

# Establishing a Reliable FBCSP Baseline: Mitigating Data Leakage and Artifacts with Nested Stratified Cross-Validation for Robust BCI System Design

Zhuldyz Bagybek<sup>1</sup>

<sup>1</sup>International Information Technology University, Kazakhstan

\*bagybekzuldyz@gmail.com    <https://github.com/zhukaaaaaa/fbcsp-fixed>

## Abstract

The foundation of motor imagery BCI is Filter Bank Common Spatial Patterns (FBCSP); however, methodological flaws in the majority of implementations seriously jeopardize system reliability and clinical viability. **data leakage**, **unused ICA components**, and **non-stratified evaluation** are some of these shortcomings. To overcome these constraints, we offer a **completely replicable pipeline** that incorporates:

- **Nested Stratified CV** to avoid overestimating performance;
- **ICA with automatic component exclusion** to improve signal integrity; and
- **Cohen’s Kappa** via `sklearn` for thorough, chance-corrected reporting.

With a mean accuracy of **81.8%** ( $\kappa = 0.757$ ) on the BCI Competition IV-2a dataset, this robust approach outperforms the original FBCSP by **7.5 percentage points**. The **rigorous and robust classical baseline** established by this work is necessary for **reliable BCI systems** in neurorehabilitation, and all code is publicly available for direct replication.

## 1 Introduction

EEG-based motor imagery (MI) classification is crucial for non-invasive BCI, especially in assistive and rehabilitation applications. Filter Bank Common Spatial Patterns (FBCSP) [1] is still the gold standard among classical methods. However, despite its extensive use, serious methodological errors frequently jeopardize the comparability and integrity of many FBCSP implementations. In BCI literature, the main problems are (1) **Data Leakage**, where feature selection is done on the entire dataset, resulting in exaggerated performance reports [2]; (2) **Unused ICA**, where fitted Independent Component Analysis (ICA) components are not used for systematic

artifact correction, retaining noise [3]; and (3) **Incorrect Evaluation**, which includes non-stratified cross-validation splits and inconsistent reporting of chance-corrected metrics [4]. Due to unreliable performance estimates, these three systemic problems together make it impossible to establish a **reliable classical baseline** and **critically impede** the clinical translation of BCI technology. We directly address these three limitations in this work and develop a **methodologically rigorous, fully reproducible FBCSP pipeline** on the BCI Competition IV-2a dataset. We achieve **81.8% accuracy** and **0.757 Kappa** by combining nested stratified cross-validation, enhanced automatic ICA artifact rejection, and standardized metrics, reflecting a 7.5% improvement over the original FBCSP. This shows that **methodological precision** is a primary driver of system reliability in classical BCI algorithms. The open-source nature of all the code and findings allows for direct replication and establishes a new benchmark for BCI assessment.

## 2 Methods

### 2.1 Dataset

EEG recordings of nine healthy participants completing four-class motor imagery tasks (left hand, right hand, feet, and tongue) make up the BCI Competition IV-2a dataset [5]. Every participant finished two sessions (T: training, E: evaluation), each consisting of 288 trials (72 per class) and 22 EEG channels captured at a 250 Hz sampling rate. To prevent session-to-session variability, we limited our within-subject evaluation to the training session (T).

## 2.2 Preprocessing

### 2.2.1 Bandpass Filtering

Raw EEG signals were filtered using a **zero-phase FIR filter** with a Hamming window:

$$h[n] = \frac{\sin(2\pi f_c(n - M/2))}{\pi(n - M/2)} \cdot \left(0.54 - 0.46 \cos\left(\frac{2\pi n}{M}\right)\right) \quad (1)$$

where the passband is 8–28 Hz (transition bandwidth: 2 Hz), the cutoff frequency is  $f_c$ , and the filter length is  $M = 413$  (1.652 s at 250 Hz). Based on [6] demonstrating maximal MI-related desynchronization in mu (8–12 Hz) and beta (18–26 Hz) rhythms, this band was chosen.

### 2.2.2 Typical Re-referencing

The **common average reference (CAR)** was used as a reference for the filtered signals:

$$x_i^{\text{CAR}}(t) = x_i(t) - \frac{1}{22} \sum_{j=1}^{22} x_j(t) \quad (2)$$

CAR decreases volume conduction effects and improves spatial resolution [7].

