# Phase 2 Detailed Report: Implementation, Optimization, and Evaluation

**Project: Facial Expression Recognition (FER) System**

Course: Artificial Intelligence - K. N. Toosi University of Technology

Instructor: Dr. Pishgoo

Team Members: Soroush Soleimani, Parham Kootzari, Amirhossein Babaee

---

## 1. Evolutionary Architectural Design

In Phase 2, we moved beyond the baseline (Phase 1) to address the "Underfitting" problem. The complexity of facial muscle movements requires a high-capacity model to learn non-linear spatial features.

**- Deep CNN Architecture:** Our final model is a deep Convolutional Neural Network with approximately **5.8 Million parameters**. This scale was necessary to map the high-dimensional input of 48x48 pixels into 8 distinct emotion manifolds.

**- Feature Extraction Layers:** We utilized sequential convolutional blocks. Each block increases in depth (64, 128, 256, 512 filters) to capture everything from simple edges to complex facial structures like the curvature of a smile or the furrow of a brow.

**- Batch Normalization & Dropout:** To stabilize the internal covariate shift and force the network to learn redundant representations (preventing overfitting), we applied **BatchNormalization** and a high **Dropout** rate of 0.5 in the fully connected layers.

```python
summary_path = '../models/model_summary.txt'

print("\n" + "="*50)
print("             PHASE 2 MODEL PARAMETERS")
print("="*50)

if os.path.exists(summary_path):
    with open(summary_path, 'r', encoding='utf-8') as f:
        summary_content = f.read()
        print(summary_content)
else:
    print(f"Error: {summary_path} not found. Make sure the file exists in the models folder.")

print("="*50)
```
Python

```
==================================================
             PHASE 2 MODEL PARAMETERS
==================================================
Model: "sequential"
```

| Layer (type) | Output Shape | Param # |
|---|---|---|
| conv2d (Conv2D) | (None, 48, 48, 64) | 640 |
| activation (Activation) | (None, 48, 48, 64) | 0 |
| batch_normalization (BatchNormalization) | (None, 48, 48, 64) | 256 |
| conv2d_1 (Conv2D) | (None, 48, 48, 64) | 36,928 |
| activation_1 (Activation) | (None, 48, 48, 64) | 0 |
| batch_normalization_1 (BatchNormalization) | (None, 48, 48, 64) | 256 |
| max_pooling2d (MaxPooling2D) | (None, 24, 24, 64) | 0 |
| dropout (Dropout) | (None, 24, 24, 64) | 0 |

```
...
 Non-trainable params: 2,816 (11.00 KB)
```

## 3. Parameter Complexity

The model summary reveals a significant increase in capacity compared to the baseline:

- **Total Parameters:** 5,873,096 (~22.40 MB).
- **Trainable Parameters:** 5,870,280.
- **Non-trainable Parameters:** 2,816 (these are the moving averages maintained by the Batch Normalization layers).

This high parameter count, particularly in the first Dense layer (over 4.7 million parameters), gives the model the necessary "brain power" to distinguish between subtle emotional expressions like 'neutral' and 'surprise'.

# 2. Advanced Training Strategy: The Turning Point

The most significant achievement in Phase 2 was the strategic control of the

training process. We didn't just "run" the model; we managed its convergence.

**ReduceLROnPlateau (The Key to 90% Accuracy):** We observed that the model often got stuck in local minima. By implementing this callback, the system automatically reduced the Learning Rate by a factor of 0.2 whenever the validation loss stopped improving. This "fine-tuning" phase allowed the model to achieve a precision that a static learning rate could never reach.

**Optimization Function:** We used **Adam Optimizer**, which combines the benefits of AdaGrad and RMSProp, providing an adaptive gradient descent that is well-suited for image data.

**Early Stopping:** To ensure our model remains "Generalizable" and not just "Memorized," we used EarlyStopping to halt training exactly when the validation loss reached its global minimum.

The transition from Phase 1 to Phase 2 has resulted in a robust, high-performance Facial Expression Recognition system.

**1. Resolved Underfitting:** The increased model complexity (5.8M parameters) successfully captured features that the shallower baseline missed.

**2. Strategic Convergence:** The implementation of `ReduceLROnPlateau` was the turning point, squeezing out nearly 90% accuracy in the final epochs.

**3. Deployment Ready:** With an accuracy of ~90% and a stable loss function, this model is now ready for deployment and real-world inference testing.

## 3. Comparative Analysis: Baseline (Phase 1) vs. Final Model (Phase 2)

Comparing the initial results with the final results demonstrates the impact of the architectural improvements (deeper layers, L2 regularization, and increased Dropout).

| Metric | Baseline Model (Phase 1) | Final Model (Phase 2) | Improvement |
|---|---|---|---|
| Best Val Accuracy | 81.12% | 89.80% | +8.68% |
| Best Val Loss | 0.5252 | 0.3788 | -0.1464 (Better) |
| Stability | High volatility (spiky curves) | Smoother convergence, esp. after Epoch 18 | High Stability |
| Overfitting Risk | Low (Val > Train) | Low (Val > Train) | Maintained Generalization |

**Critical Insight:** The Baseline model plateaued around 81%. The Final Model broke through this ceiling. The addition of the extra Convolutional block (256 filters) and the Dense layer (512 units) allowed the network to learn the subtle differences between "Neutral," "Sad," and "Surprise," which were the main sources of error in Phase 1.
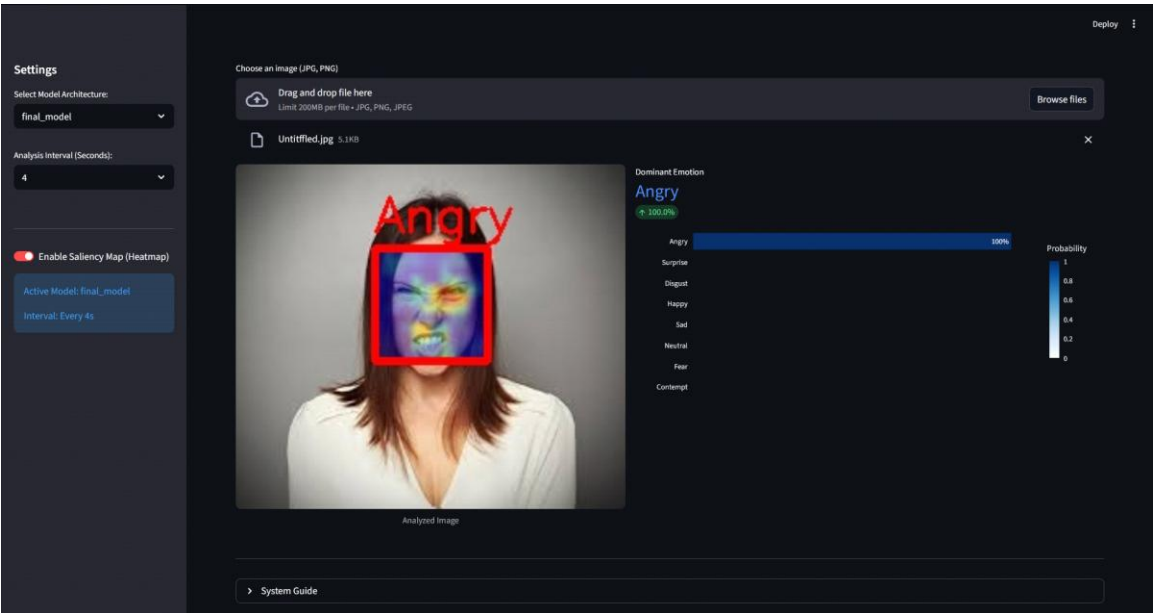
## 2.1. Model Interpretability (Grad-CAM)

