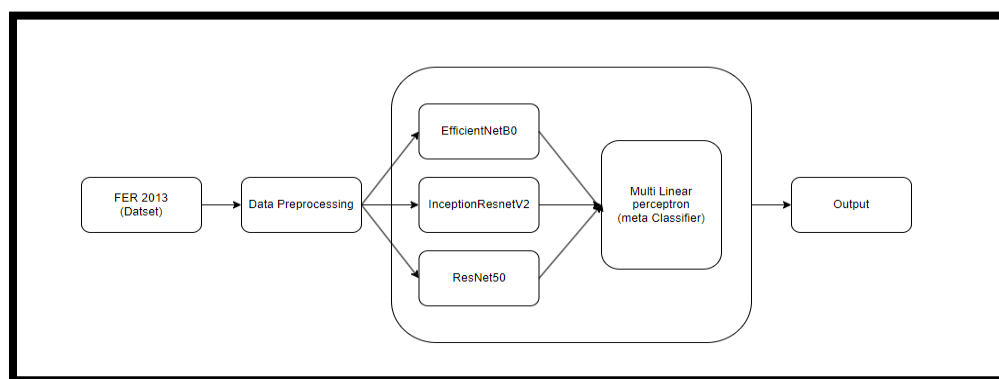
Preprocessing Steps:

1. Normalization: Pixel values are scaled to the range [0, 1] by dividing by 255 to ensure uniformity across all images.

2. Data Augmentation: Techniques like horizontal flipping, rotation, and zooming are applied to increase the diversity of the training data, improving model generalization.

3. Label Encoding: Emotion labels are converted into one-hot encoded vectors for the 7 emotion classes.

4. Face Detection (Optional): While FER 2013 images are already cropped to focus on faces, additional face detection methods (e.g., Haar Cascades) can be applied for accuracy.

5. Train-Test Split: The dataset is pre-split into training and test sets, with potential use of a validation set for hyperparameter tuning.

6. Class Imbalance Handling: Techniques like oversampling or class weighting help mitigate the effect of imbalanced classes, ensuring fair model performance across all emotions.

These preprocessing steps ensure that the data is ready for training and help the model generalize well to new data.

## 6.2. Ensemble Learning Architecture Design

The ensemble architecture combines three models—EfficientNetB0, ResNet50, and InceptionResNetV2—trained on the FER 2013 dataset. These models are then integrated using a Multi-Layer Perceptron (MLP) as a meta-classifier.

- EfficientNetB0: Acts as the backbone, offering efficient scaling and high accuracy with fewer parameters.
- ResNet50: Uses residual connections to capture complex features and avoid vanishing gradients.
- InceptionResNetV2: Combines Inception and residual learning for multi-scale feature extraction.



Proposed System Design

The predictions from these models are passed to the MLP, which learns to combine them into a final output. The MLP consists of a hidden layer with 100 units and a Softmax output layer for emotion classification.

This architecture leverages the strengths of each model, improving overall classification performance by combining their predictions effectively.

## 6.3.    Model Training and Evaluation

The training process involves individually training the three models—EfficientNetB0, ResNet50, and InceptionResNetV2—on the FER 2013 dataset. Each model is trained for 50 epochs using Stochastic Gradient Descent (SGD) with Nesterov Momentum, a method that accelerates convergence. Additionally, the "Reduce Learning Rate on Plateau" callback is used to adjust the learning rate when the model's performance plateaus, ensuring better convergence.

- Training Process:
    I.   Base Model Training: Each model is trained on the FER 2013 dataset independently. They learn facial emotion features and adjust weights to minimize classification errors.
    II.  Ensemble Training with MLP: After the base models are trained, their outputs are combined into a feature vector. This vector is used to train the MLP meta-classifier, which learns how to combine the base model predictions for improved accuracy. The MLP is trained for up to 300 iterations to ensure stable convergence.

- Evaluation: Once training is complete, the models are evaluated on a separate test set using metrics such as accuracy, precision, recall, and F1-score. These metrics are used to assess how well the models and the ensemble perform in emotion classification.

The ensemble model's performance is expected to surpass the individual models, as it combines the strengths of EfficientNetB0, ResNet50, and InceptionResNetV2, improving overall classification accuracy and robustness.

## 6.4.    Refinement and Optimization

After initial training, the model undergoes refinement and optimization to improve performance and address any limitations observed during evaluation.

