

DATASETS REPORT

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: 23/04/2025

Version: 2.0

Table of contents

1. Dataset Analysis	2
1.1 Content of the Dataset.....	2
1.2 Dataset Format and Collection Method	4
1.3 Quantitative and Qualitative Analysis.....	4
Emotion Pain Dataset:.....	5
Left-side pain dataset.....	5
Right-side-dataset.....	6
Acute Stroke Dataset:	7
2. Dataset Visualization.....	8
Sample image plotting	9
Histogram analysis.....	10
Facial landmark mapping.....	11
3. Cleaning and Preprocessing of Dataset.....	16
3.1 Image Processing Steps	16
3.2 Preprocessing Considerations.....	17
Conclusion.....	17

1. Dataset Analysis

1.1 Content of the Dataset

This project is built upon two core datasets that serve distinct but complementary purposes:

1.1.1. Emotion Pain Dataset (Healthy Individuals)

The Emotion Pain Dataset consists of facial images of healthy individuals intentionally subjected to painful stimuli (e.g., limb movements, light exposure) to evoke natural pain-related expressions. Each image is labeled using Prkachin and Solomon Pain Intensity (PSPI) scores derived from Facial Action Units (AUs), which indicate the intensity and authenticity of expressed pain.

To align with the project's clinical objective—detecting pain from one half of the face—the dataset was manually preprocessed and split into two sub-datasets:

- Right-Side Pain Dataset: Contains only the right half of each image, extracted using facial landmark detection and vertical splitting.
- Left-Side Pain Dataset: Contains the left half of the same images, processed similarly.

Each half-face is paired with its corresponding .txt file that includes facial AUs and the computed PSPI score. These two subsets were used to train separate models for detecting pain when the stroke leaves either the left or the right side expressive. This design allows flexible deployment depending on which side of the face is unaffected.

1.1.2. Acute Stroke Dataset

The Acute Stroke Dataset contains high-quality facial images of patients clinically diagnosed with acute stroke. The dataset was used to train a facial asymmetry detection model that identifies which side of a patient's face is affected by stroke-related impairments.

This model plays a crucial role in the pipeline. It processes a full-face image and determines whether the left or right side is impaired. Based on this prediction, the system automatically selects the unaffected side and passes it through the appropriate pain-detection model (left or right), trained on the corresponding subset of the Emotion Pain Dataset.

Integration

By combining the Emotion Pain Dataset (split by side) and the Stroke Patient Dataset, the system is able to:

- Automatically detect stroke-affected facial asymmetry
- Select the unaffected side of the patient's face
- Accurately estimate the PSPI pain score using deep learning on a healthy-expression benchmark

1.2 Dataset Format and Collection Method

The datasets include both color and grayscale images in widely used formats such as JPEG and PNG, ensuring compatibility with standard image processing pipelines.

Images in the Emotion Pain dataset were collected through controlled laboratory experiments, where healthy individuals were asked to perform specific body movements or were exposed to light-based stimuli to evoke facial expressions of pain. These experiments were designed to capture a broad range of pain intensities while maintaining uniformity in background and lighting wherever possible.

In contrast, the Stroke Patient dataset consists of images taken from clinical or real-world healthcare settings, capturing a more natural and diverse representation of facial expressions. These images often include minor variations in head pose, lighting, and background, reflecting the unpredictability of real patient assessments.

All images from the Emotion Pain dataset were further split into left and right facial halves using automated landmark detection and vertical cropping. This allowed us to build two distinct datasets - one for the left side of the face and one for the right side - each labeled with PSPI scores indicating the level of perceived pain.

Facial Action Unit information (from .txt files) was collected alongside each image to ensure pain annotations were precise and machine-readable. These files are essential for training and evaluating the pain detection models.

1.3 Quantitative and Qualitative Analysis

Features:

Each image consists of either full-face or half-face portraits. In the Emotion Pain dataset, subjects express clear facial reactions to controlled pain stimuli. For the Stroke dataset, natural expressions are captured from real patients in clinical conditions.

Classes:

The data is grouped into two primary categories:

1. Healthy individuals with pain-related facial expressions (annotated with PSPI scores)
2. Stroke patients, for whom the goal is to identify the affected side and analyze pain on the unaffected side

Data Types:

- Pixel-based image data (PNG, JPEG)

- Structured annotation files (.txt) containing Facial Action Units and calculated PSPI values (for healthy pain datasets)

Range of Values:

- PSPI scores vary from neutral (0) to intense pain (7+), ensuring a representative spectrum
- Images include a wide range of facial orientations, lighting conditions, and skin tones, aiding generalization

Each database selected will have relevance to achieving the result. At the start it was aimed to work correctly to process the data and big data batches were avoided. Since however the initial dataset was not diverse enough, it was expanded.

Emotion Pain Dataset:

The emotion pain dataset will validate the result obtained from the stroke dataset. The dataset had to be analyzed by splitting it into left side affectation poses and right-side affectation poses. The sections also share common characteristics irrespective of the database chosen. Link to the [Dataset](#).