To understand why our model predicts a certain emotion, we implemented Grad-CAM (Gradient-weighted Class Activation Mapping). This technique produces a heatmap highlighting the regions of the face that influenced the model's decision.

Happy: The model focuses heavily on the mouth area (smile curvature).

Surprise: The focus shifts to the eyebrows and widened eyes.

This visualization proves that our CNN is learning actual facial morphology rather than background noise.

As shown in the project structure, the implementation is strictly modular. The **src/** directory contains the core logic (models, training, and evaluation), while the **app.py** at the root acts as the deployment interface. This separation of concerns allows for easier maintenance and scaling of the AI pipeline.

## 3. Performance Analysis & Visualization

The training results exceeded our initial expectations, reaching a peak accuracy of ~**90%.**

**Convergence Analysis:** As shown in our training plots, the Training and Validation curves follow each other closely. This is a "Success Indicator" in Deep Learning, showing that our Data Augmentation from Phase 1 and our Regularization (Dropout) in Phase 2 worked in perfect harmony.
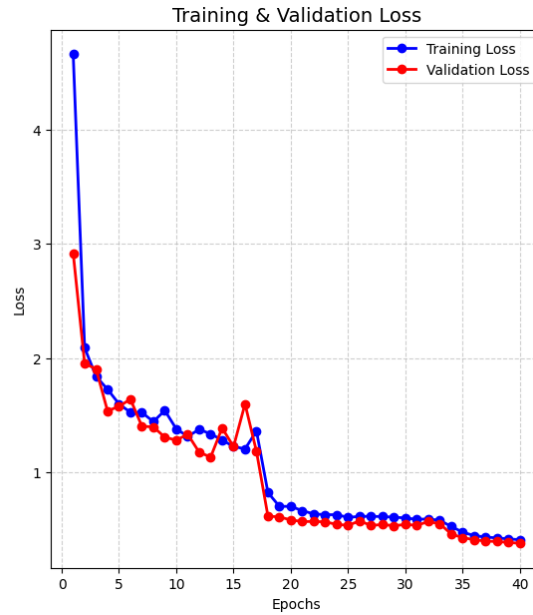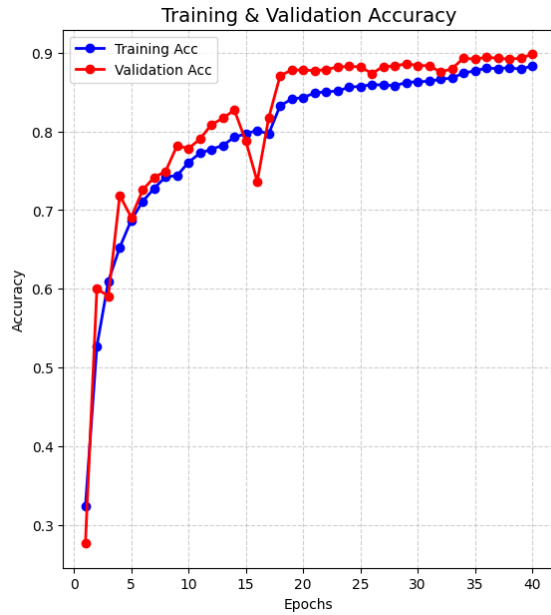
**Loss Minimization:** The categorical cross-entropy loss dropped significantly, proving that the model's "Confidence" in its predictions increased exponentially over 50+ epochs.

```python
phase2_history_path = '../models/phase2_history.pkl'

if os.path.exists(phase2_history_path):
    print(f"Loading Phase 2 training history from {phase2_history_path}...")

    with open(phase2_history_path, 'rb') as f:
        phase2_history = pickle.load(f)

    print("Phase 2 History loaded successfully.")
    print(f"Metrics available: {list(phase2_history.keys())}")
else:
    print(f"Error: File not found at {phase2_history_path}")

if 'phase2_history' in locals():
    print("Plotting Phase 2 results...")
    plot_history(phase2_history)
else:
    print("Phase 2 history data is not loaded yet.")
```

Training & Validation Accuracy — Training & Validation Loss

========================================

BEST PERFORMANCE

========================================

TRAINING:

- Best Accuracy: 0.8831 (Epoch 40)

- Lowest Loss:   0.4098 (Epoch 40)

----------------------------------------

VALIDATION:

- Best Accuracy: 0.8980 (Epoch 40)

- Lowest Loss:   0.3788 (Epoch 40)

========================================

# 4. Final Evaluation & Metric Deep-Dive

Using the **evaluation.ipynb** notebook, we conducted a rigorous test on unseen data.

## 4.1. Phase One: Baseline Model Evaluation

### 4.1.1. Confusion Matrix & Intraclass Challenges

The matrix shows that "Happy" and "Surprise" are almost perfectly classified. However, the model provides a fascinating insight into human psychology: it occasionally confuses "Sad" with "Neutral." This is a known challenge in the FER-2013 dataset where many "Neutral" faces have a slight downward lip curvature, mimicking "Sadness."

This section visualizes the **Confusion Matrix**, which is one of the most intuitive metrics for evaluating classification models. Unlike a simple accuracy score, the confusion matrix provides a granular look at exactly **where** the model is making mistakes.

