

# Mixed-Precision Quantization and Parallel Implementation of Multispectral Riemannian Classification for Brain–Machine Interfaces

Xiaying Wang\*, Tibor Schneider\*, Michael Hersche\*, Lukas Cavigelli†, Luca Benini\*‡

\*ETH Zürich, D-ITET, Switzerland †University of Bologna, DEI, Italy ‡Huawei Technologies, Zurich RC, Switzerland

**Abstract**—With Motor-Imagery (MI) Brain–Machine Interfaces (BMIs) we may control machines by merely thinking of performing a motor action. Practical use cases require a wearable solution where the classification of the brain signals is done locally near the sensor using machine learning models embedded on energy-efficient microcontroller units (MCUs), for assured privacy, user comfort, and long-term usage. In this work, we provide practical insights on the accuracy-cost trade-off for embedded BMI solutions. Our proposed Multispectral Riemannian Classifier reaches 75.1% accuracy on 4-class MI task. We further scale down the model by quantizing it to mixed-precision representations with a minimal accuracy loss of 1%, which is still 3.2% more accurate than the state-of-the-art embedded convolutional neural network. We implement the model on a low-power MCU with parallel processing units taking only 33.39 ms and consuming 1.304 mJ per classification.

**Index Terms**—brain–machine interface, edge computing, parallel computing, machine learning, deep learning, motor imagery.

## I. INTRODUCTION

Motor-Imagery (MI) Brain–Machine Interfaces (BMIs) use Electroencephalography (EEG) signals recorded from the brain to decode a movement imagined by the subject. The decoded information can be used to control an external device, such as a drone [1] or a wheelchair [2], [3], or for stroke rehabilitation [4]. It is especially useful for individuals with physical disabilities to regain independence [4], [5]. However, the high variability across subjects and among different recording sessions poses big challenges to an accurate MI-BMI. Moreover, recording and labelling EEG data is expensive, time consuming, and prone to errors, resulting in scarce amounts of data available for training complex models with large numbers of parameters. In fact, many studies using Convolutional Neural Networks (CNNs) acknowledge the fact that overfitting is the biggest issue for these types of models [6], [7], [8].

On the other hand, successful methods have been proposed to extract discriminative, domain-specific features from EEG signals. The well-known Common Spatial Patterns (CSP) learns spatial filters that discern between different MI tasks [9]. An improved algorithm, called Filter-Bank CSP, that accounts for multiple frequency bands achieved better accuracy [10]. More recent studies have proposed Riemannian methods to extract more comprehensive features also in absence of labeled data [11]. The unsupervised feature calibration enables online adaptation of the classifier to combat the large inter-session variance in MI-BMIs [12]. So far, these methods are believed to be the most promising feature extractors for several kinds of BMI paradigms [13], [14], [15].

Traditional BMI systems adopt offline, remote processing of the sensor data, raising concerns over data privacy, latency,

high energy consumption, and battery lifetime. A promising solution is to bring the processing near the sensor, i.e. on the body of the user, using low-power low-cost microcontroller units (MCUs), allowing the data to be processed locally [16]. However, these devices suffer from limited on-board resources in terms of memory and computational capabilities. Hence, researching compact yet accurate algorithms [7], [17] and designing low-power processors with high capabilities [18] has become an emerging trend. Most of the MI-BMI models, particularly CNNs, are too demanding for low-power MCUs [19]. TPCT [20] reached the state-of-the-art (SoA) accuracy of 88.87% on the BCI Competition IV-2a dataset [21]. The model consists of around 7.78 M parameters. Other similar CNNs reach 81.1% with 240 k parameters [22] or 75.8% with 155 k parameters [23]. A notable exception is EEGNET [17] with only few thousands parameters, i.e. three orders of magnitude less demanding, but still achieving around 70% accuracy on 4-class MI classification. By virtue of its compactness, it has been successfully quantized with Q-EEGNET [24] and implemented on a low-power System-on-Chip (SoC) based on RISC-V called Mr. Wolf [18]. It has proven to be three orders of magnitude more energy efficient than an implementation on commercially available MCUs based on ARM Cortex-M architecture [25], making it the SoA embedded CNN in terms of energy efficiency, compact model size yet accurate performance. Another effort for embedded BMI has been made by Belwafi et al. [26] implementing a CSP-based classifier on a FPGA device. The multispectral and multiscale Riemannian classifiers proposed in [27], [28] outperform both EEGNET and CSP-based models by around 5% and 2% higher accuracy, respectively. However, their proposed models are still very challenging for embedded deployment on low-power resource-constrained MCUs due to large memory footprint and high computational complexity.