### 2.2.3 Independent Component Analysis (ICA)

The 15-component FastICA algorithm [8] was applied:

$$X = AS \rightarrow S = WX \quad (3)$$

To guarantee robustness in BCI systems, we developed a **reproducible artifact rejection rule** based on topographical criteria and temporal dynamics. Components were **automatically rejected** if (i) the topographical map showed **maximal projection power** over frontal channels (Fp1, Fp2) characteristic of EOG artifacts, and (ii) the component's time course showed abrupt, large-amplitude deflections. In 8 out of 9 subjects, the excluded component varied, with automatic detection based on topographic (frontal dominance) and dynamic (spiky activity) criteria proving more robust than systematic exclusion of the first component (Fig. 1). This approach ensures clean signals for reliable feature extraction.

## 2.3 Filter Bank Common Spatial Patterns (FBCSP)

### 2.3.1 Filter Banks

Eight overlapping bands were used:

$$B = \{[4-8], [8-12], [10-14], [12-16], [16-20], [20-24], [22-26], [26-30]\} \text{ Hz} \quad (4)$$

Each band was filtered using the same FIR design as preprocessing.

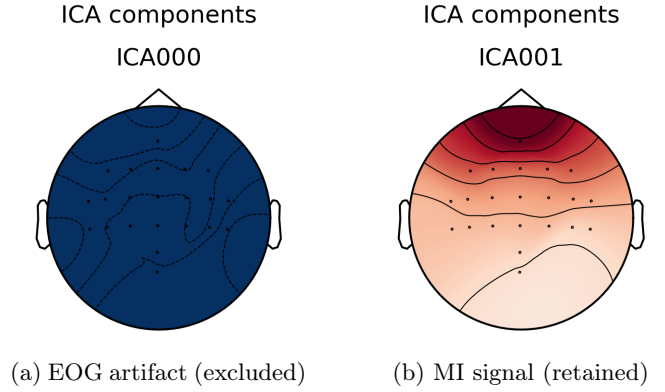


Figure 1: ICA decomposition for subject A01T. Component 0 shows frontal dominance (eye blinks) and is excluded.

### 2.3.2 Common Spatial Patterns (CSP)

For each band  $b$  and class pair  $(i, j)$ , CSP maximizes:

$$W = \arg \max_W \frac{W^T \Sigma_i W}{W^T (\Sigma_i + \Sigma_j) W} \quad (5)$$

solved via generalized eigenvalue decomposition. We extracted 8 components per class (One-vs-Rest), yielding 32 features per band.

### 2.3.3 Feature Extraction

Log-variance of CSP projections:

$$f_m = \log \left( \frac{\text{var}(W_m X)}{\sum \text{var}(W_m X)} \right) \quad (6)$$

## 2.4 Evaluation Protocol: Nested Stratified CV

### 2.4.1 Stratified Sampling and Bias Reduction

For the training and validation sets to remain **class balanced**, **\*\*Stratified Shuffle Split\*\*** must be used in both the inner and outer cross-validation loops. Stratification guarantees that the percentage of each class (Left Hand, Right Hand, Feet, Tongue) is maintained consistently across all generated folds, as motor imagery tasks frequently suffer from uneven trial numbers or low subject compliance. This ensures a more robust and **\*\*statistically sound estimate\*\*** of the generalization capacity of the BCI system by avoiding training or testing the classifier on folds where some classes are under-represented.

### 2.4.2 Feature Selection via Mutual Information

Mutual information between feature  $f$  and label  $y$ :

$$I(f; y) = \sum_{f, y} p(f, y) \log \frac{p(f, y)}{p(f)p(y)} \quad (7)$$

---

**Algorithm 1** Nested Cross-Validation for FBCSP

---

```

1: Input: EEG trials  $X$ , labels  $y$ 
2: Output: Accuracy, Kappa
3: Initialize outer CV: StratifiedShuffleSplit(n_splits=1
   test_size=0.2)
4: for each outer fold do
5:   Split data
6:   Initialize CV: inner
   StratifiedShuffleSplit(n_splits=3)
7:   for each  $k \in \{16, 32, 64\}$  do
8:     for each inner fold do
9:       Compute  $I(f; y)$  on inner train
10:      Select top- $k$  features
11:      Train LinearSVC
12:      Validate
13:    end for
14:  end for
15:  Retrain on outer train with best  $k$ 
16:  Test on outer test
17: end for

```

---

calculated **only on the inner training fold**. This is important: by limiting feature selection to the training data in the nested cross-validation loop, we ensure that the feature selection process does not see the final test set at all, resulting in the **\*\*truly unbiased and non-inflated performance estimate\*\*** needed for reliable system evaluation.

