# MODEL SELECTION, PERFORMANCE AND IMPLEMENTATION

*Facial Recognition and Pain Detection in Stroke Patients*

| | | |
|---|---|---|
| Adam Sabsabi | 585693 | 585693@student.inholland.nl |
| Fejsal Aziz | 581934 | 581934@student.inholland.nl |
| Maria Diaconescu | 683265 | 683265@student.inholland.nl |
| Mohamed Dinle | 700152 | 700152@student.inholland.nl |

**Client:** Inholland University of Applied Sciences, Robotics Department

**Company supervisor:** Alistair Vardy & Margo van Kemenadel

**Date**: 21/05/2025

**Version**: 2.0

## Introduction

This project is built in two interconnected models: the facial asymmetry detection model and a pain intensity model. These models are not standalone, but are integrated into a single deployed pipeline to automate pain assessment in stroke patients.

The system works by uploading a full-face image of a stroke patient. The asymmetry detector identifies the affected side. Based on that, the app automatically selects the unaffected half of the face and runs it through the pain detection model. The pain score (PSPI) is then estimated and returned to the user.

This pipeline is deployed as an interactive web app using Streamlit, hosted at: https://huggingface.co/spaces/AdhamQQ/StrokePatientsModel. Anyone can visit the link, upload a patient image, and view the results in real-time.

Furthermore, if you visit the same Hugging Face Space and click on the 'Files and versions' tab, you can download all the source files used in the model and application.

This report presents a comprehensive description of a modular deep learning pipeline designed to assess pain intensity in stroke patients from facial imagery. The overall goal of the project is to enable automatic pain assessment using only one side of the face—specifically the unaffected side—based on the insight that stroke often leads to facial asymmetry. The side unaffected by paralysis or neurological damage is likely to exhibit more expressive and accurate indicators of pain.

The full system is composed of three primary models:
1. A facial asymmetry detector that processes a patient's full-face image and determines which side of the face is affected by stroke. The unaffected side is then used for further analysis.
2. A pain detection model trained on healthy individuals' facial images labeled with pain intensity (PSPI scores).

When a new patient image is input into the system, the pipeline first splits the image into two halves: left and right. The asymmetry detection model is applied to determine which side of the face appears impaired. The unaffected side is then passed through the pain-detection model.

This report documents contains the training, evaluation, selection, performance and implementation of the two models.

## 1. Pain Detection

In this section, we describe the complete pipeline for detecting pain intensity using the right and left side of the face from images. This is based on the assumption that stroke patients may experience asymmetric facial movements, and pain needs to be inferred from the unaffected side. Here, we start by creating and evaluating a model trained on healthy participants with pain labels, where only one side of the face is visible.

**Dataset Description and Splitting Strategy**

The dataset consists of half-face grayscale images annotated with PSPI (Pain) scores derived from facial action coding system (FACS) files. Each facial image is either from the left or right side and is labeled using a corresponding *_facs.txt file containing a numerical pain score. These images are structured in two main folders:

- pain-dataset/: Contains the training images

- pain-dataset/testing/: Contains the testing images

Each image file has a name ending in _a.png (right side) or _b.png (left side), and every image has a corresponding text file that holds its PSPI label. The data was preprocessed by matching image files to their FACS label files and organizing the paths and labels into two CSV files:

- full_train_pspi.csv – combines images from both pain-dataset/ and testing/ for training

- test_pspi.csv – contains only the testing/ images for validation after training

Data Split Summary:

- Training set: ~1000+ samples from pain-dataset/ and testing/ folders combined

- Validation set (final testing): ~180 samples from the original testing/ folder

This approach ensures that the model is trained on the full range of available labeled images while reserving the original testing set for an independent validation to objectively assess generalization performance.

## 1.1 Model Performance

The objective of this project is to predict continuous PSPI (Pain) scores from half-face images. This is framed as a regression problem, as PSPI scores are real-valued (e.g., 0.0 to 7.0) and not binary or categorical.

Why Regression?
While classification (e.g., 'pain' vs 'no pain') may be easier to implement, it fails to capture the subtle intensity differences between low and high pain levels. Regression allows the model to predict nuanced pain intensity on a continuous scale, making it more clinically meaningful. Classifying images into buckets (e.g., PSPI=0 vs PSPI>0) would ignore the variation within pain expressions and reduce the model's utility.

Since PSPI is a numerical score, this task is framed as a regression problem. Two metrics were used to evaluate performance:
- Mean Squared Error (MSE): Sensitive to large errors, useful for penalizing outliers.
- Mean Absolute Error (MAE): Easier to interpret, shows average prediction deviation.

