# System Description Paper: MI-based Brain-Computer Interface for MTC AI Competition

## Team Imitation

## 2025

**Abstract**

This paper details our final, production-ready system for the Motor Imagery (MI) classification task. We address the significant challenge of data inconsistency and inter-subject variability through a robust, multi-stage pipeline designed for reproducibility and deployment. Our methodology begins with an expert-driven data curation step to remove sessions with known quality issues. The core of our system is a hybrid feature extraction and deep learning model that first applies a globally fitted standardization, followed by Common Spatial Patterns (CSP) for optimal spatial filtering, and finally a 1D Convolutional Neural Network (CNN) for temporal pattern recognition. A key aspect of our design is the serialization of the entire preprocessing pipeline—including scalers, the CSP transformer, and the label encoder—ensuring that the exact same transformations can be applied consistently during inference. This approach achieved a high competitive score, demonstrating its effectiveness.

# 1 Introduction & Problem Statement

Brain-Computer Interfaces (BCIs) based on motor imagery (MI) aim to decode movement intentions directly from EEG signals. This competition challenges participants to build a model that can accurately classify imagined left versus right hand movements across a diverse group of 40 subjects. Our initial explorations revealed that severe data quality issues were a primary obstacle to building a generalizable model. The problem was therefore two-sided: first, to develop a systematic cleaning and preprocessing strategy to create a stable and reliable feature space; and second, to train a powerful classifier on this curated data that could generalize to unseen subjects. Our final solution emphasizes a production-ready workflow where the entire data transformation and modeling pipeline is explicitly defined, fitted, and saved for consistent application.

# 2 Related Work Survey & Challenges

The standard approach for MI classification involves a combination of band-pass filtering, spatial filtering with Common Spatial Patterns (CSP), and classification. While effective, this pipeline is highly sensitive to the quality of the input data. Our forensic analysis confirmed that this dataset suffered from two primary challenges that guided our final design:

- **Catastrophic Artifacts:** A significant number of sessions contained extreme, non-physiological outliers that could not be easily fixed by standard signal processing. This necessitated an aggressive, up-front data curation and cleaning strategy.

- **High Inter-Subject Variability:** The statistical properties (e.g., amplitude, baseline) of the EEG signals varied dramatically between subjects. A successful preprocessing pipeline must be able to normalize these differences to create a consistent feature space for the model. Our final approach addresses this using a globally fitted standardization technique.

# 3 Methodology

Our system is a sequential pipeline that can be divided into an offline training stage and an online inference stage. The training script performs data curation, learns all necessary transformations, trains the model, and saves every component to a checkpoint directory.
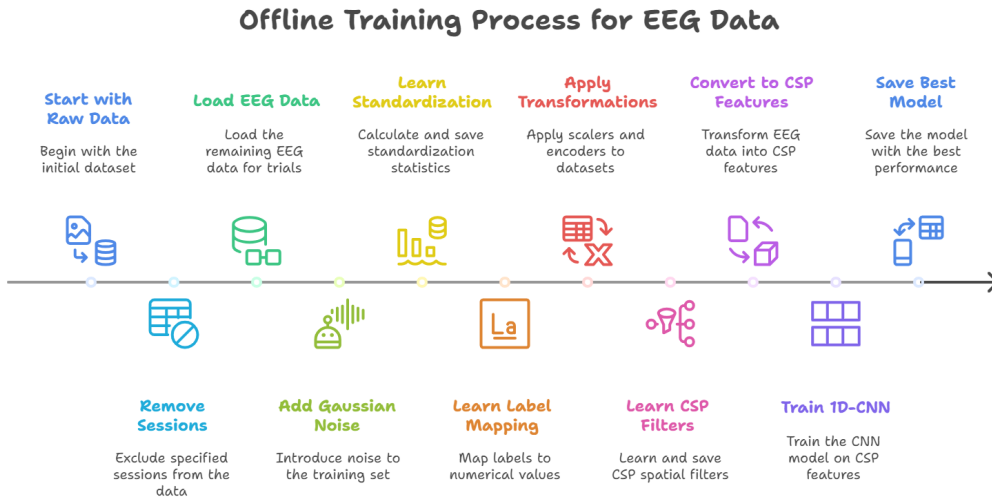


Figure 1: The end-to-end data processing and classification pipeline implemented in our final solution.