For the first time in literature, we propose an embedded MI-BMI based on a Riemannian classifier [27]. The main contributions of this paper are: (a) We tailor the model for better embedded deployment by reducing its size and complexity, i.e., the number of frequency bands and temporal windows, while at the same time keeping comparable classification accuracy by introducing regularization (75.1% ours vs. 75.5% [27]). (b) We further quantize the Multispectral Riemannian Classifier (MRC) from full precision (32-bit float) to a mixture of precisions with 8-, 16-, 32-bit fixed- and floating-point representations, to maximize efficiency on low-power MCUs by enabling the use of fixed-point SIMD instructions while maintaining a minimal accuracy loss. The quantization yields 1% accuracy drop which is still 3.2% more accurate than the embedded CNN-based EEGNET (74.1% ours vs. 70.9% [24]).

(c) We efficiently implement the mixed-precision model on Mr. Wolf by exploiting the underlying hardware architecture, i.e., custom Instruction Set Architecture (ISA) extensions and concurrent execution on multiple cores and measure the performance on-board. Experimental measurements show that the proposed model takes only 33.30ms and consumes 1.304mJ per inference. (d) Our work provides a practical accuracy-cost trade-off between MRC, a discriminative feature-based approach, and EEGNET as a CNN-based approach, supported by an actual implementation and measurement results. Being the first embedded implementation of Riemannian covariance kernels and the most accurate embedded MI-BMI, it opens the path for other BMI paradigms deploying Riemmanian methods, e.g., steady-state visual evoked potential [15] and P300 [29]. Finally, we release open-source code<sup>1</sup>.

## II. DESIGN AND QUANTIZATION

MRC [27] consists of a non-linear feature extraction applying the Riemannian covariance method [30] on multiple frequency bands and temporal windows, followed by a linear Support Vector Machine (SVM), depicted in Fig. 1. First, the input data is filtered using  $f$  different Infinite Impulse Response (IIR) bandpass filters. Then, the covariance matrix is estimated and regularized with the parameter  $\rho$ . The next block, called *Whitening*, multiplies from the left and right with a reference matrix  $\mathbf{C}_{\text{ref},k}^{-1/2}$ , that is computed for each frequency band  $k$  independently during training. Afterwards, the matrix logarithm is computed with the help of Eigendecomposition (EVD). Then, the function  $\text{vect}(\mathbf{L}_k)$  vectorizes the symmetric matrix  $\mathbf{L}_k$  by concatenating the diagonal values and the upper right non-diagonal elements. To preserve the norm, the off-diagonal elements are scaled with  $\sqrt{2}$ . Finally, the SVM classifier predicts the MI class.

We quantize the feature extraction to a mixture of 8-, 16-, and full precision 32-bit fixed- and floating-point representations and the SVM to 8-bit fixed-point, summarized in Fig. 2. The decision on the precision depends on the trade-off between energy efficiency and accuracy preservation. With 8- or 16-bit fixed-point numbers, it is possible to exploit the Single Instruction, Multiple Data (SIMD) instructions. However, not all the parts of the MRC can be quantized due to numerical instability and significant accuracy loss.

1) *IIR Bandpass Filters*: The input data  $\mathbf{X} \in \mathbb{R}^{N_{ch} \times N_s}$  with dimensions number of EEG channels  $N_{ch}$  and number of time samples  $N_s$ , is quantized to 8 bits. Each channel is filtered with  $f$  IIR bandpass filters. The filters can become unstable, especially with quantization. The internal accumulators can diverge, even if the output remains bounded. We implement the Direct-Form I defined in [31], since it does not experience numerical overflow in the internal signals, because all internal registers store either the input or the output of the filter [31]. A typical approach for quantizing an IIR filter is to express them as a cascade of Second-Order Sections (SOSs), each of which can be quantized with different dynamic ranges, thus minimizing the effect of quantizing the filter coefficients on the impulse response. With 8-bit fixed-point quantization, the impact is significant, while with 12 bits these effects are minimal. Therefore, we choose 12 bits for the filter coefficients to prevent overflows that would occur with 16 bits. We re-scale

the intermediate results in between the SOSs to remain in the same dynamic range and accumulate them with 16-bit registers in order to use SIMD operations for the following iteration. All dynamic ranges for all sections are chosen independently and forced to be a power of two to implement simple bit-shifts instead of expensive divisions.