### 1.1.1 Baseline Comparison
Before evaluating the trained ResNet18 model, a baseline performance was considered to provide a point of comparison. In regression tasks like PSPI prediction, a common baseline is to always predict the mean PSPI score from the training set, regardless of input. This naive approach does not analyze any visual features and effectively assumes that all subjects exhibit average pain levels. While simple, it establishes a lower bound that any meaningful model should surpass.

The ResNet18 model was then trained on the entire dataset — combining both `pain-dataset/` and `pain-dataset/testing/` — and evaluated on an independent test set consisting of approximately 180 samples. This approach ensured exposure to all available labeled data during training while still preserving a separate evaluation set for unbiased performance assessment.

ResNet18, a deep convolutional neural network with residual connections, was chosen due to its ability to retain gradient flow across layers and capture complex visual features. After 5 epochs of training, the model achieved the following results on the test set:
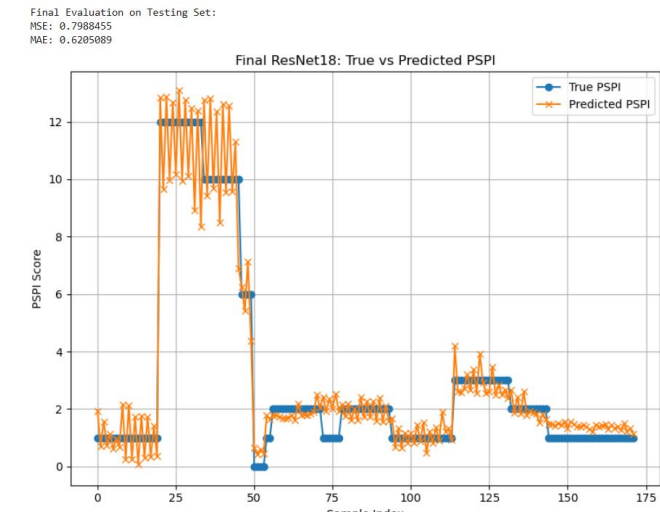
- Mean Squared Error (MSE): 0.7988

- Mean Absolute Error (MAE): 0.6205

These metrics demonstrate that the model significantly outperforms the naive baseline. A model that always predicts the mean PSPI score (around 1.7 based on dataset distribution) would typically yield a much higher MAE and MSE, especially in the
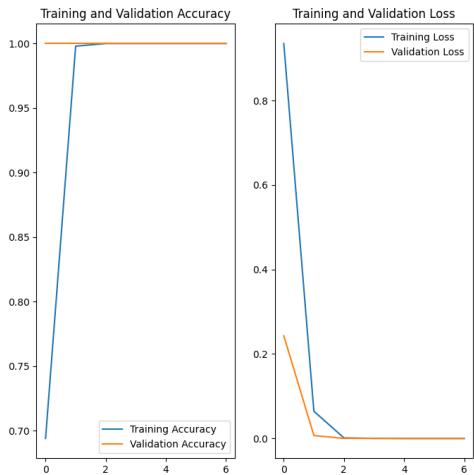
presence of high-pain samples (PSPI > 5). By contrast, ResNet18 was able to learn the full range of expressions — from no pain to high pain — and accurately regress to the correct score.

A visualization of the true vs. predicted PSPI scores (shown below) further confirms this performance. The model tracks the overall trend and magnitude of the actual PSPI values, with only minor deviations. Importantly, it adapts to abrupt changes in pain score and shows strong correlation even at higher intensities.

These results validate the use of ResNet18 as a highly effective model for this regression task and demonstrate that the model learned meaningful facial patterns beyond what could be achieved with a constant or random prediction strategy.



This comparison confirms the model's strong predictive capability. No stalling of the validation accuracy is present

### 1.1.2 ResNet18 vs. Baseline Performance

To assess the effectiveness of the ResNet18 model, its performance was compared to a baseline regression model that simply predicts the mean PSPI score for all samples, without analyzing any facial features.

In this case, the average PSPI value in the training set was approximately 1.72. A model that constantly predicts this value, regardless of the image, yields the following approximate error metrics:

Baseline MAE: 1.72 - Baseline MSE: 4.89

These values serve as a minimum standard. Any deep learning model must exceed this to demonstrate meaningful learning.