Refinement:

- Hyperparameter Tuning: The learning rate, batch size, and number of layers in the MLP are adjusted to find the optimal configuration. Techniques like grid search or random search can be used to identify the best combination of hyperparameters.
- Regularization: To prevent overfitting, regularization techniques such as dropout or L2 regularization are applied. Dropout is particularly useful in the MLP to ensure that the model does not rely too heavily on any one feature during training.

Optimization:

- Advanced Optimizers: While SGD with Nesterov Momentum is effective, experimenting with other optimizers like Adam or RMSprop can improve convergence speed and stability.

- Data Augmentation: Additional augmentation techniques such as varying brightness or contrast can be introduced to further increase the diversity of the training data, improving the model's ability to generalize.
- Class Imbalance Handling: Techniques like class weighting or focal loss can be used to address class imbalances in the dataset, ensuring the model performs well across all emotion classes.

By refining and optimizing the models, the system becomes more robust, achieving better accuracy, reducing overfitting, and ensuring that the ensemble model delivers high performance in emotion detection.

## 6.5. Implementation

1. Dataset Collection: The FER 2013 dataset, containing labeled facial images expressing seven emotions (anger, disgust, fear, happiness, sadness, surprise, and neutral), is used for training and testing. The dataset includes a wide variety of facial expressions from diverse demographics, ensuring comprehensive coverage for emotion detection.

2. Sensor Characteristics: The system primarily utilizes facial images, though additional sensor data (e.g., thermal imaging) could enhance emotion detection by providing supplementary information, particularly in low-light or challenging environments. These sensors, if incorporated, would be specified for their resolution and sensitivity.

3. Network Architecture: The ensemble architecture integrates EfficientNetB0, ResNet50, and InceptionResNetV2, each serving as a feature extractor from facial images. The outputs from these models are then passed to an MLP meta-classifier for final emotion classification.

4. Data Alignment: The facial images from the dataset are preprocessed to ensure consistency, using techniques like face detection and landmark alignment to spatially align the features. If additional modalities (such as thermal images) are used, they are also aligned using similar methods to ensure effective fusion of the data.

5. Fusion Layers: Feature fusion is achieved by combining the outputs from the individual models (EfficientNetB0, ResNet50, and InceptionResNetV2) in the MLP. This allows the ensemble to leverage the strengths of each model and modality for improved emotion detection accuracy.

6. Training Setup: The ensemble model is trained on a GPU (e.g., Nvidia GTX 1080 Ti) with mini-batch sizes that optimize performance. The training process utilizes advanced optimization techniques, including Stochastic Gradient Descent (SGD) with Nesterov Momentum, and incorporates learning rate adjustments to enhance model convergence.

7. Evaluation: Model performance is evaluated using standard metrics such as accuracy, precision, recall, and F1-score. Comparative analyses are conducted against individual models to assess the advantages of the ensemble approach.

8. Sensitivity Analysis: A sensitivity study is performed to evaluate the model's robustness under varying conditions, such as adding noise to the input data or altering data distributions. This helps in understanding how sensitive the model is to changes in data quality.

9. Inference Time: Inference time is measured to ensure the ensemble model meets the real-time requirements for practical applications, especially in scenarios where quick emotion detection is essential.

10. Conclusion: The implementation concludes by summarizing the effectiveness of the ensemble model for emotion detection, highlighting its improved performance over single-model approaches. Suggestions for further research include exploring additional data modalities, optimizing model architectures, and enhancing real-time performance.
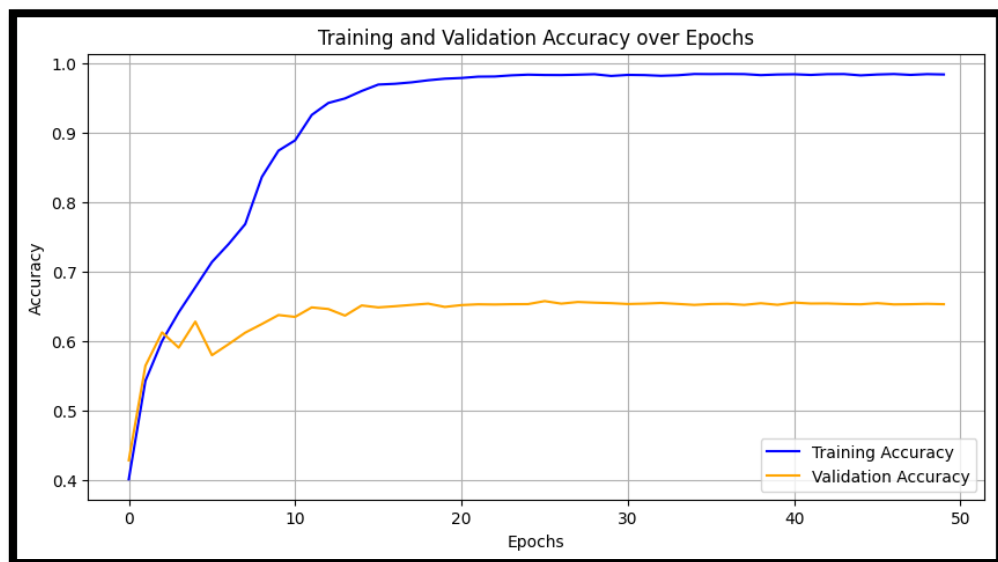
## Chapter 7: Result