2) *Covariance Matrix and Whitening*: Recall, that the covariance matrix  $\mathbf{C} \in \mathbb{R}^{n \times n}$ , in our case  $n = N_{ch}$ , including regularization, is computed as

$$\mathbf{C} = \mathbf{X}\mathbf{X}^T + \rho\mathbf{I} \quad (1)$$

and Whitening is defined as

$$\mathbf{W} = \mathbf{C}_{\text{ref}}^{-1/2} \mathbf{C} \mathbf{C}_{\text{ref}}^{-1/2}, \quad (2)$$

with  $\mathbf{C}_{\text{ref}}^{-1/2}$  being the reference matrix, computed by averaging the covariance matrices of all the training trials. For quantization, we define  $n_c$  and  $n_{\text{ref}}$  to be the number of bits to represent  $\mathbf{C}$  and  $\mathbf{C}_{\text{ref}}^{-1/2}$ , respectively. Since we can exploit either 4- or 2-way SIMD operations, we test both  $n_c = n_{\text{ref}} = 8$  and 16. However, the former yields a significant accuracy drop, while the latter causes overflows. Hence, we reduce  $n_{\text{ref}}$ , until training completes without overflow, resulting in  $n_{\text{ref}} = 11$ . Our experiments have shown that using  $n_c = 16$  and  $n_{\text{ref}} = 11$  yields similar accuracy to the full-precision version. Moreover, we force the scaling factor for the covariance matrix computation to be a power of two to exploit bit-shifts, while the dynamic range for the Whitening depends on the quantization of  $\mathbf{C}$  and  $\mathbf{C}_{\text{ref}}^{-1/2}$ . Finally, for the intermediate and final results of Eq. 2, we keep the full dynamic range with 32 bits since the input to the matrix logarithm is very sensitive to quantization errors, as explained next.

3) *Matrix Logarithm*: The matrix logarithm of a square, positive definite matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is defined in terms of its EVD, as

$$\log(\mathbf{A}) = \mathbf{Q}^{-1} \log(\mathbf{D}) \mathbf{Q}, \quad (3)$$

where  $\mathbf{A} = \mathbf{Q}^{-1} \mathbf{D} \mathbf{Q}$ , and the logarithm of a diagonal matrix  $\mathbf{D}$  is computed by applying the logarithm to its diagonal elements. The whitened covariance matrix  $\mathbf{W}$  in MRC is dense and symmetric, allowing us to optimize the EVD. We first compute the tridiagonal decomposition to obtain a tridiagonal matrix  $\mathbf{T}$  similar to the original one, i.e. the Eigenvalues are preserved. Then the EVD can be computed on  $\mathbf{T}$  requiring less computational effort. The final transformation is

$$\mathbf{W} = \mathbf{Q}_t^T \mathbf{T} \mathbf{Q}_t = \mathbf{Q}_t^T \mathbf{Q}_d^T \mathbf{D} \mathbf{Q}_d \mathbf{Q}_t, \quad (4)$$

where  $\mathbf{Q}_t$  is the orthogonal matrix for the tridiagonal transformation and  $\mathbf{Q}_d$  the one for the EVD.  $\mathbf{Q}_d \mathbf{Q}_t$  is an orthogonal matrix containing the Eigenvectors of  $\mathbf{W}$ . To compute the tridiagonal matrix, we use the Householder transformation [32]. The complexity of the transformation can be reduced by rearranging the operations and exploiting the sparsity of the vectors [32]. For computing the diagonal matrix  $\mathbf{D}$  from the tridiagonal symmetric matrix  $\mathbf{T}$ , we use the QR algorithm with implicit Wilkinson Shift [33]. The matrix logarithm only exists if the matrix is positive definite, meaning that all the Eigenvalues are positive. In full-precision MRC, the input of the matrix logarithm is always positive definite, while with quantization the Eigenvalues change and in some cases even become negative, making it impossible to compute real logarithm. We address this issue by (a) making use of the entire 32-bit dynamic range for the inputs, and (b) clipping all Eigenvalues  $\lambda_k$  to  $\max\{\lambda_k, \lambda_{\min}\}$  by introducing a threshold