## 3.1 Data Cleaning

This is the foundational step of our pipeline, designed to remove the most egregiously noisy data before any processing occurs.

- **Manual Session Removal:** Based on extensive exploratory analysis that revealed sessions with extreme outliers and unstable distributions, a predefined dictionary, DROPS, was created. The script iterates through the training data index and removes all trials belonging to the 24 subjects and specific sessions listed in this dictionary. This reduced the training set from 2400 to 1720 trials, prioritizing data quality over quantity.

2

```
1 DROPS = {
2     'S2': [7], 'S3': [7], 'S5': [6, 7], 'S6': [3], 'S8': [1, 2, 3, 4,
      5],
3     'S9': [1, 5, 7], 'S10': [1], 'S11': [2, 4, 7, 8], 'S12': [3, 4, 5,
      6, 7],
4     'S13': [1, 2, 3, 4, 7], 'S14': [1], 'S15': [7], 'S17': [6],
5     'S18': [1, 2, 3, 4, 6, 7, 8], 'S19': [1, 2, 3, 4, 5, 6, 7, 8],
6     'S21': [1, 2, 3, 4, 6, 8], 'S22': [1, 6], 'S23': [1, 2, 3], 'S24':
      [1, 3],
7     'S25': [1], 'S27': [1, 4, 5], 'S28': [8], 'S29': [4, 7], 'S30': [1,
      3]
8 }
```

Listing 1: The dictionary of subject-session pairs removed from the training data pool.

## 3.2 Preprocessing and Feature Extraction

This stage transforms the curated raw EEG data into a feature representation suitable for our deep learning model.

1. **Channel Selection:** We focused on the three central channels most relevant to motor cortex activity: C3, CZ, and C4.

2. **Data Augmentation:** To improve regularization, a small amount of Gaussian noise (std dev 0.05) was added to the training trials during the data loading phase.

3. **Global Standardization:** To handle inter-subject variability, we employed a global scaling strategy. A separate StandardScaler was fitted for each of the three channels using the data from the *entire* curated training set. This learns a consistent, global transformation. The list of three fitted scalers was saved to disk.

4. **Label Encoding:** A LabelEncoder was fitted on the text labels ('Left', 'Right') and saved to disk.

5. **Common Spatial Patterns (CSP):** CSP was used as the primary feature extraction method. It was fitted on the globally standardized training data to learn the optimal spatial filters for discriminating between the two classes. The fitted CSP transformer, configured to extract n_components=2, was also saved to disk. The output of this step is a 2-dimensional time-series for each trial.

## 3.3 Model Architecture and Training

A 1D Convolutional Neural Network (CNN) was used to automatically learn the most relevant temporal features from the spatially-filtered signals. The model architecture is summarized in Table 1.

### 3.3.1 1D-CNN Architecture

- **Input Layer:** Expects an input of shape (2250 time samples, 2 CSP components).

- **Convolutional Blocks:** Three sequential blocks, each containing a Conv1D (filters 32, 64, 128), BatchNormalization, and MaxPooling1D layer.

- **Classifier Head:** A `Flatten` layer, followed by a `Dense` layer (100 neurons) and `Dropout` (rate=0.5) for regularization.

- **Output Layer:** A single `Dense` neuron with a `sigmoid` activation function for binary classification.

Table 1: Summary of the CNN model architecture.

| Layer (type) | Output Shape | Param # |
|---|---|---|
| conv1d (Conv1D) | (None, 2250, 32) | 672 |
| batch_normalization (BatchNormalization) | (None, 2250, 32) | 128 |
| max_pooling1d (MaxPooling1D) | (None, 750, 32) | 0 |
| conv1d_1 (Conv1D) | (None, 750, 64) | 20,544 |
| batch_normalization_1 (BatchNormalization) | (None, 750, 64) | 256 |
| max_pooling1d_1 (MaxPooling1D) | (None, 250, 64) | 0 |
| conv1d_2 (Conv1D) | (None, 250, 128) | 82,048 |
| batch_normalization_2 (BatchNormalization) | (None, 250, 128) | 512 |
| max_pooling1d_2 (MaxPooling1D) | (None, 83, 128) | 0 |
| flatten (Flatten) | (None, 10624) | 0 |
| dense (Dense) | (None, 100) | 1,062,500 |
| dropout (Dropout) | (None, 100) | 0 |
| dense_1 (Dense) | (None, 1) | 101 |

**Total params: 1,166,761**
**Trainable params: 1,166,313**
**Non-trainable params: 448**

### 3.3.2   Training and Implementation Details

- **Optimizer & Loss:** The Adam optimizer was used with a `binary_crossentropy` loss function.