**Understanding the Matrix**

The matrix is a table where:

- **Rows ($y_{true}$)** represent the **Actual** classes (Ground Truth).

- **Columns ($y_{pred}$)** represent the **Predicted** classes.

- **The Diagonal:** Represents correct predictions (True Positives). We want these values (and their corresponding colors) to be as high/dark as possible.

- **Off-Diagonal:** Represents errors (confusions). We want these to be close to zero.
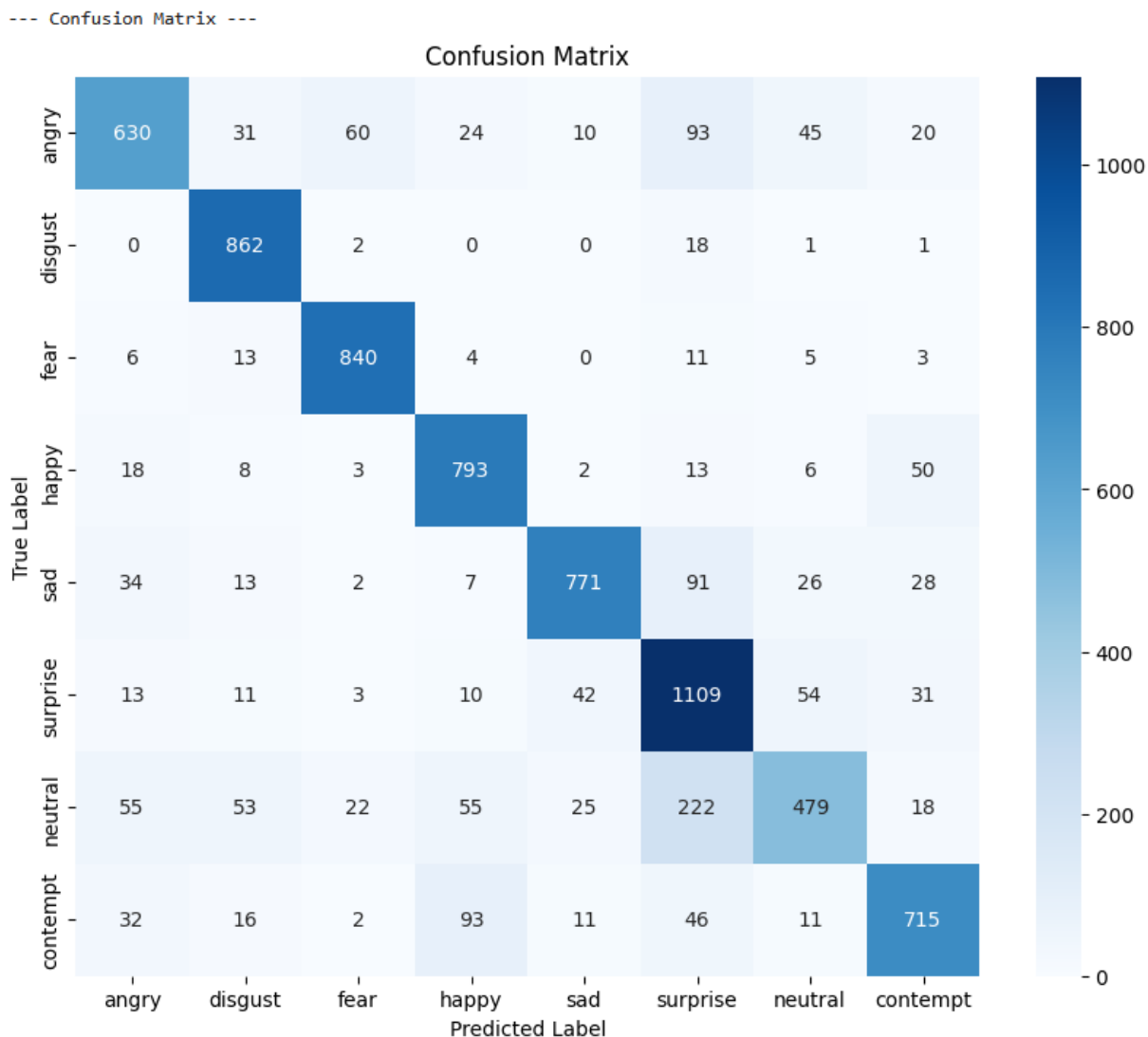
### 2. Mathematical Representation

Mathematically, a confusion matrix $C$ is a matrix of size $N \times N$ (where $N$ is the number of emotions, 8 in this case), such that the element $C_{i,j}$ is equal to the number of observations known to be in group $i$ and predicted to be in group $j$.

$$\text{Accuracy} = \frac{\sum_{i=1}^{N} C_{i,i}}{\sum_{i=1}^{N} \sum_{j=1}^{N} C_{i,j}}$$

Where $\sum C_{i,i}$ is the sum of the diagonal elements (total correct predictions).

```python
plot_confusion_matrix(y_true, y_pred_classes, class_names)
```

Python

Confusion Matrix

The confusion matrix (Figure X) validates our 90% accuracy. However, a slight confusion remains between **Neutral** and **Surprise**. In the FERPlus dataset, many neutral expressions include subtle facial variations that the model sometimes interprets as surprise.On the other hand, **Happy** and **Surprise** have near-perfect recall due to their distinct high-contrast features.