<sup>1</sup><https://github.com/pulp-platform/multispectral-riemannian>

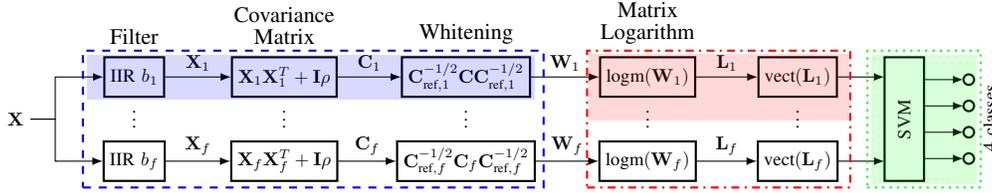


Fig. 1: Multispectral Riemannian Classifier with  $n = 18$  frequency bands and one time window.

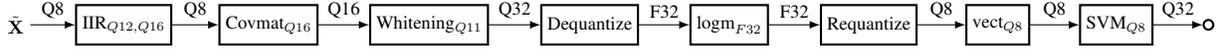


Fig. 2: Quantized MRC of a single frequency band, showing the representation of each intermediate signal.

$\lambda_{\min} = 10^{-3}$  to ensure all Eigenvalues remain above zero. Its value is chosen based on the smallest Eigenvalue occurring while training the full precision MRC. Moreover, both Householder transformation and QR algorithm are computed with 32-bit floating-point values. Finally, we convert the results back to 8-bit fixed-point format using the dynamic ranges after training. We do not rescale the output of the SVM because the prediction is made based on the relative largest output value. Hence, the weight vector can use the entire range available with 8 bit, reducing the quantization error.

4) *Support Vector Machine (SVM)*: The final classifier in MRC is a SVM, which we train on the quantized features. The weights and biases are then quantized with bit-width  $n_w = 8$  and  $n_b = 32$ , respectively, by determining the dynamic ranges after training. We do not rescale the output of the SVM because the prediction is made based on the relative largest output value. Hence, the weight vector can use the entire range available with 8 bit, reducing the quantization error.

### III. IMPLEMENTATION

We implement the mixed-precision MRC on Mr. Wolf [18] which has a SoC domain and a compute cluster with 8 parallel RISC-V-based processors called RI5CY, or CV32E40P, implementing RV32IMFC ISA with custom XPULPV2 extensions for Digital Signal Processing (DSP), e.g., SIMD instructions, hardware loops, post-incremental load and store [34]. The cluster cores have two shared Floating Point Units (FPUs) and 64 kB of shared L1 memory via the Tightly Coupled Data Memory (TCDM) interconnect. More memory can be accessed via a Direct Memory Access (DMA) unit from the shared L2 memory (448 kB) present in the SoC domain.

Our MRC implementation is divided into three main blocks framed with blue, red, and green lines in Fig. 1, respectively: (a) computation of the frequency bands until Whitening: each frequency band, highlighted with blue rectangle, is computed using 8 cores as described in the following paragraphs; (b) computation of the matrix logarithm and vectorization: every core computes one matrix logarithm followed by the vectorization concurrently with the other cores, i.e. 8 matrix logarithms, colored with red rectangle, are computed at the same time; (c) SVM computed with a single core, colored in green.

1) *IIR Filter*: As described in Section II-1, we set the bit-width of the coefficients to  $n_a = n_b = 12$ , and the bit-width of the internal registers to  $n_i = 16$ . Each SOS contains three Multiply Accumulates (MACs) for the forward accumulation and two MACs for the backward accumulation. This enables the usage of SIMD instructions with bit-width 16. We compute the filtered output of different EEG channels on separate cores of the cluster to utilize the concurrent capabilities of Mr. Wolf.

2) *Covariance Matrix*: The computation of the covariance matrix is a matrix-matrix multiplication (MMM), as shown in

Eq. (1), which results in a symmetric matrix. Therefore, we only compute the upper right triangle and copy the remaining elements. Since  $\mathbf{X}_k$  is the filtered input data of band  $k$ , packed to 8 bits, the implementation makes use of SIMD instructions to improve the performance significantly. The computation is implemented concurrently by splitting the upper right part of the output matrix among all processing units.

3) *Whitening*: Whitening consists of two MMMs, as described in Eq. (2). Based on the quantization scheme described in Section II-2, the first multiplication is computed in 16 bit, and the second in 32 bit. For the first multiplication, we use 2-way SIMD instructions. We use the concurrent implementation found in the DSP library [35] for PULP, where each core computes a part of the matrix.