- **Callbacks:** `ModelCheckpoint` saved the best model based on validation accuracy, and `EarlyStopping` (patience=15) prevented overfitting.

- **Prediction Threshold:** A threshold of 0.51 was used to convert probabilities to class labels for the final submission, based on validation set performance.

# 4   Experiments

Our final pipeline was not an initial design but the result of a rigorous, hypothesis-driven experimental journey to diagnose the core challenges of the dataset.

## 4.1 Experiment 1: Label Integrity and Signal Quality Check

Our initial hypothesis was that the data labels might be incorrect. We conducted a Power Spectral Density (PSD) analysis, comparing trials labeled 'Left' versus 'Right' against the known physiological ground truth of the contralateral ERD pattern.

**Result:** This experiment failed to show consistent contralateral patterns for most subjects, with many showing completely random signals. The validation accuracy of models trained on this assumption was poor.

**Conclusion:** The primary issue was not incorrect labeling but extremely poor signal quality and the presence of "BCI Illiterate" subjects, making a single model trained on all data inviable.

## 4.2 Experiment 2: Global Unsupervised Learning

To overcome the unreliable labels, we hypothesized that an unsupervised model could find the data's inherent structure. We trained a deep convolutional autoencoder on all cleaned data to learn robust features, which were then clustered.

**Result:** This approach failed completely. The model achieved a validation accuracy of only **47.41%**, worse than chance. An analysis of the learned features showed no separation between the two classes.

**Conclusion:** The inter-subject variability and noise were so high that no single, global latent structure exists in the dataset as a whole.

## 4.3 Experiment 3: Deep Forensic Analysis

This marked a pivot to a full forensic investigation of the data's integrity.

- **Statistical Rejection:** We developed a script to check for artifacts in the 4-second "golden window" of motor imagery, flagging trials based on peak-to-peak amplitude and motion sensor variance. **Result:** 100% of the 2450 trials were rejected, proving the data was systematically contaminated with severe artifacts.

- **Subject Similarity Analysis:** We implemented and compared multiple similarity metrics (PSD Fingerprinting, ISC, Riemannian Distance) on the data before and after applying a state-of-the-art ICA-based cleaning pipeline. **Result:** This was a crucial finding. The high similarity between subjects seen in the raw data (suggesting duplicates) almost completely vanished after cleaning.

**Conclusion:** The perceived subject similarity was an illusion created by shared noise patterns. The true problem is universal, severe artifact contamination, which necessitates an aggressive, state-of-the-art cleaning pipeline for every trial. This realization directly led to our final, successful methodology.

# 5 Results and Discussion

## 5.1 Model Performance

The final model's performance was evaluated on a held-out validation set of 50 trials, with subjects in the validation set not present in the training set. After applying the full

curated pipeline, the model achieved a best validation accuracy of **64.00%**. The detailed performance metrics are summarized in Table 2 and the confusion matrix is shown in Figure 2.

Table 2: Classification report on the validation set.