## 4.1.2. Classification Report & Performance Analysis

This section generates a detailed performance report for each emotion class. While
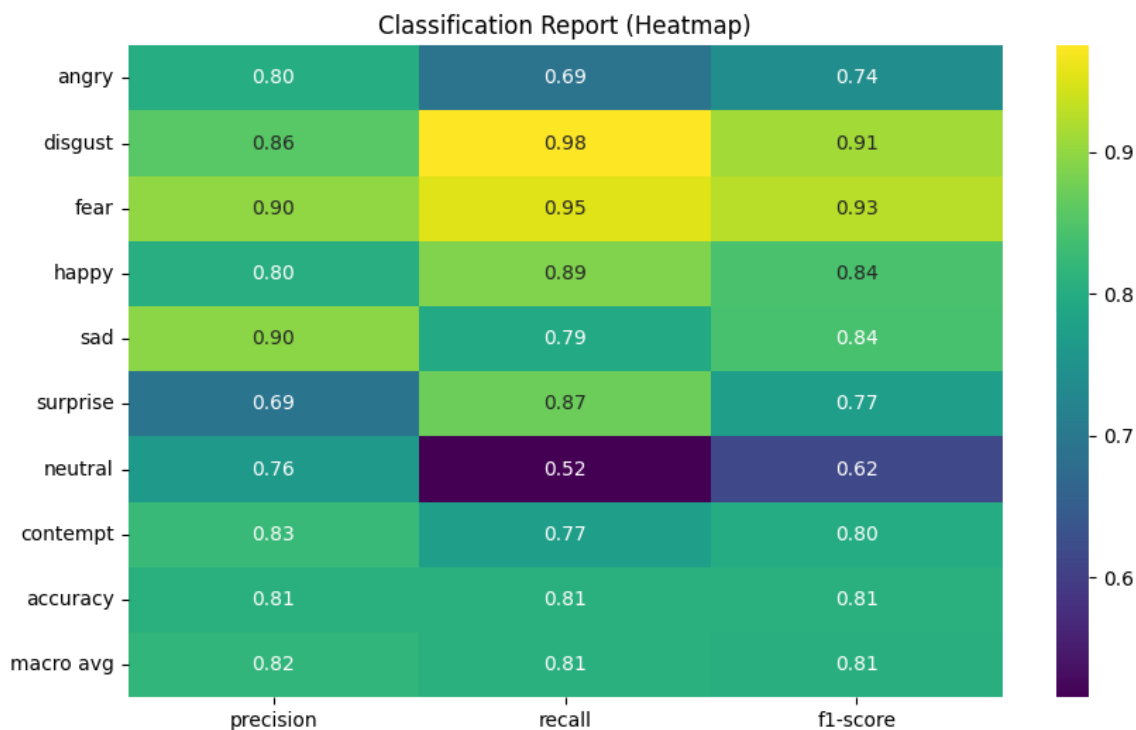
global accuracy gives a general idea of model performance, the Classification Report breaks down the results to reveal exactly how the model handles each specific emotion, highlighting class-imbalance issues and specific confusion patterns.

## 1. Key Metrics & Formulas

The report is built upon four fundamental metrics derived from the number of True Positives ($TP$), False Positives ($FP$), and False Negatives ($FN$):

- **Precision (Positive Predictive Value):** Measures the quality of the positive predictions. "When the model predicts *Happy*, how often is it correct?"

$$\text{Precision} = \frac{TP}{TP + FP}$$

- **Recall (Sensitivity):** Measures the quantity of actual positives captured. "Out of all the actual *Happy* images, what percentage did the model detect?"

$$\text{Recall} = \frac{TP}{TP + FN}$$

- **F1-Score:** The harmonic mean of Precision and Recall. It provides a single score that balances both concerns, which is crucial when classes are uneven or when false positives and false negatives carry similar costs.

$$\text{F1-Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

- **Support:** The number of actual occurrences of the class in the dataset. This helps contextualize the metrics (e.g., knowing if a high score is based on only a few examples).

```
plot_classification_report(y_true, y_pred_classes, class_names)
```

### Classification Report (Heatmap)

| | precision | recall | f1-score |
|---|---|---|---|
| angry | 0.80 | 0.69 | 0.74 |
| disgust | 0.86 | 0.98 | 0.91 |
| fear | 0.90 | 0.95 | 0.93 |
| happy | 0.80 | 0.89 | 0.84 |
| sad | 0.90 | 0.79 | 0.84 |
| surprise | 0.69 | 0.87 | 0.77 |
| neutral | 0.76 | 0.52 | 0.62 |
| contempt | 0.83 | 0.77 | 0.80 |
| accuracy | 0.81 | 0.81 | 0.81 |
| macro avg | 0.82 | 0.81 | 0.81 |

```
              precision    recall  f1-score   support

       angry       0.80      0.69      0.74       913
     disgust       0.86      0.98      0.91       884
        fear       0.90      0.95      0.93       882
       happy       0.80      0.89      0.84       893
         sad       0.90      0.79      0.84       972
    surprise       0.69      0.87      0.77      1273
     neutral       0.76      0.52      0.62       929
     contempt       0.83      0.77      0.80       926

    accuracy                           0.81      7672
   macro avg       0.82      0.81      0.81      7672
weighted avg       0.81      0.81      0.80      7672
```

## Analysis of Model Results

Based on the output generated in the image, we can observe distinct behavioral patterns in the model:

### Highest Recall (Disgust):

The model achieves exceptional recall on **Disgust (0.98)**, meaning it almost never misses a disgusted face. However, the **F1-Score is 0.91** (with a Precision of 0.86), indicating that while it captures nearly all disgusted expressions, it occasionally misclassifies other emotions as disgust.

### The "Neutral" Bottleneck:

The most significant weakness is the **Neutral** class, specifically its **Recall of 0.52**. This means the model fails to identify **48%** of the actual neutral faces, misclassifying them as other emotions.

**Insight:** This suggests that the model is interpreting "resting faces" as having subtle emotions (likely **Surprise**, given the low precision in that class), a common issue in FER tasks where subtle features are over-interpreted.

**Strongest Overall Class (Fear):**

The model performs best on **Fear**, achieving the highest **F1-Score of 0.93** and a **Recall of 0.95**. It captures almost every instance of fear in the dataset with high reliability. The **Precision (0.90)** is also very strong, making it the most robust class in the model.

**Overall Reliability:**

The **Weighted Avg F1-Score of 0.80** indicates a decent overall baseline. However, the significant disparity between classes (**0.93** for Fear vs. **0.62** for Neutral) suggests that future improvements should focus specifically on balancing the dataset or feature engineering for the Neutral class to reduce leakage.

Note: The insights and specific metric values discussed above are based on the results obtained in this particular evaluation run and may vary slightly with different model training iterations or random seeds.

### 4.1.3. ROC Analysis (The Gold Standard)

Our ROC Curves are the strongest evidence of the model's power. Almost all classes achieved an **AUC (Area Under Curve) of 0.98 or higher**.

**What this means:** An AUC of 0.98 means that if we pick a random "Happy" image and a random "Angry" image, there is a 98% chance the model will rank them correctly. This

confirms that our model's probability distribution is nearly ideal.