Left-side pain dataset

The database with people experiencing pain on the left side of the body shows people requested to do movement with parts of the body towards left (moving eyes and limbs). It is proved that performing these causes them to have painful face expressions because some of the images have associated pain-intensity scores which indicate pain.

The left-side pain dataset contains a sample of images and facial action unit files indicating the Prkachin and Solomon Intensity score for pain (PSPI).

The dataset was downloaded from a recommended online database, and it contains a folder with more subjects tested left-side for pain reactions. The size of the folder goes around 600 PSPI files and their associated images. The dataset is relevant because it displays certain type of stimuli response which is common for stroke patients (increased sensitivity to light and movement).

The images from the dataset were chosen based on the criteria that the person would show pain on the left side of the face expression. For the criteria to match the real affected side, the principle that the affected side is the inverse side of the face was applied from the viewer's perspective.

-Left-side affected by pain

The part of the dataset with a PSPI score greater than 0 and which thus indicates pain, includes the left side of each image after it was split in left and right. These images are part of the directory pspi (fig. 1).

-Left side not affected by pain

The part of the dataset with a PSPI score equal to 0 and which thus does not indicate pain, includes the right side of each image after it was split in left and right. These images are part of the directory no-pspi (fig. 1).

The left-side model is to be trained with left-side no pain images and left-side pain images. The training technique involves creating the image dataset from the directory training, with the labels inferred. These labels represent the classes pspi and no-pspi.

```
training/
├── right-side-pain/
│   └── right-side-pspi/
├── left-side-pain/
│   ├── training/
│   │   ├── pspi/
│   │   └── no-pspi/
│   ├── imagetoclassify/
│   └── left-side-pspi/
├── labels.csv
├── right-side-pain.csv
├── right-side-pain-vs_no_pain.csv
├── left-side-pain-vs_no_pain.csv
├── left-side-pain.csv
└── stroke_data.zip
```

Figure 1. Structure of left-side dataset directory

The dataset (fig. 2) is a good set to use because of the following reasons:

- it includes images of subjects of both male and female, of varied ages including younger aged adults and middle-aged adults. There is a younger aged male adult and a younger aged female adult. There is also a middle-aged male and a middle-aged female.
- it includes face expressions with values of the pain intensity score which support the recognition of pain using Prkachin and Solomon Intensity score values in the range 0-16
- it includes face expressions corresponding to various severity levels, which can help categorize correctly the pain experienced by a patient

The aspect of pain intensity score diversity was concretized in a table which contains information about the dataset, images, test method, pain severity and PSPI score. The dataset information is available [here](#).

Right-side-dataset

The right-side pain dataset was built from the same original database as the left-side dataset. However, instead of focusing on the left side of the subjects' expressions, we extracted the right half of the face from each image for this part of the project.

This dataset features individuals performing body movements and gestures aimed at the right side (e.g., turning eyes and limbs to the left). From a neurological perspective, these stimuli indirectly affect the right side of the face, often resulting in pain-related expressions. This is supported by PSPI scores (Prkachin and Solomon Pain Intensity), which indicate visible discomfort in some of the recorded expressions.

The right-side dataset contains a collection of facial images and their corresponding facial action unit text files. The images were vertically split using facial landmarks, and the right half of each face was extracted and stored. The folder structure and size match the left-side set, comprising around 600 images and 600 PSPI .txt files.

The selection criteria focused on ensuring that visible expressions of pain were clearly represented on the right half of the face, maintaining visual quality and consistent lighting. From a clinical standpoint, this mirrors the situation where a stroke affects the left hemisphere of the brain, and the right side of the face remains expressive and capable of displaying pain reactions.

This dataset was critical for training the right-side pain detection model, which is used in the overall system when the left side is detected as the stroke-affected side.

Right-side affected by pain

The part of the dataset with a PSPI score greater than 0 — indicating visible pain — includes the right side of each image after it was split into left and right halves. These images are organized under the directory `pspi`.

Right-side not affected by pain

The portion of the dataset with a PSPI score equal to 0 — indicating a neutral or no-pain expression — includes the right side of each image after the same splitting process. These images are stored in the directory `no-pspi`.

The right-side model is trained using both right-side pain images and no-pain images. The training approach involves constructing the dataset from the training directory, where the class labels (`pspi` and `no-pspi`) are inferred from the folder structure. These labels serve as targets for classification or regression, depending on the model's configuration.

Acute Stroke Dataset:

The Acute Stroke Dataset consists of 1,259 full-face images of individuals clinically diagnosed with stroke. This dataset was provided as part of the official assignment resources (`opdracht`) and plays a central role in developing and testing the facial asymmetry detection component of the system.

This dataset is particularly valuable because it offers real-world clinical data - images taken from stroke patients in varying conditions, as opposed to lab-controlled or synthetic datasets. Each image captures natural facial expressions and includes

diversity in lighting, head orientation, age, gender, and skin tone, making it an ideal testbed for evaluating model robustness.

Although not curated for facial symmetry research specifically, the dataset is diverse and sufficiently large to support effective training and validation of machine learning models that aim to detect stroke-affected facial asymmetry. It enables the development of a model that can generalize across different patients, improving the practical relevance and clinical applicability of the proposed pain assessment system.