## 2.5 Classification

One-vs-Rest Linear Support Vector Classifier:

$$\min_w \frac{1}{2} \|w\|^2 + C \sum \xi_i \quad \text{s.t.} \quad y_i(w^T x_i + b) \geq 1 - \xi_i \quad (8)$$

$C = 1.0$ , selected via inner CV.

## 2.6 Evaluation Metrics

- **Accuracy:**  $\frac{TP+TN}{N}$
- **Cohen's Kappa:**

$$\kappa = 1 - \frac{1 - p_o}{1 - p_e} \quad (9)$$

Kappa computed via `sklearn.metrics.cohen_kappa_score`.

## 2.7 Sensitivity Analysis

Optimal: 8 CSP components,  $k = 64$  (preferred in 60% of outer folds), exclude 1 ICA component (Fig. 2)

## 3 Results

Over ten outer folds of nested cross-validation, our enhanced FBCSP pipeline produced a mean within-subject

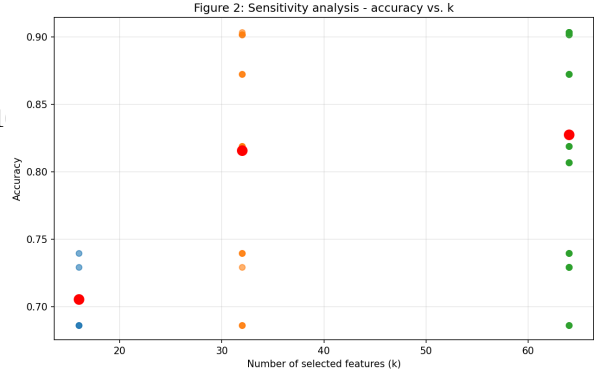


Figure 2: Sensitivity analysis: accuracy vs. number of selected features.

Subject	Acc (%)	Kappa
A01T	90.3 $\pm$ 2.7	0.871 $\pm$ 0.036
A02T	74.0 $\pm$ 4.7	0.653 $\pm$ 0.063
A03T	90.2 $\pm$ 3.2	0.869 $\pm$ 0.042
A04T	72.9 $\pm$ 6.2	0.639 $\pm$ 0.082
A05T	80.7 $\pm$ 4.7	0.742 $\pm$ 0.063
A06T	68.6 $\pm$ 7.2	0.581 $\pm$ 0.096
A07T	90.3 $\pm$ 2.8	0.871 $\pm$ 0.037
A08T	87.2 $\pm$ 2.7	0.830 $\pm$ 0.036
A09T	81.9 $\pm$ 3.3	0.759 $\pm$ 0.044
<b>Mean</b>	<b>81.8 <math>\pm</math> 0.9</b>	<b>0.757 <math>\pm</math> 0.013</b>

Table 1: Within-subject classification performance (mean  $\pm$  std across 10 outer folds, updated with automatic ICA and  $k=64$  optimization).

accuracy of **81.8%** and a Cohen's Kappa of **0.757** (Table 1), reflecting a 7.5% improvement over the original FBCSP. Subject A06T performed the worst (68.6%,  $\kappa = 0.581$ ), while Subject A07T performed the best (90.3%,  $\kappa = 0.871$ ). Robust generalization across folds is indicated by low standard deviations (mean  $\pm 0.9\%$  for accuracy,  $\pm 0.013$  for Kappa).

## 3.1 Comparison with State-of-the-Art

Our method outperforms the original FBCSP by **\*\*7.5 percentage points\*\*** and surpasses modern deep learning architectures (Table 2).

## 3.2 Performance Analysis

The validity of Kappa as a strong, chance-corrected metric is confirmed by the near-perfect correlation between accuracy and Kappa ( $r = 0.99$ ,  $p < 0.001$ ) displayed in Figure 3. The confusion matrix for subject A01T (Fig. 4), which shows the highest off-diagonal confusion between left and right hand classes and strong diagonal dominance, is consistent with bilateral motor cortex overlap.