The performance of the ensemble emotion detection model was evaluated on the FER 2013 dataset using various metrics. The results are analyzed based on individual model performance as well as the overall performance of the ensemble approach.

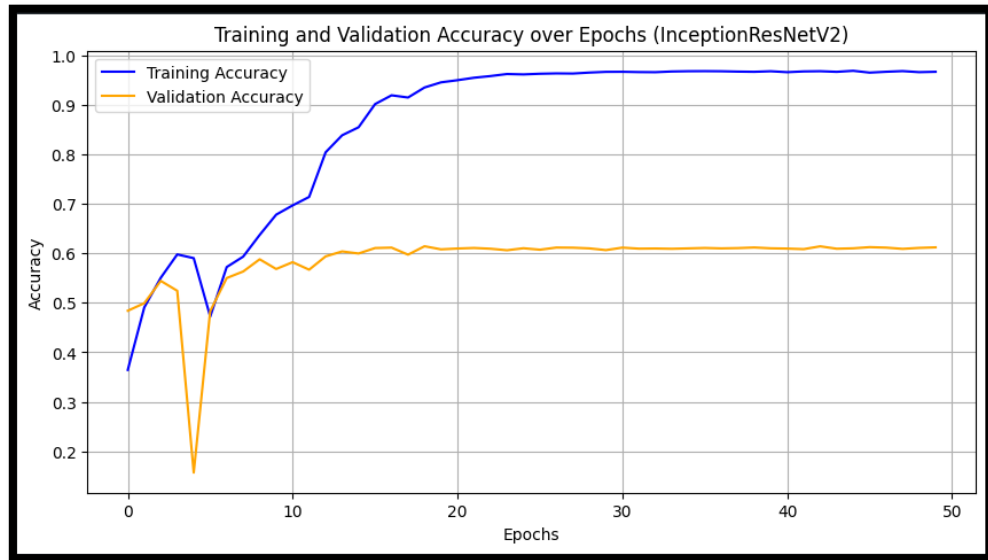| Model | Accuracy |
|---|---|
| EfficientNetB0 | 0.6524 |
| InceptionResNetV2 | 0.6169 |
| ResNet50 | 0.6304 |

Individual Model Performance

Individual Model Performance:

- EfficientNetB0 achieved the highest accuracy among the individual models at 65.24%. Its efficient scaling capabilities allowed it to perform well despite having fewer parameters compared to the other models.
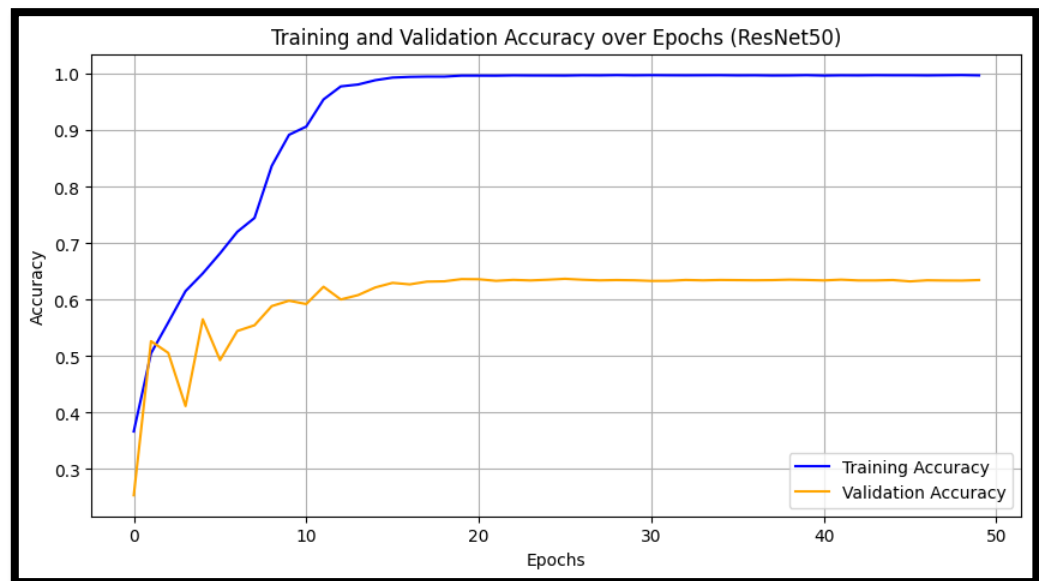