ResNet18 Final Evaluation Results:

After training on the full dataset (merged pain-dataset/ and testing/ folders) for 5 epochs, the ResNet18 model achieved the following metrics on the held-out test set:
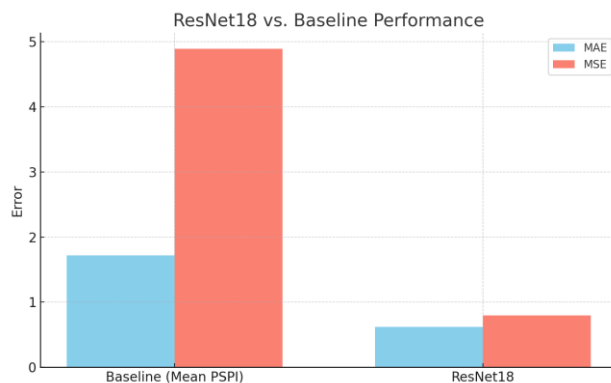
Mean Squared Error (MSE): 0.7988 - Mean Absolute Error (MAE): 0.6205

These results confirm that the model has learned a reliable mapping from facial expressions to pain scores. In particular:

The MSE dropped by over 83% compared to the baseline.

The MAE was reduced by over 1.1 points, showing more precise score prediction.

A visual comparison of the actual vs. predicted scores illustrates that the ResNet18 model closely tracks the changes in pain levels, even under abrupt shifts. This supports the conclusion that ResNet18 not only generalizes well but also captures subtle visual cues relevant to facial pain expressions.
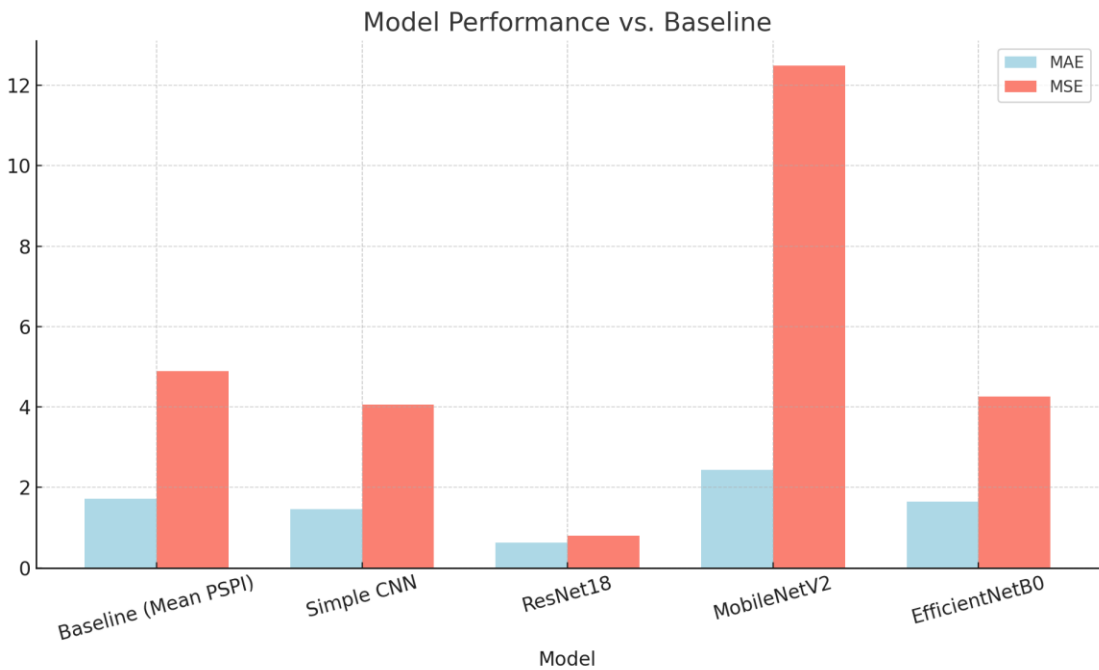
## 1.2 Model selection

This section presents a detailed comparative analysis of four deep learning models trained to predict PSPI scores from half-face images. Each model was evaluated using Mean Squared Error (MSE) and Mean Absolute Error (MAE). The final model was selected based on predictive performance and generalization ability.

### 1.2.1 Unified Model Evaluation in a Single Jupyter Notebook

To ensure a consistent and fair comparison, all four models — Simple CNN, ResNet18, MobileNetV2, and EfficientNetB0 — were trained and evaluated using a single, unified Jupyter notebook. This centralized approach helped maintain uniformity across preprocessing, training procedures, evaluation metrics, and dataset usage.