4) *Matrix Logarithm*: For computing the EVD, we implement both the basic version of Householder transformation and the improved version [32] for speedup analyses. The computation of the rotation matrix required for the Givens rotation [36] of each QR step is done exclusively with multiplications, divisions, and additions, without using expensive trigonometric functions [37]. For parallel implementation, every core is assigned with a frequency band and computes the Householder transformation and QR algorithm.

5) *Support Vector Machine (SVM)*: The matrix-vector product of the SVM is computed using 8-bit SIMD instructions. We implement it on a single core, since it accounts for a negligible portion of the computation of the entire model.

### IV. EXPERIMENTAL RESULTS AND DISCUSSION

We apply our methods on the BCI Competition IV-2a dataset [21] with 22 EEG channels and 4 MI classes from 9 different subjects. There are 288 trials for each of the training and testing sets. Each trial lasts 6 s and is sampled at 250 Hz.

Table I reports the classification accuracy of our proposed models compared to related work with different MRC configurations and EEGNET. MRC can be scaled to use more or fewer frequency bands and temporal windows. Hersche et al. [27] have shown that  $f = 43$  frequency bands and a single temporal window  $t = 1$  can already achieve comparable accuracy (74.8% on average) to the full MRC (75.5%) while requiring  $3\times$  fewer features. In this work, we use only one temporal window  $t = 1$  of 3.5 s and further scale down the number of frequency bands. Our results show that with  $2.4\times$  less frequency bands, i.e.  $f = 18$ , of bandwidth 2 Hz between 4 and 40 Hz, our full precision model achieves slightly higher accuracy by introducing the regularization with the hyperparameter  $\rho = 1$ . Comparing to EEGNET, which is known to be a compact CNN for BMI applications [17],

our full precision MRC is 3.8% more accurate. Regarding the quantization, EEGNET can be quantized down to 8-bit precision for the entire network with Q-EEGNET [24] without significant loss in accuracy (0.4%). However, our proposed mixed-precision MRC is still 3.2% more accurate. The minimal loss in accuracy of 1% from full to mixed-precision can be attributed mainly to the quantization at the input of the matrix logarithm. Regarding the memory footprint, Q-EEGNET requires 68.15 kB, while our MRC implementation uses approximately 84 kB, i.e.  $2 \cdot 22 \cdot 876$  for 8-bit input and output of IIR filters,  $18 \cdot (22+1) \cdot 22/2$  for  $W_k$  in 32-bits and reused for  $L_k$ ,  $18 \cdot (22+1) \cdot 22/2$  for the model parameters  $C_{ref,k}^{-1/2}$  in 16 bits, and 4554.4 for SVM weights in 8 bits.

Table II shows the computation time and the performance impact of the optimizations and Fig. 3 depicts the measured power trace. The first 18 peaks are measured when the frequency bands are calculated using 8 cores, framed with blue dashed line. The IIR filter implementation achieves 3.77 MACs per cycle with  $7.26 \times$  parallel speedup. Here, each output sample requires 10 MACs, 3 shuffle operations, and 4 bit-shifts, resulting in a theoretical maximum of 5 MACs per cycle. The covariance matrix computation reaches 8.14 MACs per cycle with concurrent execution yielding a speedup of  $7.10 \times$  using 8 cores. The parallel speedup of the Whitening is  $4.98 \times$  due to the parallelization overhead that is more visible with smaller matrix sizes (here  $22 \times 22$ ). However, it is not the bottleneck part of the MRC. The improvements of the Householder transformation have a significant impact on the performance yielding a speedup of  $3.6 \times$  on the computation of the matrix logarithm compared to the baseline, while the parallel speedup is  $5.67 \times$  compared to the single core computation and  $20.64 \times$  compared to the baseline. 18 matrix logarithms are computed, distributed to the 8 cores on a first-come first-served schedule, i.e. twice 8 matrix logarithms are computed on 8 cores, then the remaining 2 on two cores, as reflected on the power trace, framed with red dashdotted line. This workload unbalance contributes negatively to the parallel speedup. However, the performance would not increase significantly with a more balanced distribution since the ideal speedup would be  $6 \times$  with six parallel cores. Moreover, the maximal number of Floating Point Operations (FLOPs) per cycle is 2, of which we reach 1.69, limited by the iteratively computed divisions and square root operations. Finally, the SVM accounts for a minimal part of the execution with 0.15 ms, highlighted with green frame in Fig. 3. For comparison, the embedded