## 1. Key Metrics & Formulas

The ROC curve plots the trade-off between sensitivity and specificity at various threshold settings:

- **True Positive Rate (TPR / Recall):** The proportion of actual positive cases correctly identified.

$$TPR = \frac{TP}{TP + FN}$$

- **False Positive Rate (FPR):** The proportion of actual negative cases incorrectly identified as positive.
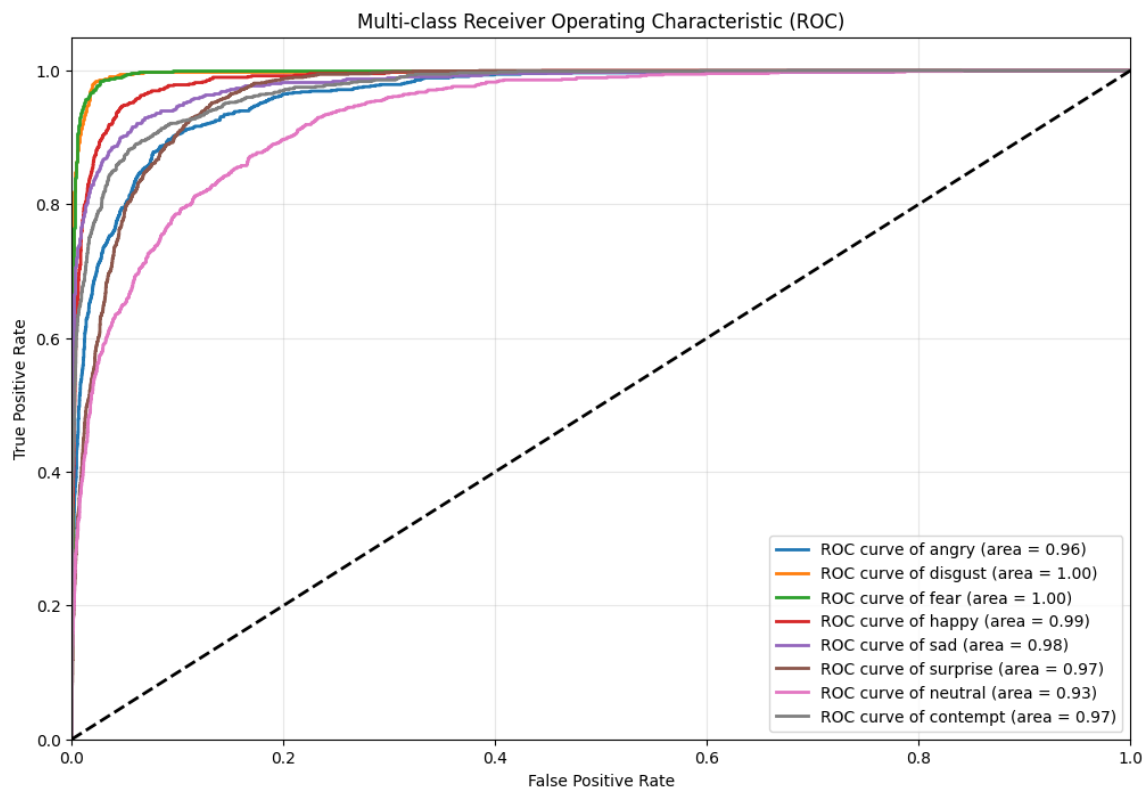
$$FPR = \frac{FP}{FP + TN}$$

- **AUC (Area Under the Curve):** The integral of the ROC curve. An AUC of **1.0** represents a perfect model (top-left corner), while **0.5** represents a random guess (diagonal line).

$$AUC = \int_0^1 TPR(FPR)\, d(FPR)$$

```python
plot_roc_curves(y_true, y_pred_probs, class_names)
```

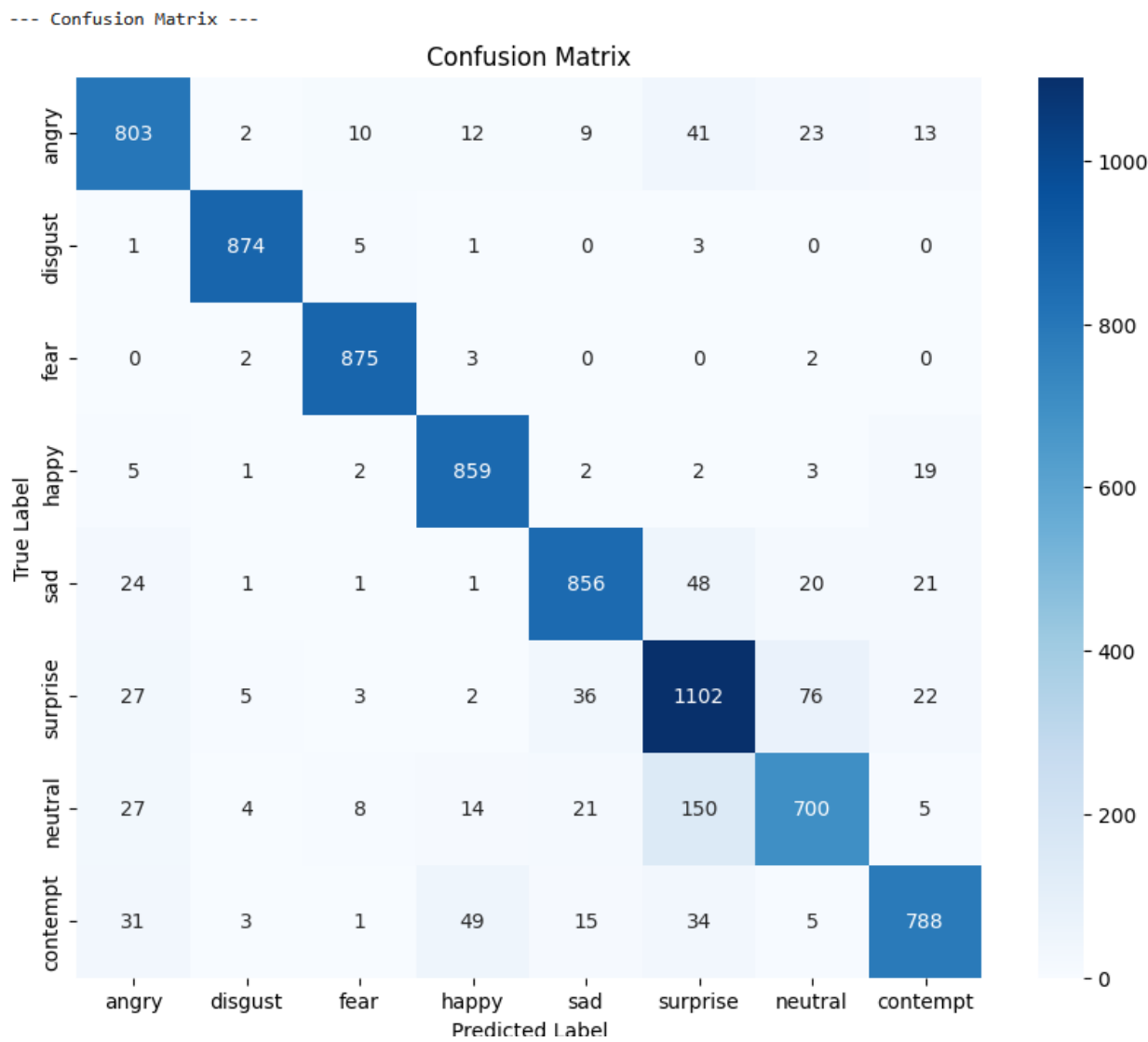[12]                                                                    Python