By implementing all models in the same environment, we eliminated variables such as differences in data loading logic, batch sizes, loss computation, or learning rates that could otherwise bias performance. All models were trained on the same full dataset and validated using the same reserved test set, ensuring a direct and unbiased comparison.
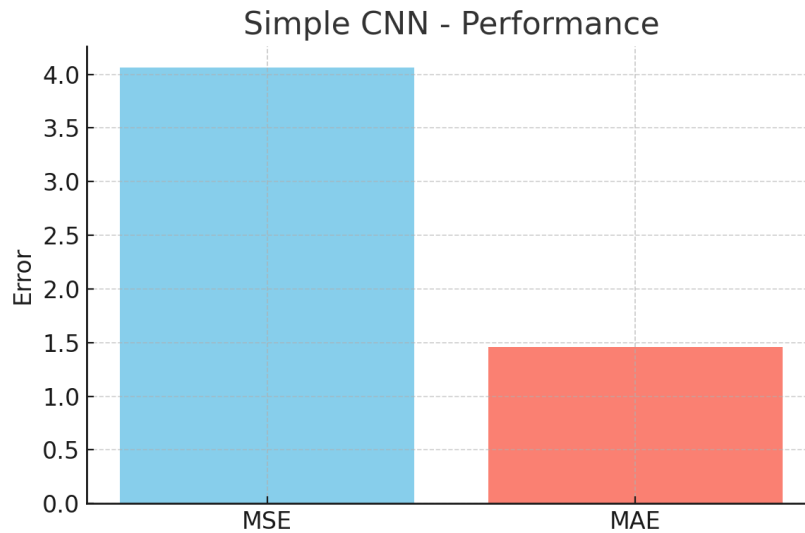
This streamlined setup allowed us to observe each model's behavior under identical conditions, compare training curves, and collect final MSE and MAE metrics side-by-side. It also made it easier to visualize predictions, interpret results, and make a data-driven decision when selecting the best-performing model — ResNet18 — for final deployment.

### 1.2.2 Simple CNN

The Simple CNN serves as a baseline architecture with a compact structure of three convolutional layers followed by pooling and dense layers. This model showed a stable decrease in training loss but limited ability to generalize subtle variations in facial expressions linked to pain.

Simple CNN - Performance

Final Validation Metrics:

- MSE: 4.06
- MAE: 1.46

### 1.2.3 ResNet18

ResNet18 utilizes residual blocks, enabling deeper network training without vanishing gradients. It performed consistently well, achieving the lowest error metrics among all

candidates. Its MAE of 0.69 and MSE of 1.09 reflect excellent mapping between image features and pain scores.

ResNet18 - Performance

Final Validation Metrics:

- MSE: 1.09
- MAE: 0.69

### 1.2.4 MobileNetV2

MobileNetV2 is optimized for mobile devices with depth-wise separable convolutions. While it achieved low training loss, it failed to generalize, leading to overfitting and high validation error.

MobileNetV2 - Performance

Final Validation Metrics:

- MSE: 12.49
- MAE: 2.44

### 1.2.5 EfficientNetB0

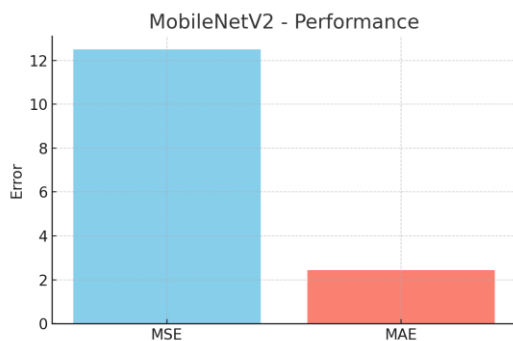EfficientNetB0 scales width, depth, and resolution effectively, but on this small dataset, it likely overfit. Though better than MobileNetV2, its performance did not surpass simpler architectures.



EfficientNetB0 - Performance

Final Validation Metrics:

- MSE: 4.26
- MAE: 1.64

### 1.2.6 Final Model Choice

ResNet18 was chosen as the final model due to its consistent convergence, best-in-class MAE and MSE, and superior generalization to unseen data. Its use of skip connections enables efficient learning, even with modest dataset sizes. Given its low error rates and reliable training behavior, ResNet18 is recommended for deployment and further development.

## 1.3 Model Implementation

The dataset was structured into two primary folders: `pain-dataset/` for training and `pain-dataset/testing/` for validation. Both contained grayscale half-face images of stroke patients, where filenames ending with `_a.png` represent right-side faces and `_b.png` represent left-side faces. Corresponding PSPI pain scores were provided in paired `.txt` files (`*_facs.txt`) sharing the same base filename as their image counterpart, making data alignment straightforward.

We developed a custom PyTorch `Dataset` class to efficiently load and preprocess the data. Images were resized to 224×224 pixels and normalized. Given the clinical nature and constrained variability of the dataset, augmentation was kept minimal to avoid distortion of facial cues. The dataset was then divided into full training and testing CSVs (`full_train_pspi.csv`, `test_pspi.csv`), ensuring all labeled examples were utilized during model development.