2. Dataset Visualization

Visualization plays a key role in understanding dataset distribution and facial expression characteristics. We employed:

- Histogram analysis of pixel intensity values to examine contrast variations.
- Sample image plotting to showcase differences between healthy pain-expressed faces and stroke patient faces.
- Facial landmark mapping to visualize key feature points used for splitting images.

Emotion Pain Dataset:

Sample image plotting

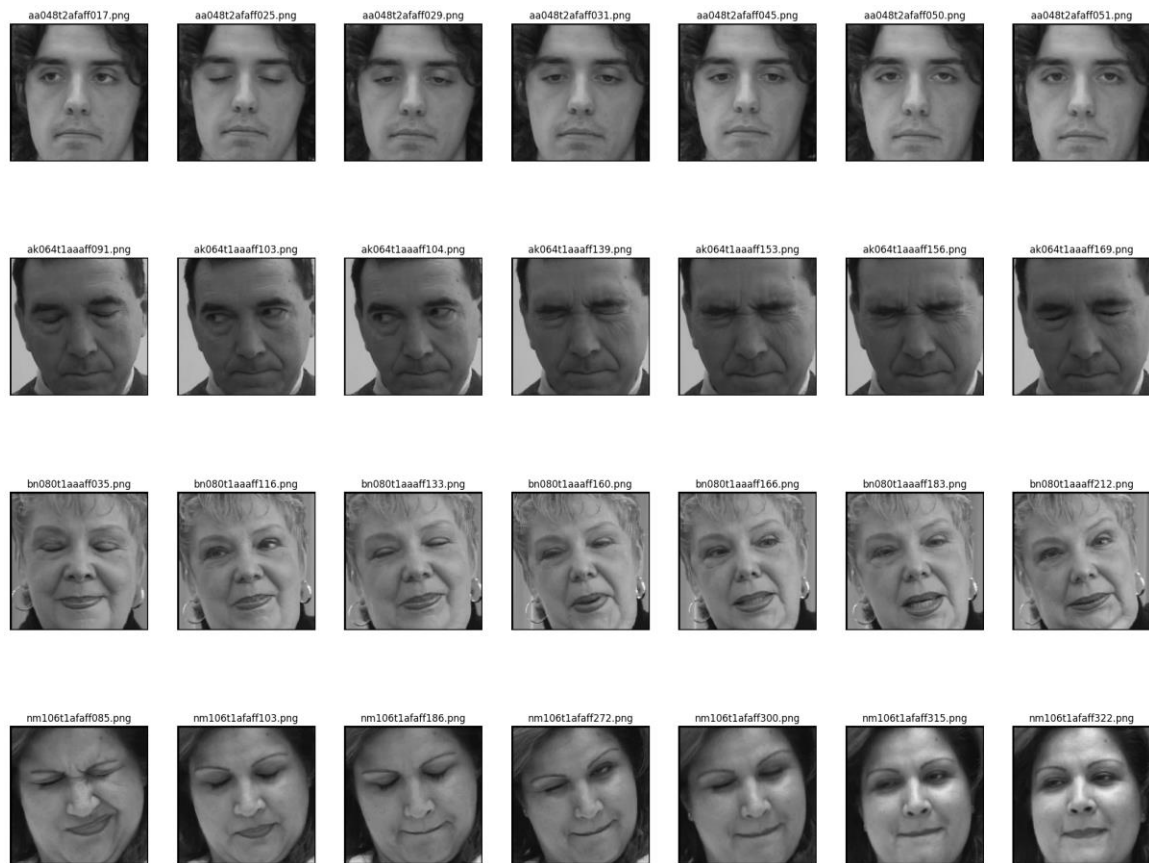


Figure 2. Sample of left and right pain dataset

Histogram analysis

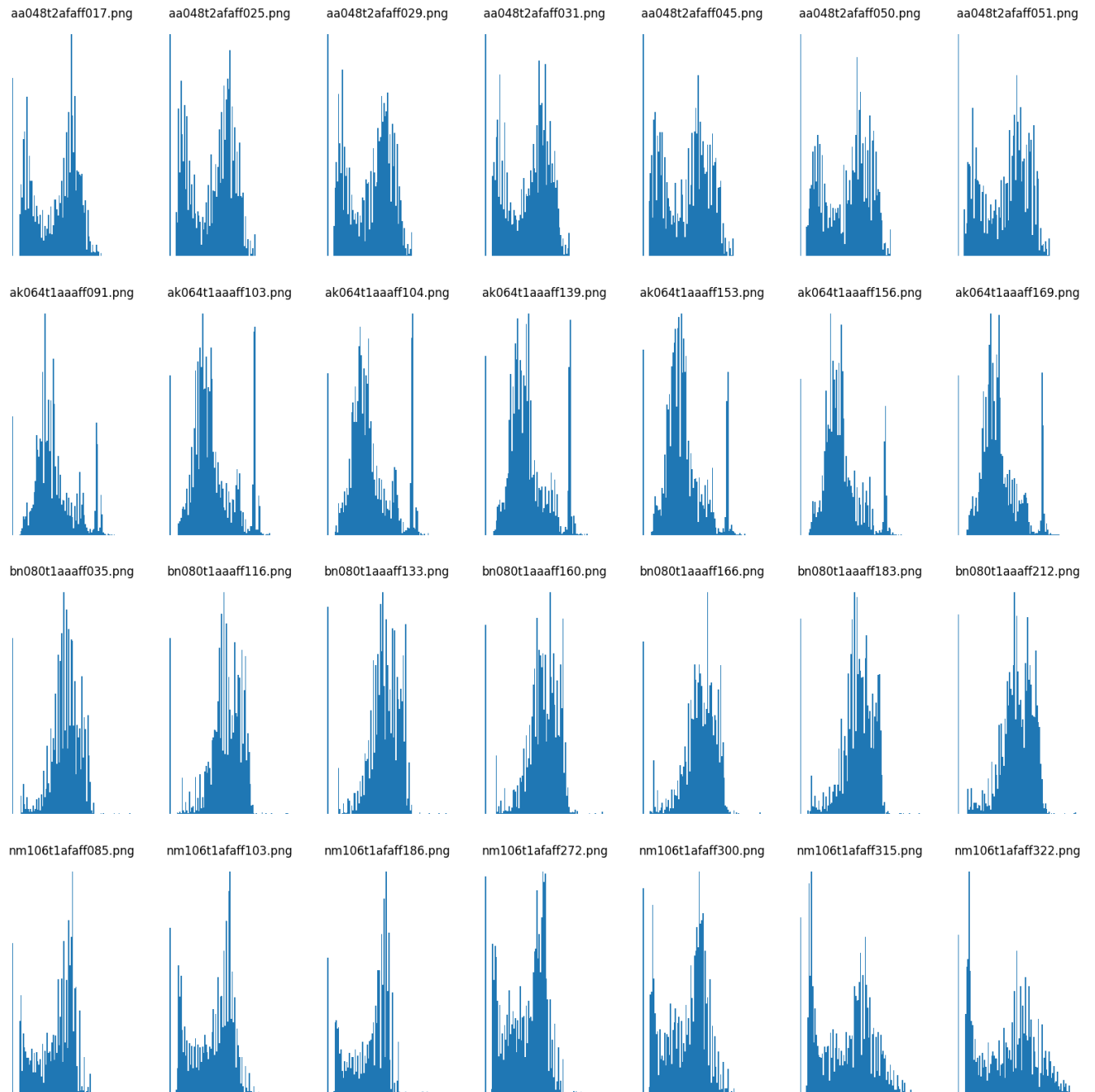


Figure 3. Histogram analysis of sample pain dataset

The histogram analysis shows that the pixel intensities are not stretched across the entire range, which is 0 to 255. Thus that the image contrast is not optimal. The images could be enhanced with techniques such as Contrast Limited Adaptive Histogram Equalization.