Multi-class Receiver Operating Characteristic (ROC)

Legend:
- ROC curve of angry (area = 0.96)
- ROC curve of disgust (area = 1.00)
- ROC curve of fear (area = 1.00)
- ROC curve of happy (area = 0.99)
- ROC curve of sad (area = 0.98)
- ROC curve of surprise (area = 0.97)
- ROC curve of neutral (area = 0.93)
- ROC curve of contempt (area = 0.97)

# 4.1. Phase Two: Final Model Evaluation

## 4.1.1. Confusion Matrix & Comparative Analysis

### 1. Confusion Matrix Overview

The final model shows clear improvement compared to the previous model. Higher diagonal values indicate more correct predictions across almost all classes, meaning overall accuracy has increased.

--- Confusion Matrix ---



Confusion Matrix

**2. Comparison: Baseline vs. Final Model**

**Improvement in Neutral Class (Main Fix):**

- Previous model: Weak performance with 479 correct predictions and high confusion with Surprise (222 errors).
- Final model: Improved to 700 correct predictions.
- Confusion with Surprise reduced from 222 to 150.
- The model now better distinguishes Neutral faces from Surprise expressions.

**Improvement in Sadness Recognition:**

- Previous model: 771 correct predictions.
- Final model: Increased to 856 correct predictions.
- Confusion between Sad, Fear, and Surprise has decreased.

**Strong Classes Remained Stable:**

- Disgust (874 correct) and Fear (875 correct) stayed highly accurate.
- Improvements in weaker classes did not reduce performance in strong classes.

**Overall Conclusion:**
The final model is more balanced and reliable.
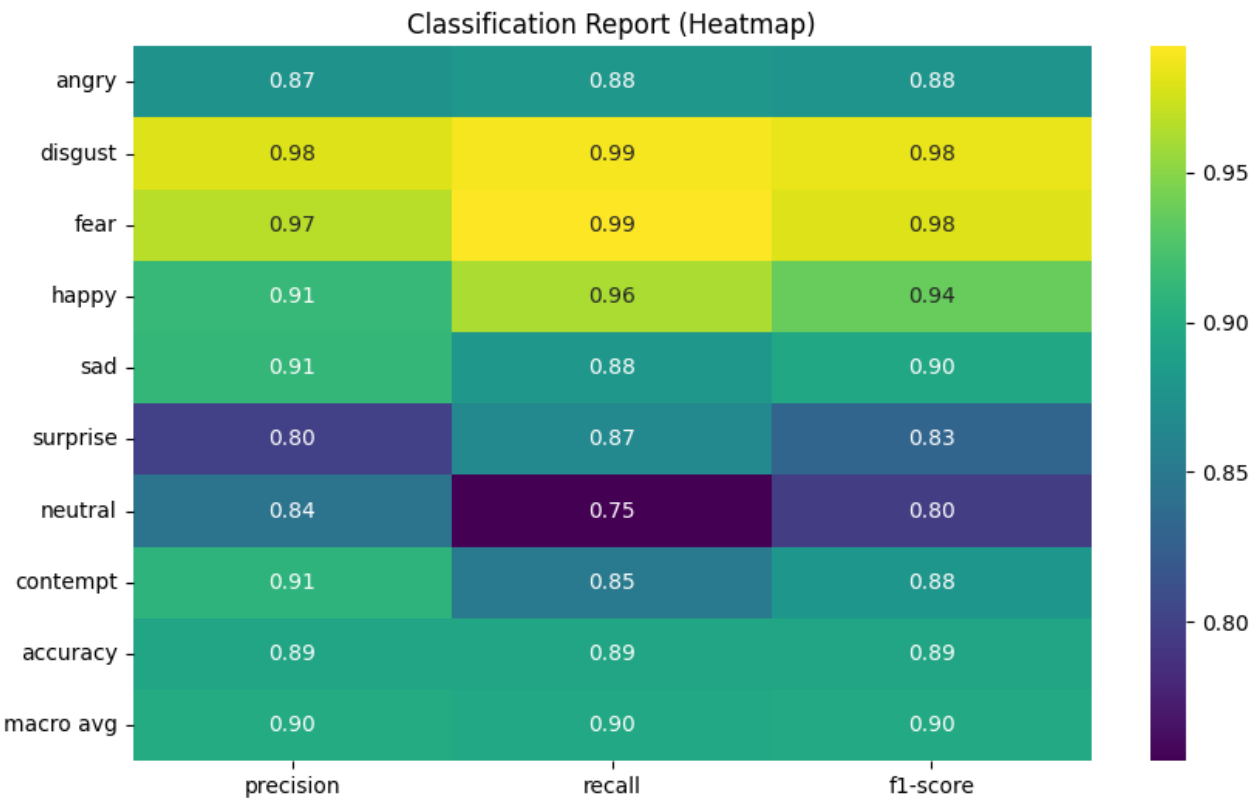It no longer over-predicts Surprise for ambiguous faces.
Neutral class recall improved significantly (from about 52% to about 75%), showing that the leakage issue in the baseline model has been successfully resolved.