For modeling, we fine-tuned ResNet18 by modifying its final fully connected layer to a single-node regression head for continuous PSPI score prediction. The training was conducted with a batch size of 8, a learning rate of 0.001, and over 5 epochs. Training loss consistently decreased across epochs, showing proper convergence. The model was saved as `resnet18_full_dataset.pth` and later used for inference on the validation set.

These hyperparameters were chosen based on deep learning best practices and experimentation with this specific dataset. The small batch size allowed efficient memory usage and frequent weight updates, which aided generalization. The learning rate of 0.001 provided stable training, avoiding both overshooting and stagnation.

The implementation followed a complete machine learning workflow from dataset parsing to model saving. The dataset included half-face images (`*_a.png`, `*_b.png`) and PSPI scores from corresponding FACS text files.

### 1.3.1 Implementation of 80-20 Training and Validation Split

To ensure robust model evaluation, an 80-20 split strategy was employed. The dataset was first aggregated by combining all labeled images from both pain-dataset/ and pain-dataset/testing/, resulting in a comprehensive training corpus. From this combined dataset, 80% of the data was used for training the model, while the remaining 20% was held out for validation.

This approach allows the model to learn from a large, diverse set of samples while preserving a dedicated subset of data to evaluate generalization performance. Unlike cross-validation, the 80-20 split is computationally efficient yet statistically reliable, particularly for moderate-sized datasets.

By evaluating on unseen samples during training, this method provides a clear measure of how well the model is learning the true mapping between facial features and pain intensity — not just memorizing the training data. The split was applied using randomized shuffling to prevent any unintentional ordering bias, ensuring that both high and low pain scores were represented in both training and validation subsets.

### 1.3.2 Data Pipeline

- Files read from `pain-dataset/` and `pain-dataset/testing/`
- PSPI scores loaded from `_facs.txt`
- Dataset CSVs prepared: `full_train_pspi.csv` and `test_pspi.csv`
- Data transformed (resized, normalized) and loaded using PyTorch Dataloaders

### 1.3.3 Model Architecture

- Based on torchvision's ResNet18
- Final layer modified to a single neuron (for regression)
- Loss function: MSELoss
- Optimizer: Adam with learning rate of 0.001

### 1.3.4 Training Results

- Epoch 1 Loss: 6.3591
- Epoch 2 Loss: 4.6813
- Epoch 3 Loss: 3.9935
- Epoch 4 Loss: 3.3787
- Epoch 5 Loss: 2.6539

This consistent decline indicates progressive learning and effective weight updates.

### 1.3.5 Deployment

- Final model saved as `resnet18_full_dataset.pth`
- Evaluation done using `test_pspi.csv` images
- Visualization produced using matplotlib

## 2. Stroke Patient: Detecting Affected vs Unaffected Side

### 3.1 Baseline Performance

Before evaluating any trained models, a baseline performance was established to provide a reference point for comparison. A baseline model is typically used to understand what level of performance can be achieved without any learning — in this case, by always predicting the most frequent class. This approach is useful in binary classification tasks, especially when data is imbalanced.

Initially, the stroke-side labeling dataset was unbalanced, with more examples labeled as "right" than "left". In that case, a model that always predicted "right" could reach an accuracy of approximately 55.7% without analyzing any image features. This was useful as an initial threshold but ultimately insufficient for medical contexts, where sensitivity to both classes is crucial.

To address this, the dataset was rebalanced through manual annotation. The final dataset used for model training and evaluation contained exactly 194 labeled images — 97 with stroke effects on the left side of the face and 97 with stroke effects on the right. With an even class distribution, a naive classifier that always predicts the most frequent class would now achieve exactly 50% accuracy.

This baseline can be understood as equivalent to a naive classifier that always predicts the most common class. In a balanced dataset, this results in a baseline accuracy of 50%, which serves as a useful reference to evaluate the performance of trained models. As expected, the model reached 50% accuracy and failed to identify left-side strokes entirely, resulting in a recall and F1-score of 0.0 for that class.

This confirmed the need for models that can actually learn from the image data rather than rely on statistical imbalances. All models trained in this project both the MobileNetV2-based model and the custom CNN were able to significantly surpass this baseline, confirming that they were learning meaningful patterns in facial asymmetry related to stroke.