Facial landmark mapping

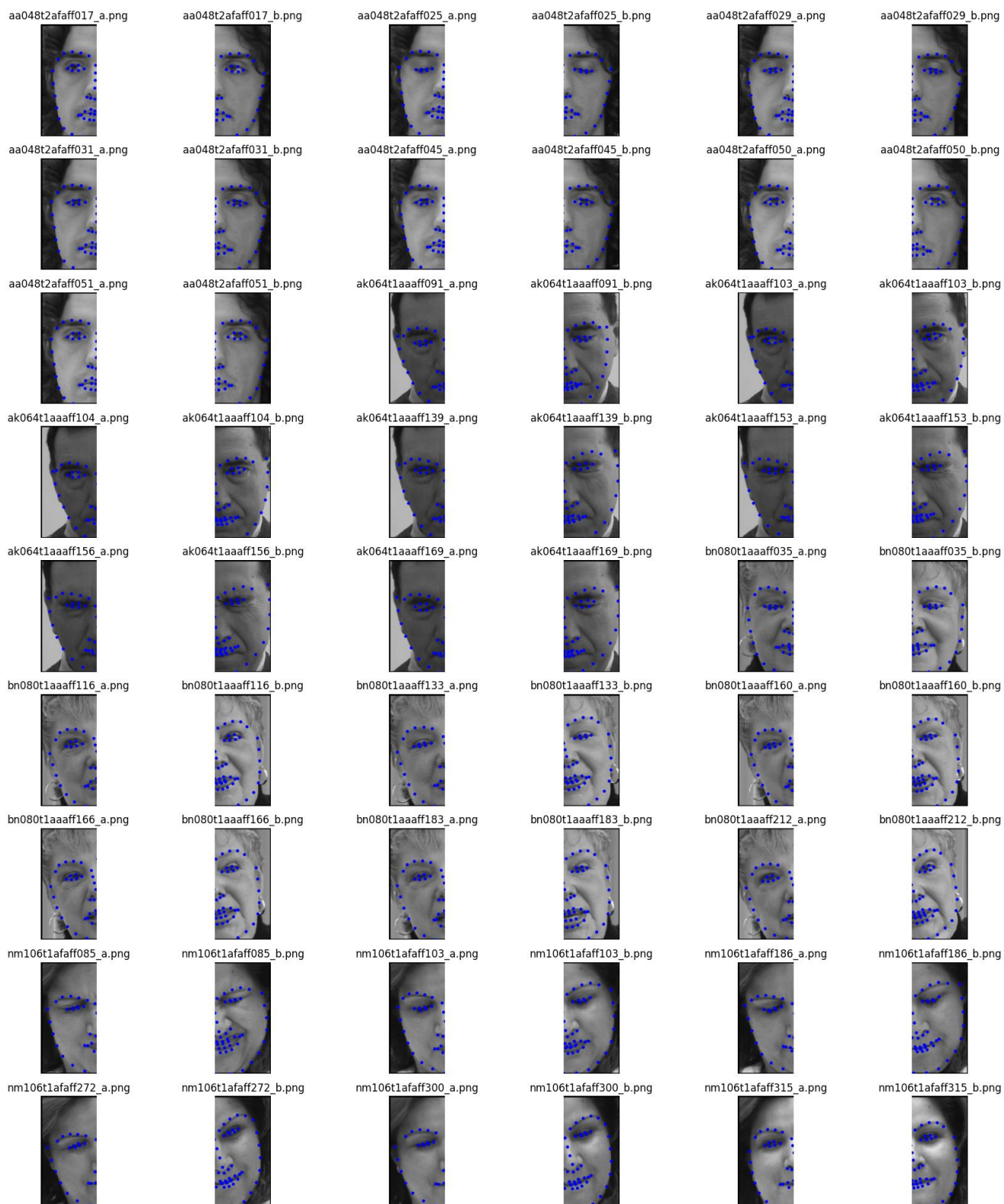


Figure 4. Facial landmarks mapping which shows symmetric split

The dataset visualisations strategies which were useful include:

- Dataset visualisations to determine how well the datasets are balanced and how much work will be required to achieve balanced datasets

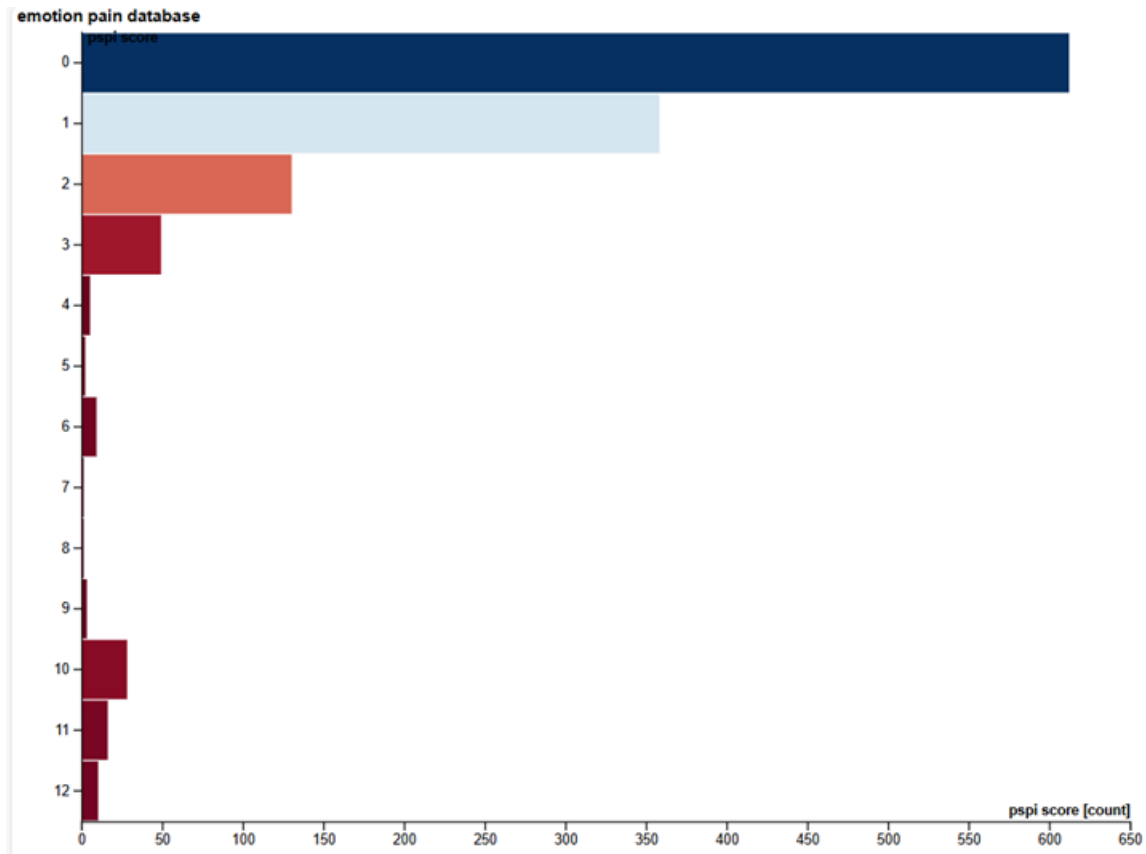


Figure 5. PSPI bar chart

The PSPI scores diversity is shown in the PSPI bar chart (fig. 5). It shows the PSPI scores and related amounts. Each bar represents a category and the width is proportional to the quantitative dimension.