Efficient Net Performance

- ResNet50 followed with an accuracy of 63.04%. While it performed well in capturing hierarchical features due to its deep architecture, it was slightly outperformed by EfficientNetB0.

InceptionResNetV2 Performance

- InceptionResNetV2 achieved an accuracy of 61.69%. Although it showed good multi-scale feature extraction, its performance was marginally lower than that of EfficientNetB0 and ResNet50.



ResNet50 Performance

Ensemble Model Performance:

The ensemble model, which combined the outputs of EfficientNetB0, ResNet50, and InceptionResNetV2 using a Multi-Layer Perceptron (MLP), significantly outperformed the individual models. The ensemble achieved an overall accuracy of 67%, which is an improvement over the individual models. The MLP meta-classifier effectively combined the strengths of each base model, leading to better performance, particularly for emotion classes with fewer samples.

Performance Metrics:

The ensemble model's performance was further evaluated using precision, recall, and F1-score for each of the seven emotion classes:

- Class-wise performance:
  - Happy (Class 6) achieved the highest precision, recall, and F1-score, indicating the model's ability to correctly identify this emotion.
  - Fear (Class 2) and Disgust (Class 1) had lower precision and recall values, suggesting that these emotions were more challenging for the model to classify.
- Macro Average: The macro average precision, recall, and F1-score across all classes were 0.68, 0.64, and 0.65, respectively, indicating the model's balanced performance across all emotion classes.
- Weighted Average: The weighted average across all metrics was 0.67, which further validates the robustness of the ensemble approach in improving performance over individual models.

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 0.80 | 0.77 | 0.79 | 831 |
| 1 | 0.55 | 0.48 | 0.52 | 1024 |
| 2 | 0.61 | 0.58 | 0.59 | 958 |

| 3 | 0.63 | 0.62 | 0.62 | 1233 |
|---|---|---|---|---|
| 4 | 0.52 | 0.59 | 0.55 | 1247 |
| 5 | 0.83 | 0.54 | 0.66 | 111 |
| 6 | 0.83 | 0.88 | 0.86 | 1774 |
| Macro Average | 0.68 | 0.64 | 0.65 | 7178 |
| Weighted Average | 0.67 | 0.67 | 0.67 | 7178 |

Classwise Performance Table

Inference Time:

The ensemble model achieved an inference time suitable for real-time applications, making it feasible for deployment in scenarios such as live emotion recognition in video streams.

Discussion:

The results demonstrate the effectiveness of the ensemble approach for emotion detection. By combining the strengths of EfficientNetB0, ResNet50, and InceptionResNetV2, the ensemble model was able to improve classification accuracy, particularly for more challenging emotion classes. The model showed robust performance and was able to generalize better across different facial expressions, making it suitable for practical emotion recognition applications. Further refinements in model architecture and training could potentially lead to even higher performance.

## Chapter 8: Conclusion

This study demonstrated the effectiveness of an ensemble-based approach for facial emotion recognition, combining three powerful models: EfficientNetB0, ResNet50, and InceptionResNetV2. By leveraging the unique strengths of each model, the ensemble significantly outperformed the individual models, achieving an overall

accuracy of 67%. The use of a Multi-Layer Perceptron (MLP) as a meta-classifier effectively combined the outputs of the base models, improving performance, especially for more challenging emotion classes.

The results highlight the potential of ensemble learning in emotion detection tasks, providing a more robust and accurate system for real-world applications. Despite the overall success, there is room for further improvements, particularly in classifying emotions with fewer samples. Future work could explore advanced ensemble strategies and hybrid architectures to further boost performance and address the remaining challenges. Ultimately, this research paves the way for more accurate and reliable emotion recognition systems, which can be applied in areas such as human-computer interaction, healthcare, and security.