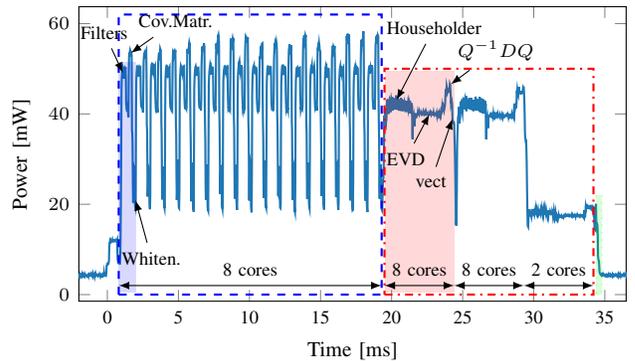


Fig. 3: End-to-end power measurement. The colors match the compute blocks in Fig. 1 explained in Sec. III.

BMI in [26] consumes 0.7 W and takes around 0.4 s, more than an order of magnitude more in terms of both, power consumption and execution time—or two orders of magnitude worse in terms of energy efficiency. We also compare to the Q-EEGNET implementation in [24] that is publicly available. We run both Q-EEGNET and MRC on Mr. Wolf at 100 MHz and 1.1 V. The former takes 13.64 ms consuming 0.678 mJ while the runtime of MRC lays within the same order of magnitude with 33.39 ms and consumes 1.304 mJ. It is up to the user to decide on the trade-off between accuracy and cost depending on the application scenario.

## V. CONCLUSION

This paper presents an improved MRC with reduced model size while keeping comparable accuracy (75.1% vs. 75.5% [27]), allowing accurate low-power embedded BMI. We further scale down the model by quantizing and proposing a mixed-precision implementation yielding a minimal accuracy loss of 1%, which is still 3.2% more accurate than the SoA embedded CNN for BMI named Q-EEGNET [24]. We propose a parallel implementation on a low-power MCU called Mr. Wolf, which takes only 33.39 ms and consumes 1.304 mJ. The higher accuracy compared to Q-EEGNET comes at the cost of a  $2.4 \times$  longer execution time and a  $1.9 \times$  higher energy consumption. However, it is still two orders of magnitude more energy efficient than other embedded solutions [26]. We provide an insight on accuracy-cost trade-off for embedded BMI models with actual implementation and measurements.

TABLE I: Classification accuracy (%) on 4-class MI.

Ref.	Q-EEGNET		MRC			
	[24] full	[24] 8-bit	[27] <sup>†</sup> full	[27] <sup>‡</sup> full	Ours <sup>◊</sup> full	Ours <sup>◊</sup> mixed
$t/f/\rho$			3/43/0	1/43/0	1/18/1	1/18/1
Subj. 1	81.0	81.0	90.0	91.8	91.8	90.7
Subj. 2	57.6	53.1	55.5	51.6	53.7	51.2
Subj. 3	87.9	91.2	81.3	83.5	83.5	81.0
Subj. 4	61.6	58.1	71.9	73.3	73.7	74.1
Subj. 5	70.6	68.4	69.6	63.4	68.8	63.0
Subj. 6	53.4	50.1	56.7	58.6	56.7	56.3
Subj. 7	75.7	75.2	85.6	86.7	84.1	58.9
Subj. 8	77.4	81.2	83.8	81.6	81.5	82.7
Subj. 9	76.7	79.7	84.9	82.6	82.2	81.8
Avg. Acc.	71.3	70.9	75.5	74.8	<b>75.1</b>	<b>74.1</b>
Std.	11.5	14.3	12.8	13.9	12.2	13.2

TABLE II: Computation time for MRC on Mr. Wolf with a frequency of 100 MHz at 1.1 V.

	baseline	improved EVD	concurrent	parallel speedup	ops/c <sup>†</sup>
Filter	66.67 ms	66.67 ms	9.18 ms	7.26	3.77
Cov. matrix	34.80 ms	34.80 ms	4.90 ms	7.10	8.14
Whitening	24.29 ms	24.29 ms	4.88 ms	4.98	0.79
Matrix logm.	309.76 ms	85.18 ms	15.01 ms	5.67	1.69
SVM	0.15 ms	0.15 ms	0.15 ms	-	1.25
Total	439.48 ms	206.93 ms	<b>33.39 ms</b>		
MACs/cycle <sup>‡</sup>	0.62	0.62	<b>4.11</b>		
FLOPs/cycle <sup>◊</sup>	0.08	0.30	<b>1.69</b>		
insn/cycle	0.907	0.837	0.788		

<sup>‡</sup> Number of fixed-point MACs over number of cycles w/o matrix logarithms.