The PSPI bar chart shows a diverse range of values with the potential to support a model for pain detection.

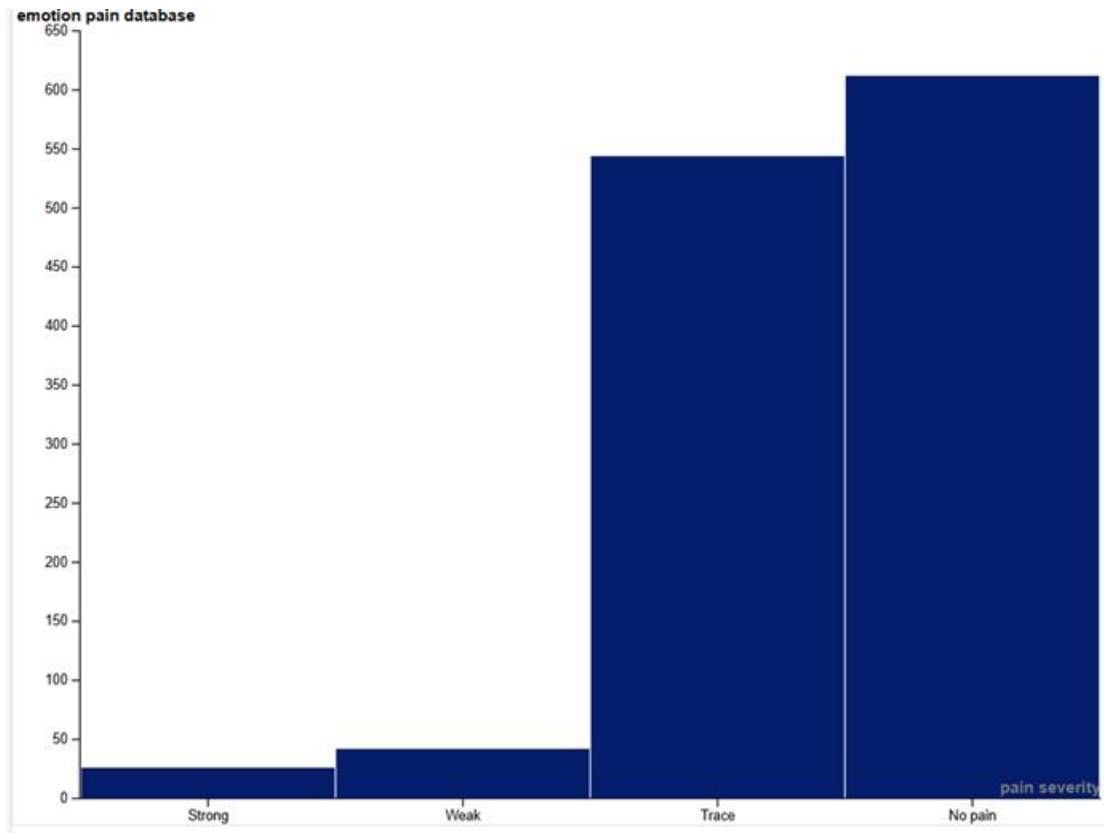


Figure 6 . Pain severity bar chart

The pain severity bar chart (fig. 6) shows the categories of pain severity. It displays the quantitative dimensions that are related to categories.

Thus it is shown that the pain dataset can be used for pain detection, with the capability to distinguish precisely various pain severity levels.