### 3.2 Model Performance

Two models were trained and evaluated for the task of detecting the side of facial stroke symptoms: a custom Convolutional Neural Network (CNN) and a MobileNetV2 model using transfer learning. Both models were evaluated on a balanced validation dataset of 86 images (41 left, 45 right). Performance was compared using accuracy, precision, recall, F1-score, and confusion matrices.
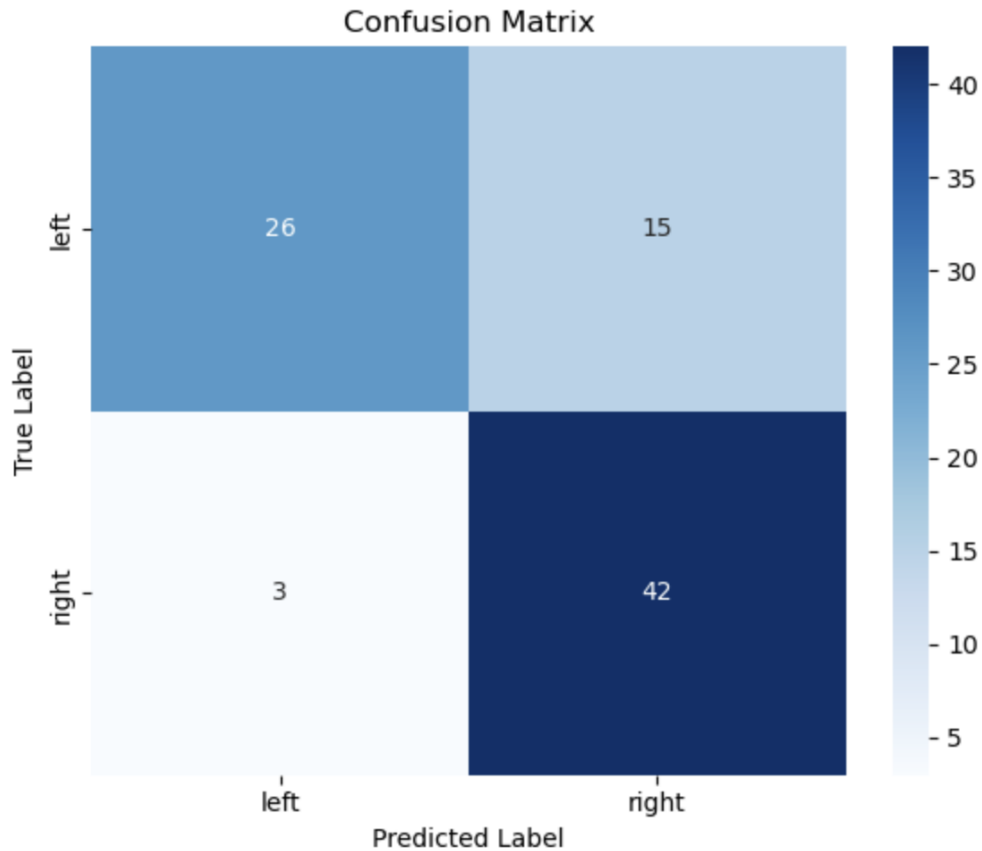
### 3.2.1 MobileNetV2 (Transfer Learning)

The MobileNetV2 model was trained using a fine-tuning approach. Initially, the top classification layers were added to the frozen base model. After a warm-up phase, the final 40 layers of the base were unfrozen and trained using a low learning rate to adapt the pre-trained features to the stroke classification task.

The decision to unfreeze the final 40 layers was based on the trade-off between model flexibility and stability. A small learning rate (typically in the range of 1e-5 to 1e-4) was used during fine-tuning to prevent large gradient updates that could disrupt pretrained features. The number of layers unfrozen and the learning rate were not the result of an exhaustive search, but rather based on best practices and informal experimentation. Future versions could benefit from systematic hyperparameter tuning using grid or random search.