<sup>◊</sup> Number of FLOPs over number of cycles during matrix logarithms.

<sup>†</sup> MACs or FLOPs per cycle for the concurrent implementation except SVM.

## REFERENCES

- [1] K. Koizumi, K. Ueda *et al.*, “Development of a cognitive brain-machine interface based on a visual imagery method,” in *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 2018, pp. 1062–1065.
- [2] Y. Yu, Z. Zhou *et al.*, “Self-Paced Operation of a Wheelchair Based on a Hybrid Brain-Computer Interface Combining Motor Imagery and P300 Potential,” *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 12, pp. 2516–2526, 12 2017.
- [3] M. Xiong, R. Hotter *et al.*, “A low-cost, semi-autonomous wheelchair controlled by motor imagery and jaw muscle activation,” in *2019 IEEE International Conference on Systems, Man and Cybernetics (SMC)*. IEEE, 2019, pp. 2180–2185.
- [4] A. A. Frolov, O. Mokienco *et al.*, “Post-stroke Rehabilitation Training with a Motor-Imagery-Based Brain-Computer Interface (BCI)-Controlled Hand Exoskeleton: A Randomized Controlled Multicenter Trial.” *Frontiers in neuroscience*, vol. 11, p. 400, 2017.
- [5] N. Kobayashi and M. Nakagawa, “BCI-based control of electric wheelchair using fractal characteristics of EEG,” *IEEJ Tran. on Electrical and Electronic Engineering*, vol. 13, no. 12, pp. 1795–1803, 2018.
- [6] J. León, J. J. Escobar *et al.*, “Deep learning for eeg-based motor imagery classification: Accuracy-cost trade-off,” *PLOS ONE*, vol. 15, no. 6, pp. 1–30, 06 2020.
- [7] H. Wu, Y. Niu *et al.*, “A parallel multiscale filter bank convolutional neural networks for motor imagery eeg classification,” *Frontiers in Neuroscience*, vol. 13, p. 1275, 2019.
- [8] R. T. Schirmeister, J. T. Springenberg *et al.*, “Deep learning with convolutional neural networks for EEG decoding and visualization,” *Human Brain Mapping*, vol. 38, no. 11, pp. 5391–5420, 2017.
- [9] F. Lotte and Cuntai Guan, “Regularizing Common Spatial Patterns to Improve BCI Designs: Unified Theory and New Algorithms,” *IEEE Transactions on Biomedical Engineering*, vol. 58, no. 2, pp. 355–362, 2 2011.
- [10] Kai Keng Ang, Zhang Yang Chin *et al.*, “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.
- [11] C. H. Nguyen and P. Artemiadis, “Eeg feature descriptors and discriminant analysis under riemannian manifold perspective,” *Neurocomputing*, vol. 275, pp. 1871 – 1883, 2018.
- [12] S. Kumar, F. Yger *et al.*, “Towards adaptive classification using riemannian geometry approaches in brain-computer interfaces,” in *2019 7th International Winter Conference on Brain-Computer Interface (BCI)*, 2019, pp. 1–6.
- [13] M. Congedo, A. Barachant *et al.*, “Riemannian geometry for eeg-based brain-computer interfaces; a primer and a review,” *Brain-Computer Interfaces*, vol. 4, pp. 1–20, 03 2017.
- [14] F. Yger, M. Berar *et al.*, “Riemannian Approaches in Brain-Computer Interfaces: A Review,” *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 10, pp. 1753–1762, 10 2017.
- [15] S. Chevallier, E. Kalunga *et al.*, “Review of riemannian distances and divergences, applied to ssvp-based bci,” *Neuroinformatics*, 06 2020.
- [16] X. Wang, M. Magno *et al.*, “FANN-on-MCU: An Open-Source Toolkit for Energy-Efficient Neural Network Inference at the Edge of the Internet of Things,” *IEEE Internet of Things Journal*, 2020.
- [17] V. J. Lawhern, A. J. Solon *et al.*, “EEGNet: a compact convolutional neural network for EEG-based brain-computer interfaces,” *Journal of Neural Engineering*, vol. 15, no. 5, p. 056013, 2018.