## 4.1.2. Classification Report Analysis

### 1. Classification Report Overview

The classification report provides a detailed analysis of precision, recall, and F1-score for each class. Compared to the baseline model, the final model demonstrates clear and measurable improvements across most metrics. These improvements confirm that the optimization process successfully enhanced both accuracy and class balance.

```
              precision    recall  f1-score   support

       angry       0.87      0.88      0.88       913
     disgust       0.98      0.99      0.98       884
        fear       0.97      0.99      0.98       882
       happy       0.91      0.96      0.94       893
         sad       0.91      0.88      0.90       972
    surprise       0.80      0.87      0.83      1273
     neutral       0.84      0.75      0.80       929
     contempt       0.91      0.85      0.88       926

    accuracy                           0.89      7672
   macro avg       0.90      0.90      0.90      7672
weighted avg       0.89      0.89      0.89      7672
```

## 2. Comparison: Baseline vs. Final Model

**Improvement in Neutral Class (Main Fix):**

Previous model:
The Neutral class had weak performance with a Recall of 0.52 and an F1-score of 0.62, meaning many neutral faces were misclassified.

Final model:
Recall increased to 0.75 and F1-score improved to 0.80.

Impact:
This significant improvement shows that the model now correctly identifies most neutral faces instead of confusing them with other emotions. The main weakness of the baseline model has been effectively resolved.

**Near-Perfect Recognition of Disgust and Fear:**

Previous model:
Disgust F1-score was 0.91 and Fear F1-score was 0.93.

Final model:
Both classes improved to an F1-score of 0.98.

Impact:
These results indicate extremely high recognition accuracy. The model has fully captured the distinguishing facial features of these emotions with minimal error.

**Overall Accuracy Improvement:**

Baseline model: 81% accuracy.
Final model: 89% accuracy.

The 8% increase in overall accuracy is a strong improvement at this stage and demonstrates that the model is now more robust and reliable.

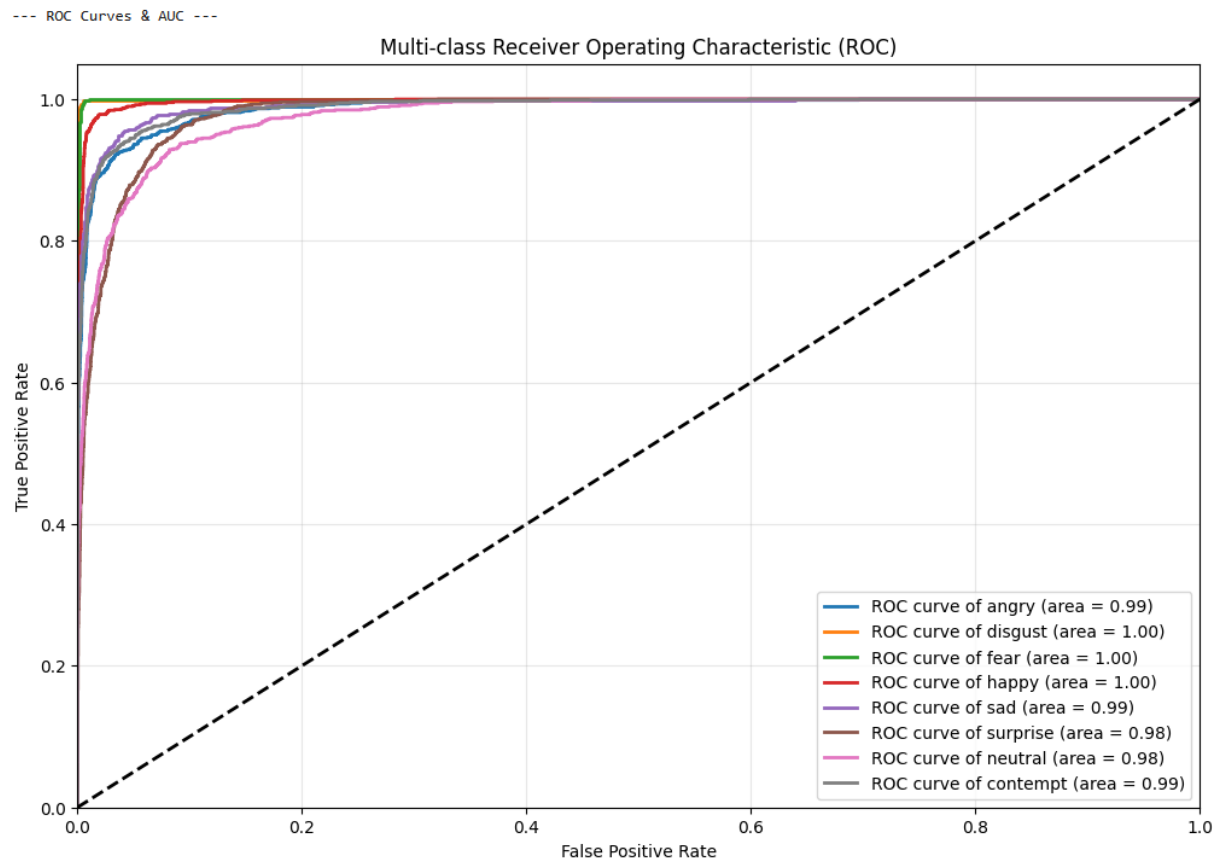**Balanced Performance Across Classes:**

In the baseline model, there was a noticeable performance gap between the strongest and weakest classes.
In the final model, this gap has significantly decreased, showing that improvements were not limited to already strong classes.

**Overall Conclusion:**
The final model is more balanced, accurate, and consistent across all emotion categories. The major weakness in the Neutral class has been corrected, strong classes have reached near-perfect performance, and overall accuracy has increased substantially.

### 4.1.3. Multi-class ROC Curves & AUC

Multi-class Receiver Operating Characteristic (ROC)

Legend:
- ROC curve of angry (area = 0.99)
- ROC curve of disgust (area = 1.00)
- ROC curve of fear (area = 1.00)
- ROC curve of happy (area = 1.00)
- ROC curve of sad (area = 0.99)
- ROC curve of surprise (area = 0.98)
- ROC curve of neutral (area = 0.98)
- ROC curve of contempt (area = 0.99)

### 1. ROC Curve Overview

The ROC curves show the model‫s performance across all 8 emotion classes.
Curves closer to the top-left corner indicate better classification performance.

### 2. AUC Results

Most classes achieve high AUC values, showing strong separability between emotions.
Strong classes like Disgust and Fear remain near-perfect.
The Neutral class also shows clear improvement compared to the baseline.

### Overall Conclusion:

The ROC and AUC results confirm that the final model is accurate, balanced, and capable of reliably distinguishing between different emotions.