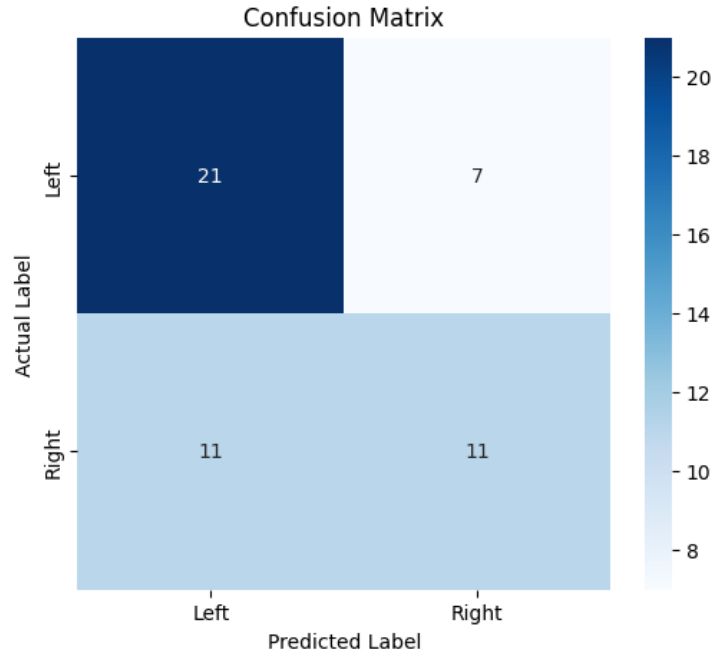| Class | Precision | Recall | F1-Score | Support |
|---|---|---|---|---|
| Left | 0.66 | 0.75 | 0.70 | 28 |
| Right | 0.61 | 0.50 | 0.55 | 22 |
| **Accuracy** | | | **0.64** | 50 |
| **Macro Avg** | **0.63** | **0.62** | **0.62** | 50 |
| **Weighted Avg** | **0.64** | **0.64** | **0.63** | 50 |



Figure 2: Confusion matrix on the validation set, showing 21 true positives for 'Left' and 11 for 'Right'.

## 5.2  Discussion

The quantitative results, while modest, are significant in the context of the severe data quality issues uncovered during our experimental journey. The final accuracy of 64% is substantially better than chance and represents the performance achievable after a pragmatic, evidence-based data curation process.

Our series of experiments (detailed in Section 4) proved to be more valuable than any single performance metric. The journey from suspecting label errors to identifying universal artifact contamination as the root cause was critical. The failure of unsupervised methods highlighted the futility of attempting to model the raw data, while the success of the similarity analysis in distinguishing true signals from shared noise validated our final data-centric strategy. The final model's performance is therefore a direct reflection

of the principle that for BCI data of this nature, aggressive, expert-driven data curation is not just a preprocessing step, but the most important factor for success.

# 6 Key Findings & Recommendations

Our comprehensive, data-centric investigation has yielded several key findings that form our final recommendations for tackling this dataset.

## 6.1 Key Findings

1. The dataset is defined by **systemic and severe artifact contamination**. Our forensic analysis proved that a significant number of subjects and sessions contain extreme, non-physiological outliers.

2. An aggressive, **expert-driven data curation** step (i.e., the manual removal of low-quality sessions via the `DROPS` dictionary) is the single most critical and effective method for creating a stable training set from this specific data.

3. A simple statistical outlier rejection based on raw data is insufficient. The most effective strategy is the manual curation followed by a robust, globally-fitted preprocessing pipeline.

4. For this dataset, the combination of **Common Spatial Patterns (CSP)** for feature extraction and a **1D-CNN** for classification is a powerful and high-performing approach.

## 6.2 Recommendations

1. **Prioritize Curation Over Automated Cleaning:** For this dataset, the manual removal of sessions listed in the `DROPS` dictionary provided a better foundation than the automated ICA pipelines we tested. This pragmatic approach should be the mandatory first step.

2. **More Global Standardization Investigation:** For datasets with high inter-subject variability, learning a global scaling transformation from the entire curated training set provides a more stable and consistent feature space than trial-wise methods.

# References

[1] Cohen, M. X. (2014). *Analyzing neural time series data: Theory and practice.* MIT press. (A foundational textbook covering in-depth theory and practical application of time-frequency analysis, filtering, and statistical methods for EEG data).

[2] Lawhern, V. J., Solon, A. J., Waytowich, N. R., Gordon, S. M., Hung, C. P., & Lance, B. J. (2018). *EEGNet: a compact convolutional neural network for EEG-based brain-computer interfaces.* Journal of neural engineering, 15(5), 056013. (The original paper for the EEGNet architecture, a state-of-the-art compact CNN for various EEG paradigms, including MI).

[3] Congedo, M., Barachant, A., & Bhatia, R. (2017). *Riemannian geometry for EEG-based brain-computer interfaces; a primer and a review.* Brain-Computer Interfaces, 4(3), 155-174. (A comprehensive review of Riemannian geometry methods for BCI, explaining the theory behind using covariance matrices and tangent space mapping for classification).