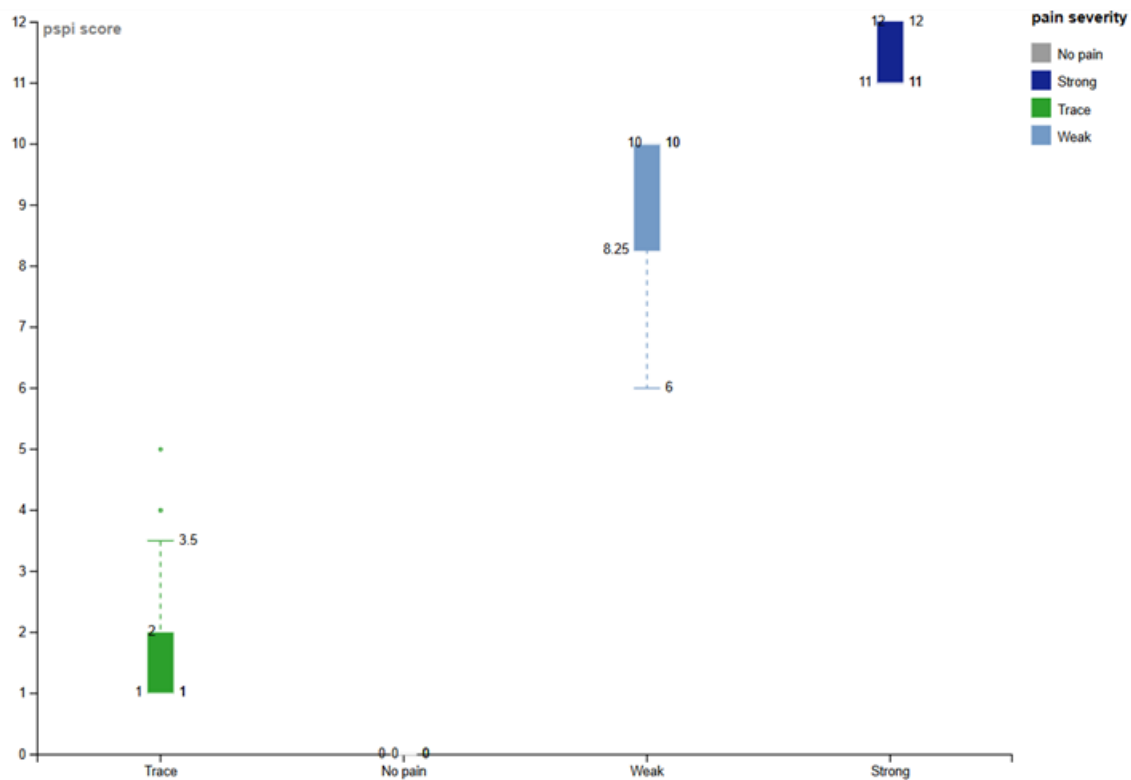


Figure 7. Pain severity box plot

The pain severity box plot shows that the PSPI scores are mostly centered around the mean of each severity category.

One problem of the dataset is about outliers. It was noticed that the dataset contains extreme values which will affect the performance of the models. One situation is the neutral expression set of values found across the dataset. Due to this fact, only some of the images, after splitting, were retained to build the no pain dataset.

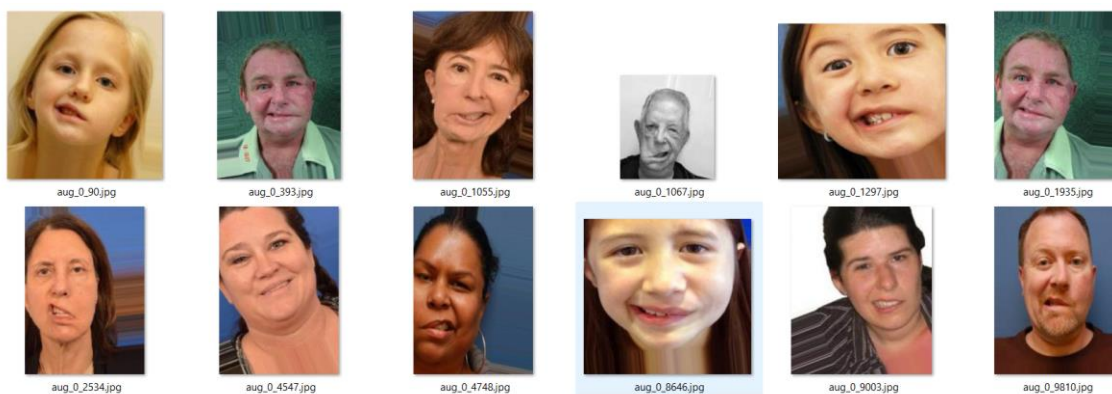
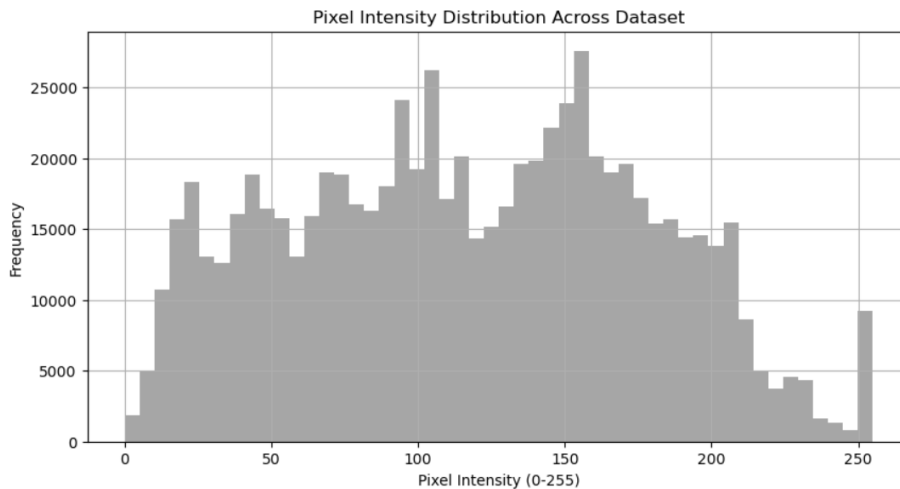


Figure 8. Sample of acute stroke dataset

A key feature of this dataset is the accompanying labels.csv file, which contains structured annotations for each image. These labels explicitly indicate which side of the patient's face is affected by the stroke — either "left" or "right". This information was crucial for supervised training of the asymmetry detection model.