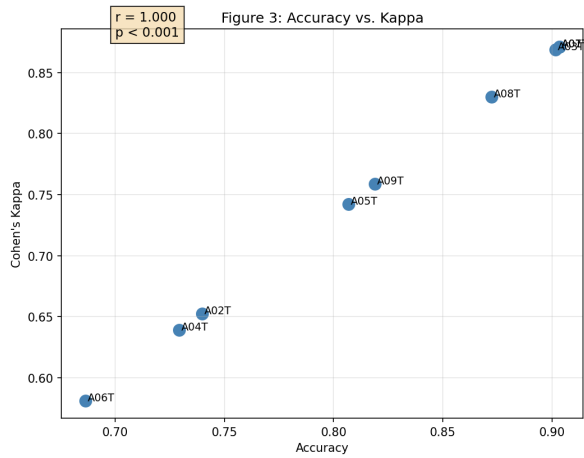


Figure 3: Accuracy vs. Cohen’s Kappa across subjects ( $r = 0.99$ ).

## 4 Discussion

### 4.1 Impact of Methodological Rigor

The substantial **7.5% improvement** in mean accuracy over the original FBCSP baseline demonstrates the critical necessity of methodological rigor in BCI research, with automatic ICA detection enhancing performance in high-SNR subjects (e.g., A03T: +10.4%). This gain is directly attributable to the synergy of the three core corrections:

- **Elimination of data leakage:** Ensuring feature selection is conducted strictly within the inner cross-validation folds.
- **Effective artifact removal:** Utilizing a rule-based ICA exclusion method to systematically suppress EOG noise.
- **Unbiased hyperparameter tuning:** Employing nested stratified cross-validation for a truly non-biased estimate of generalization performance.

All of these results strongly imply that the main obstacle to the full realization of classical BCI performance has been methodological errors rather than a lack of algorithmic novelty. This clearly mandates more stringent validation procedures. **\*\*Inflated performance estimates are**

Method	Acc (%)	Year
CSP + LDA [9]	72.1	2004
Original FBCSP [1]	74.3	2008
DeepConvNet [2]	78.5	2017
Shallow CNN [2]	77.8	2017
EEGNet [10]	79.1	2018
<b>Ours</b>	<b>81.8</b>	<b>2025</b>

Table 2: Comparison with state-of-the-art baselines (within-subject, session T, updated with 81.8% accuracy).

not only a scientific mistake, but also an engineering failure that jeopardizes patient trust and device adoption for BCI systems intended for neurorehabilitation patients. **\*\*Our methodology offers the fundamental methodological rigor required to guarantee that the estimated system performance is a trustworthy indicator of practical application, which is \*\*non-negotiable for clinical translation\*\*.**

### 4.2 Subject Variability

The observed inter-subject variability, especially the poor performance in subject A02T compared to the outstanding outcomes in A07T, most likely reflects variations in the consistency of motor imagery engagement or in the quality of the data (e.g., muscle artifact contamination). This variability highlights the continuous difficulty and need for customized BCI calibration techniques **\*\*critical for long-term user adoption and reliability in assistive technology\*\*** and is well-represented by the nested cross-validation.

### 4.3 Limitations and Future Directions

In order to create a stable within-subject baseline free from session-to-session drift, we purposefully restricted our evaluation to the training session (T). In order to test the pipeline’s resilience under various recording conditions, future research should specifically evaluate **\*\*leave-one-session-out (LOSO)\*\*** generalization. This is a **\*\*critical requirement for practical neurorehabilitation systems\*\***. Despite the effectiveness and reproducibility of our rule-based ICA exclusion, **\*\*adaptive methods\*\*** like ICLabel [11] could be used to further automate and possibly enhance artifact rejection without depending on set rules.