- [18] A. Pullini, D. Rossi *et al.*, “Mr.Wolf: An Energy-Precision Scalable Parallel Ultra Low Power SoC for IoT Edge Processing,” *IEEE Journal of Solid-State Circuits*, vol. 54, no. 7, pp. 1970–1981, 2019.
- [19] T. Ingolfsson, M. Hersche *et al.*, “Eeg-tcnet: An accurate temporal convolutional network for embedded motor-imagery brain-machine interfaces,” *arXiv:2006.00622*, 05 2020.
- [20] M.-A. Li, J.-F. Han *et al.*, “A Novel MI-EEG Imaging With the Location Information of Electrodes,” *IEEE Access*, vol. 8, pp. 3197–3211, 2020.
- [21] C. Brunner, R. Leeb *et al.*, “BCI competition 2008 - Graz data set A,” <http://bnci-horizon-2020.eu/database/data-sets>.
- [22] Y. Zhao, S. Yao *et al.*, “On the improvement of classifying EEG recordings using neural networks,” in *Proc. IEEE Big Data*, Dec. 2017, pp. 1709–1711.
- [23] H. Wu, Y. Niu *et al.*, “A Parallel Multiscale Filter Bank Convolutional Neural Networks for Motor Imagery EEG Classification,” *Frontiers in Neuroscience*, vol. 13, Nov. 2019.
- [24] T. Schneider, X. Wang *et al.*, “Q-EEGNet: an Energy-Efficient 8-bit Quantized Parallel EEGNet Implementation for Edge Motor-Imagery Brain-Machine Interfaces,” *arXiv:2004.11690v1*, Apr. 2020.
- [25] X. Wang, M. Hersche *et al.*, “An accurate eegnet-based motor-imagery brain-computer interface for low-power edge computing,” in *2020 IEEE International Symposium on Medical Measurements and Applications (MeMeA)*, 2020, pp. 1–6.
- [26] K. Belwafi, O. Romain *et al.*, “An embedded implementation based on adaptive filter bank for brain-computer interface systems,” *Journal of Neuroscience Methods*, 2018.
- [27] M. Hersche, T. Rellstab *et al.*, “Fast and Accurate Multiclass Inference for MI-BCIs Using Large Multiscale Temporal and Spectral Features,” in *2018 26th European Signal Processing Conference (EUSIPCO)*. IEEE, 9 2018, pp. 1690–1694.
- [28] P. Yang, J. Wang *et al.*, “Mlp with riemannian covariance for motor imagery based eeg analysis,” *IEEE Access*, vol. 8, pp. 139 974–139 982, 2020.
- [29] P. L. C. Rodrigues, C. Jutten *et al.*, “Riemannian procrustes analysis: Transfer learning for brain-computer interfaces,” *IEEE Transactions on Biomedical Engineering*, vol. 66, no. 8, pp. 2390–2401, 2019.
- [30] F. Yger, M. Berar *et al.*, “Riemannian approaches in brain-computer interfaces: a review,” *IEEE Transactions on Neural Systems and Rehabilitation Engineering*, vol. 25, no. 10, pp. 1753–1762, 2016.
- [31] J. O. Smith, *Introduction to Digital Filters with Audio Applications*. <http://ccrma.stanford.edu/jos/filters/>, 2020, online book.
- [32] R. Burden and J. Faires, *Numerical analysis*. Cengage Learning, 2004.
- [33] J. H. Wilkinson, *The algebraic eigenvalue problem*. Oxford Clarendon, 1965, vol. 662.
- [34] M. Gautschi, P. D. Schiavone *et al.*, “Near-threshold RISC-V core with DSP extensions for scalable IoT endpoint devices,” *IEEE Transactions on VLSI Systems*, vol. 25, no. 10, pp. 2700–2713, 2017.
- [35] X. Wang, “DSP library for PULP,” <https://github.com/pulp-platform/pulp-dsp>, 2019.
- [36] W. Givens, “Numerical computation of the characteristic values of a real symmetric matrix,” Oak Ridge National Lab., Tech. Rep., 1954.
- [37] D. Bindel, J. Demmel *et al.*, “On computing givens rotations reliably and efficiently,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 28, no. 2, pp. 206–238, 2002.