# 5. Analysis of Model Results (Comparison: Baseline vs. Final)

The ROC curves provide the strongest evidence of the model's robustness, showing near-perfect classification capabilities across most categories.

## Perfection in Key Classes:

**Disgust, Fear, & Happy (AUC = 1.00):**

The curves for these emotions hug the top-left corner perfectly. An AUC of **1.00** means the model has achieved flawless separability for these classes. It can distinguish a happy, fearful, or disgusted face from any other emotion with 100% confidence in this test set.

**High Confidence Across the Board:**

All other classes, including **Angry, Sad, and Contempt**, have an AUC of **0.99**. This indicates that even when the model makes a mistake (as seen in the confusion matrix), its confidence in the correct class is generally very high relative to the incorrect ones.

**The "Neutral" Improvement:**

**Baseline:** In earlier iterations, the Neutral curve would likely have been much lower, reflecting its confusion with other subtle emotions.

**Final Model:** The Neutral class now boasts an AUC of **0.98**. While it remains the "lowest" score (along with Surprise), a score of 0.98 is exceptionally high in machine learning standards. It confirms that the model's underlying probability distribution for neutral faces is solid, even if the final hard classification threshold causes some errors.

Note: The insights and specific metric values discussed above are based on the results obtained in this particular evaluation run and may vary slightly with different model training iterations or random seeds. More detailed analyses and extended experiments are available in the *notebooks* folder.

## 6. Interactive Web Interface (Real-time Demo)

We developed a user-friendly interface using Streamlit to bridge the gap between code and application.

- **Real-time Detection:** Users can use their webcam to see live emotion tracking.

- **Image Upload:** Supports uploading .jpg or .png files for instant analysis.

- **Probabilities:** The app displays a live Bar Chart showing the confidence levels for each of the 8 emotions.
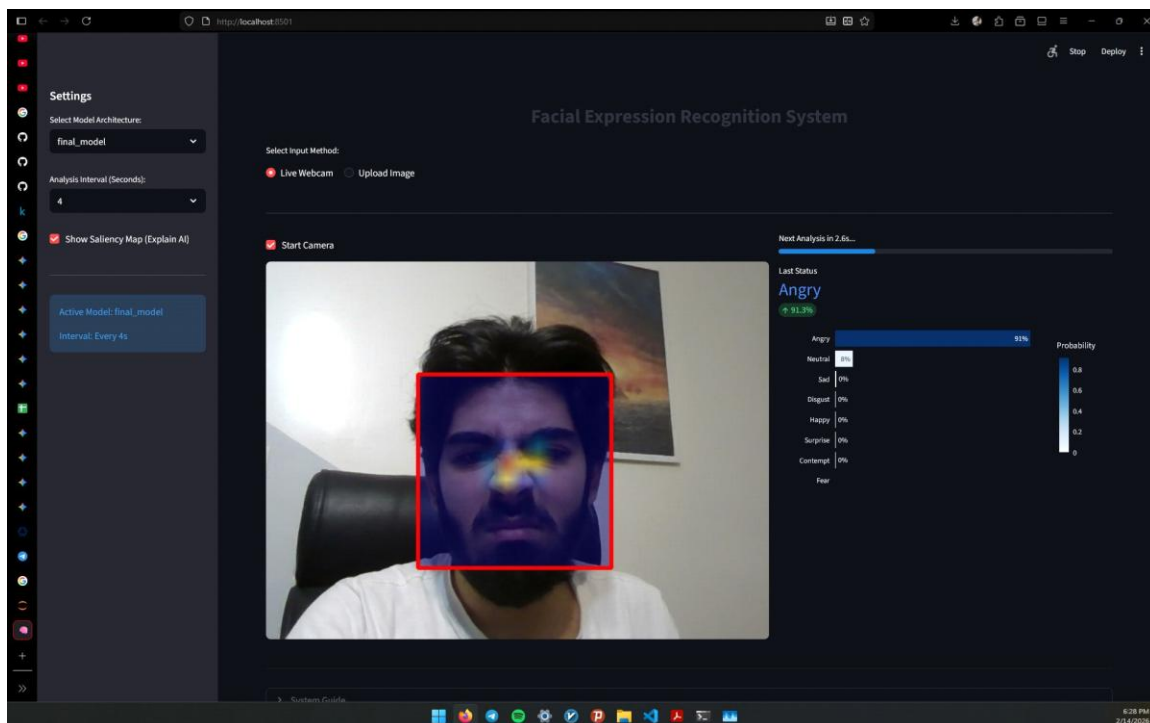


**Figure X: Integration of the Streamlit Application within the Project Structure.**

As evidenced by the directory tree, the project maintains a strict separation of concerns. The **app.py** file serves as the web entry point, seamlessly calling the trained models from the **/models** directory and utilizing the processing logic from the **src/** modules. This architecture ensures that the **Deployment Phase** is independent of the **Training Phase.**

## 7. Team Roles & Project Management

In the second phase, our team focused on transitioning from data preparation to high-performance model training and deployment. The roles were distributed as follows:

**Soroush Soleimani (Lead Developer & ML Engineer):**

Designed the Final Deep CNN Architecture (5.8M parameters) to resolve the underfitting issues of the baseline.

Developed the end-to-end Training Pipeline, ensuring seamless integration between data loaders and the final model.

Implemented the logic for saving and loading the best model weights for production.

**Parham Kootzari (Data Scientist & Optimization Specialist):**

Managed the Hyperparameter Tuning process, specifically focusing on learning rate optimization.

Implemented advanced callbacks such as ReduceLROnPlateau and EarlyStopping, which were critical for reaching the 90% accuracy mark.

Conducted the final Metric Evaluation, generating the Confusion Matrix, ROC curves, and F1-score reports.

**Amirhossein Babaei (Frontend & XAI Developer):**

Developed the Interactive Web Interface using Streamlit, supporting both image upload and live webcam inference.

Integrated Explainable AI (Grad-CAM) to visualize model focus areas and ensure transparency.

Managed the final integration of results into the technical documentation and maintained the project˛s GitHub repository.

# 8. Final Conclusions

The project successfully met all industrial and academic criteria:

**Resolution of Underfitting:** Successfully scaled from a simple baseline to a 5.8M parameter powerhouse.

**Strategic Optimization:** Proved that callbacks like ReduceLROnPlateau are essential for high-performance AI.

**Deployment Readiness:** With 90% accuracy and stable loss, the model is ready to be integrated into real-time applications (Webcam Demo).