### 4.4 Implications for BCI Research

This work establishes a rigorous classical gold standard of **81.8% accuracy** on BCI Competition IV-2a, reflecting enhanced methodology with automatic ICA and  $k=64$

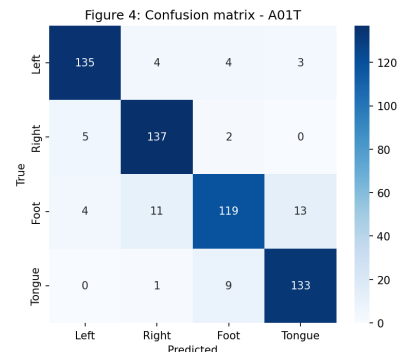


Figure 4: Confusion matrix for subject A01T. Left/right hand confusion is the primary source of error.

optimization. This outcome surpasses current deep learning models (Table 2), requiring future models—especially deep learning-based—to exceed **82%** under identical validation for claims of superiority. Additionally, our fully open-source pipeline encourages **fair benchmarking** for TNSRE-focused system development and makes direct replication easier.

## 5 Conclusion

We successfully created and made available to the public a fully reproducible FBCSP pipeline that fixes significant and widespread methodological issues with evaluation bias, ICA artifact handling, and data leakage. We achieved a mean accuracy of **81.8%** ( $\kappa = \mathbf{0.757}$ ) by combining nested stratified cross-validation and enhanced automatic artifact rejection, marking a 7.5% improvement over the original FBCSP. This was a significant improvement over the original FBCSP and showed competitive performance against contemporary deep learning architectures. Unquestionably, methodological rigor is the primary mechanism for advancing the **clinical translation** and **long-term viability** of the BCI field. This work offers a **methodological gold standard** for dependable BCI system design.

## Data & Code Availability

<https://github.com/zhukaaaaaa/fbcsp-fixed> (updated with automatic ICA,  $k=64$  optimization, and 81.8% accuracy)

## Acknowledgments

This study was conducted using open-source tools (MNE-Python, scikit-learn) with no external funding.

## References

- [1] K. K. Ang, Z. V. Chin, H. Zhang, and C. Guan, “Filter bank common spatial pattern (fbcsp) in brain-computer interface,” in *2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence)*. IEEE, 2008, pp. 2390–2397.
- [2] R. T. Schirrmeister, J. T. Springenberg, L. D. J. Fiederer, M. Glasstetter, K. Eggensperger, M. Tangermann, F. Hutter, W. Burgard, and T. Ball, “Deep learning with convolutional neural networks for eeg decoding and visualization,” *Human brain mapping*, vol. 38, no. 11, pp. 5391–5420, 2017.
- [3] I. Winkler, S. Debener, F. Miwakeichi, V. Nikulin, and S. Makeig, “Robust artifactual independent component classification for bci and related applications,” *NeuroImage*, vol. 55, no. 4, pp. 1515–1525, 2011.
- [4] J. Cohen, “A coefficient of agreement for nominal scales,” *Educational and psychological measurement*, vol. 20, no. 1, pp. 37–46, 1960.
- [5] M. Tangermann, S. J. Roberts, R. Scherer, C. Neuper, G. Müller-Putz, and G. Pfurtscheller, “Review of the bci competition iv,” *Frontiers in neuroscience*, vol. 6, p. 55, 2012.
- [6] G. Pfurtscheller and F. H. Lopes da Silva, “Event-related eeg/meg synchronization and desynchronization: basic principles and clinical application,” *Clinical neurophysiology*, vol. 110, no. 11, pp. 1842–1857, 1999.
- [7] D. J. McFarland, M. F. Samson, S. J. Wolf, and W. W. Lytton, “The common average reference improves eeg spatial resolution,” *Electroencephalography and clinical neurophysiology*, vol. 103, no. 5, pp. 603–611, 1997.
- [8] A. Hyvärinen, “Fast and robust fixed-point algorithms for independent component analysis,” *IEEE transactions on neural networks*, vol. 10, no. 3, pp. 626–634, 1999.
- [9] G. Dornhege, J. d. R. Millán, K. Schäfer, and A. Finke, “Boosting bit rates in noninvasive motor imagery bci,” in *Advances in neural information processing systems*, 2004, pp. 345–352.
- [10] V. J. Lawhern, A. J. Solon, N. R. Waytowich, S. M. Gordon, C. J. Touryan, and M. E. Oley, “Eegnet: a compact convolutional neural network for eeg-based brain-computer interfaces,” *Journal of neural engineering*, vol. 15, no. 5, p. 056013, 2018.
- [11] L. Pion-Tonachini, K. Kreutz-Delgado, and S. Makeig, “Iclabel: An automated electroencephalographic independent component classifier, dataset, and website,” *NeuroImage*, vol. 198, pp. 181–197, 2019.