Histogram analysis for acute stroke dataset

```
plt.ylabel("Frequency")  
plt.grid(True)  
plt.show()
```



Histogram analysis for acute stroke dataset

The histogram above represents the distribution of pixel intensities (grayscale values) across the Acute Stroke Dataset. Pixel values range from 0 (black) to 255 (white), and the plot illustrates how frequently each intensity occurs across all patient images.

Key Observations:

Broad Range of Intensities:

The dataset contains a wide spectrum of grayscale values, with no extreme skew toward either very dark or very light tones. This suggests that the images maintain balanced contrast, which is essential for ensuring facial features are clearly distinguishable during preprocessing and model training.

Well-Spread Midtones:

The highest frequency of pixel intensities falls between 100 and 180, which are midtone regions. This is common in clinical facial datasets, where lighting is typically neutral, and facial features (e.g., eyes, skin texture, wrinkles) are captured in detail.

Presence of Both Shadows and Highlights:

The presence of lower (0–50) and higher (200–255) intensity values indicates the dataset includes natural shadowing and lighting, adding realistic variability. These light and dark zones contribute to generalization, as the model learns to identify faces under varying illumination conditions.

Moderate Noise & Diversity:

The slight irregularities and small peaks throughout the histogram may indicate image-to-image diversity in background, lighting, and patient condition. This diversity is valuable for training robust models that generalize well to new patients.

3. Cleaning and Preprocessing of Dataset

3.1 Image Processing Steps

A custom Python script was developed to preprocess all facial images before training and testing the models. The script performs several critical steps:

1. Face Detection

The script uses Dlib's Histogram of Oriented Gradients (HOG)–based face detector to automatically locate faces within the images. This step ensures that only relevant facial regions are processed, discarding background noise and other irrelevant parts of the image.

Why it's important: This step standardizes the input and makes the subsequent landmark detection more accurate.

2. Facial Landmark Detection

Once a face is detected, the script applies a pre-trained 68-point facial landmark model from Dlib. These landmarks map key facial features such as the eyes, eyebrows, nose, mouth, and jawline.

Why it's important: These landmarks provide reference points for perfectly aligning and splitting the face down the vertical center, even if the head is slightly rotated or tilted.

3. Splitting Faces into Left and Right Halves

Using the detected landmarks (particularly the nose bridge and the center of the face), the script calculates a vertical midline and splits the face into left and right halves.

- For the Emotion Pain Dataset, both halves are saved separately to create two new datasets:

- Left-Side Pain Dataset
- Right-Side Pain Dataset

- For stroke patient images, this splitting enables the model to later focus only on the

unaffected side, once it has been detected.

Why it's important: This step supports the project's core idea — to assess pain using only the unaffected half of the face, which is assumed to remain expressive in stroke patients.

4. Saving Processed Images

After splitting, each half-face image is saved into a structured folder system:

- processed/left/ for left-side images
- processed/right/ for right-side images

The structure makes it easy to load, train, and evaluate the models using standard PyTorch Dataset and DataLoader utilities.

Why it's important: Organized datasets enable efficient training, reproducibility, and easy integration into machine learning workflows.

3.2 Preprocessing Considerations

- Normalization: Ensures each image maintains a standard format post-processing.
- Augmentation Decisions: No augmentation was applied to retain original expressions.
- Handling Missing Data: Images without clear facial detection were excluded from training.
- correcting the difference of naming of image dataset and facial action unit files
- removing most files with no pain intensity scale from the image dataset and facial action unit files
- Ada Boost classifier was tested to perform feature extraction due to being one of the validated solutions with results from research
- cv2 library was tested to perform split operation on the extracted face image due to being a recommended library with many possibilities

All the techniques were applied to the selected datasets to ensure that the datasets are relevant without augmentation. The model can be trained as soon as the datasets are merged.

Conclusion

This dataset report has demonstrated how carefully selected and preprocessed data can directly influence the performance and reliability of machine learning models in a healthcare context-specifically for pain intensity detection in stroke patients.

By combining two core datasets-the Emotion Pain Dataset (from healthy individuals expressing pain) and the Acute Stroke Dataset (from real-world stroke patients)—we established a solid foundation for building a modular deep learning pipeline.

A key innovation in this project was the use of facial symmetry-based preprocessing, where all facial images were split into left and right halves using precise landmark detection. This enabled us to develop two separate pain detection models:

One trained on left-side facial expressions

One trained on right-side expressions

This design reflects real clinical needs: stroke patients often lose expressiveness on one side of their face. By using the unaffected side, the models can still estimate pain intensity based on subtle but authentic expressions.

The Acute Stroke Dataset, enhanced with the labels.csv file indicating the affected side, was instrumental in training a model capable of identifying facial asymmetry. This allowed for automated selection of the expressive (unaffected) side of a patient's face.

Extensive data visualization, including histograms, landmark mappings, and bar plots of PSPI scores, provided key insights into dataset quality, diversity, and balance. Outlier detection and careful exclusion of mislabeled or irrelevant entries further improved data integrity.

Together, the cleaned, structured, and side-separated datasets allow for highly focused model training-maximizing performance while preserving generalization across varied patient profiles. This methodology establishes a scalable pipeline for clinical pain detection using only one side of the face, representing a novel and practical step forward in healthcare AI.