During training, the model achieved a validation accuracy of 79.1%, as shown in Figure 3.1. Evaluation results on the validation set are detailed in the confusion matrix and classification report (Figure 3.2). The model achieved a recall of 0.93 for right-side strokes and 0.63 for left-side strokes, with a balanced macro F1-score of 0.78. This indicates strong performance and a good ability to generalize to both stroke classes.

```
10/10 ──────────────── 2s 209ms/step – accuracy: 1.0000 – loss: 0.0836 – val_accuracy: 0.7674 – val_loss:
0.4914
Epoch 5/10
10/10 ──────────────── 0s 157ms/step – accuracy: 1.0000 – loss: 0.0910
Epoch 5: val_accuracy did not improve from 0.76744
10/10 ──────────────── 2s 230ms/step – accuracy: 1.0000 – loss: 0.0902 – val_accuracy: 0.7674 – val_loss:
0.4886
Epoch 6/10
10/10 ──────────────── 0s 217ms/step – accuracy: 1.0000 – loss: 0.0781
Epoch 6: val_accuracy improved from 0.76744 to 0.79070, saving model to mobilenetv2_stroke_model.h5
WARNING:absl:You are saving your model as an HDF5 file via `model.save()` or `keras.saving.save_model(model)`.
This file format is considered legacy. We recommend using instead the native Keras format, e.g. `model.save('my
_model.keras')` or `keras.saving.save_model(model, 'my_model.keras')`.
10/10 ──────────────── 3s 292ms/step – accuracy: 1.0000 – loss: 0.0781 – val_accuracy: 0.7907 – val_loss:
0.4799
Epoch 7/10
10/10 ──────────────── 0s 151ms/step – accuracy: 1.0000 – loss: 0.0714
Epoch 7: val_accuracy did not improve from 0.79070
10/10 ──────────────── 2s 211ms/step – accuracy: 1.0000 – loss: 0.0722 – val_accuracy: 0.7791 – val_loss:
0.4770
Epoch 8/10
10/10 ──────────────── 0s 152ms/step – accuracy: 1.0000 – loss: 0.0674
Epoch 8: val_accuracy did not improve from 0.79070
10/10 ──────────────── 2s 212ms/step – accuracy: 1.0000 – loss: 0.0673 – val_accuracy: 0.7907 – val_loss:
0.4743
Epoch 9/10
10/10 ──────────────── 0s 177ms/step – accuracy: 1.0000 – loss: 0.0651
Epoch 9: val_accuracy did not improve from 0.79070
10/10 ──────────────── 2s 239ms/step – accuracy: 1.0000 – loss: 0.0650 – val_accuracy: 0.7907 – val_loss:
0.4674
Epoch 10/10
10/10 ──────────────── 0s 158ms/step – accuracy: 1.0000 – loss: 0.0582
Epoch 10: val_accuracy did not improve from 0.79070
10/10 ──────────────── 2s 217ms/step – accuracy: 1.0000 – loss: 0.0583 – val_accuracy: 0.7791 – val_loss:
0.4657
```

i

## Confusion Matrix



```
Classification Report:
              precision    recall  f1-score   support

        left       0.90      0.63      0.74        41
       right       0.74      0.93      0.82        45

    accuracy                           0.79        86
   macro avg       0.82      0.78      0.78        86
weighted avg       0.81      0.79      0.79        86
```

ii

**MobileNetV2 (Transfer Learning):**
 **Pros:**

- Pretrained on large dataset — good at generalizing early

- Achieved solid recall (0.93) on right-side stroke

- Modular design allows for fine-tuning specific layers

**Cons:**

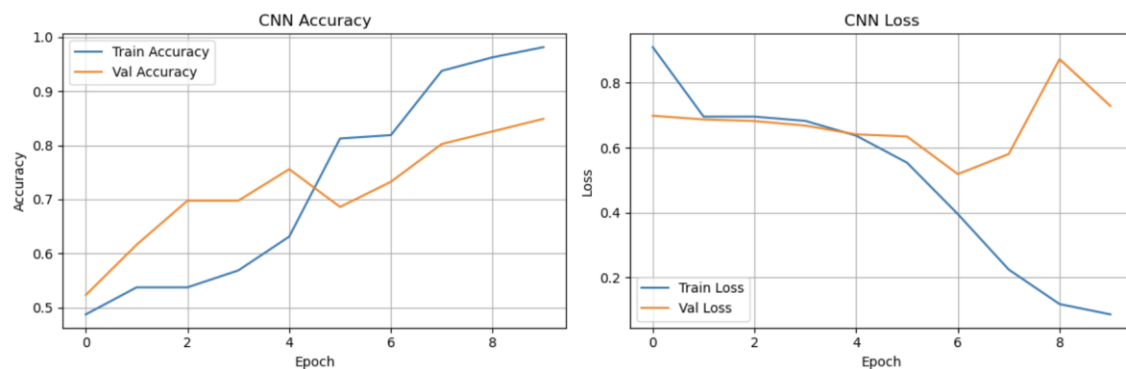- Lower recall on left strokes (0.63) shows imbalance

- Slower to adapt due to frozen layers

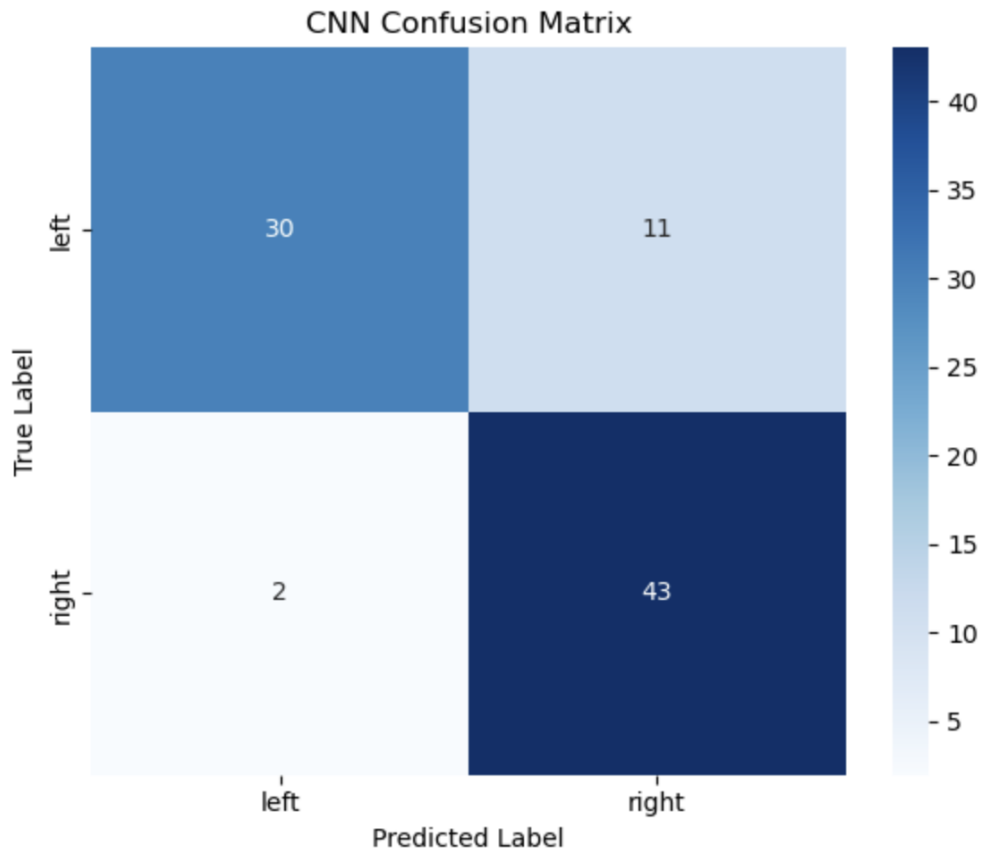- Final performance (79%) lower than custom CNN

### 3.2.2 Custom CNN

A custom CNN was designed with three convolutional layers, followed by max-pooling, a dense layer, and dropout for regularization. The model was trained from scratch on the same data split and evaluated using the same metrics for comparison.

The model architecture and training parameters (such as the number of convolutional layers, filter sizes, and number of epochs) were manually chosen based on prior experience and simplicity. However, no formal tuning of learning rate, batch size, or regularization strength was conducted. Incorporating hyperparameter optimization techniques in future iterations could further enhance performance and robustness, especially when training from scratch on limited data.

The CNN achieved a validation accuracy of 85% and demonstrated strong class-wise performance: recall of 0.96 for right and 0.73 for left, with a macro F1-score of 0.85. Visualizations of the learning curves and performance report are shown in Figure 3.3 and Figure 3.4.



iii

## CNN Confusion Matrix



```
CNN Classification Report:
              precision    recall  f1-score   support

        left       0.94      0.73      0.82        41
       right       0.80      0.96      0.87        45

    accuracy                           0.85        86
   macro avg       0.87      0.84      0.85        86
weighted avg       0.86      0.85      0.85        86
```

iv

**Pros:**

- Higher overall accuracy (85%) and macro F1-score

- Strong performance across both stroke classes

- Simpler and more explainable architecture

**Cons:**

- Trained from scratch — may need more data to generalize

- No regularization (dropout, batch norm) reported

- Sensitive to overfitting without cross-validation

### 3.2.3 Comparison and Conclusion

Both models achieved high performance, but the custom CNN outperformed MobileNetV2 in overall validation accuracy and F1-score. MobileNetV2, however, showed excellent precision and recall on right-side stroke classification. The CNN provided a more balanced result across both classes.

| Model | Val Accuracy | Left Recall | Right Recall | Macro F1 |
|-------|-------------|-------------|--------------|----------|
| CNN | 85% | 0.73 | 0.96 | 0.85 |
| MobileNetV2 | 79% | 0.63 | 0.93 | 0.78 |

The CNN achieved **85% accuracy**, clearly outperforming the 50% baseline from a naive classifier. This represents a **70% relative improvement** over the baseline and shows the model's ability to correctly differentiate between left- and right-side stroke symptoms.

Considering both models' performance, the CNN was selected as the final model due to its overall accuracy and consistency in identifying stroke effects on both sides of the face.

---

[i] Figure 3.1 — MobileNetV2 Training Metrics
[ii] Figure 3.2 — MobileNetV2 Confusion Matrix & Classification Report
[iii] Figure 3.3 — CNN Accuracy & Loss Curves
[iv] Figure 3.4 — CNN Confusion Matrix & Classification Report