

# Challenges and Opportunities of Edge AI for Next-Generation Implantable BMIs

MohammadAli Shaeri<sup>1</sup>, Arshia Afzal<sup>1,2</sup>, Mahsa Shoaran<sup>1</sup>

<sup>1</sup>Institute of Electrical and Micro Engineering, Center for Neuroprosthetics, EPFL, 1202 Geneva, Switzerland

<sup>2</sup>Department of Electrical Engineering, Sharif University of Technology, Tehran, Iran

Email: {mohammad.shaeri, arshia.afzal, mahsa.shoaran}@epfl.ch

**Abstract**—Neuroscience and neurotechnology are currently being revolutionized by artificial intelligence (AI) and machine learning. AI is widely used to study and interpret neural signals (analytical applications), assist people with disabilities (prosthetic applications), and treat underlying neurological symptoms (therapeutic applications). In this brief, we will review the emerging opportunities of on-chip AI for the next-generation implantable brain machine interfaces (BMIs), with a focus on state-of-the-art prosthetic BMIs. Major technological challenges for the effectiveness of AI models will be discussed. Finally, we will present algorithmic and IC design solutions to enable a new generation of AI-enhanced and high-channel-count BMIs.

**Index Terms**—Artificial Intelligence (AI), Machine Learning (ML), Brain Machine Interface (BMI), hardware efficiency.

## I. INTRODUCTION

Globally, millions of people suffer from severe motor disabilities such as paralysis and stroke. In order to bring disabled people back to their normal lives, a wide variety of brain-machine interfaces (BMIs) are being developed at the cutting edge of neurotechnology and neuroscience. In general, BMIs are considered as systems that close the loop from sensing to action (e.g., from vision/touch to reach/grasp), as shown in Fig. 1. To realize this goal, an implantable BMI records one or more types of neural signals from the brain. Considering the trade-off between spatiotemporal resolution and invasiveness, the intracortical and cortical recordings of brain activity are widely used in BMI applications [1], [2]. Next, data processing and artificial intelligence (AI) techniques can be used to extract task-relevant informative content in the form of a movement intention or a marker of brain malfunction (e.g., brain injury). Finally, the extracted information is used to generate an actuation command to move an artificial or natural limb, or a stimulation command to modulate the brain activity.

From the standpoint of application, BMIs can be categorized into analytical, prosthetic, and therapeutic systems (Fig. 1). *Analytical BMIs* are utilized to study brain activity, function, or connectivity. Thanks to the recent success of AI in analyzing high-dimensional data, it is widely used in such studies ranging from cell-level (e.g., spike sorting) to cognitive-level (e.g., neural coding). The research goal is to discover the brain mechanism or dynamics underlying sensory perception, or to uncover brain intention for a specific action. *Prosthetic BMIs* allow subjects to perform daily tasks such as movement [1] or typing [2]. Such BMIs can employ corticomotor activities to control natural limbs via neuromuscular or spinal cord stimulation. Another type of prosthetic BMIs stimulates the somatosensory cortex in order to restore sensory feedback [1]. Through modulation of the nervous system, *therapeutic BMIs*

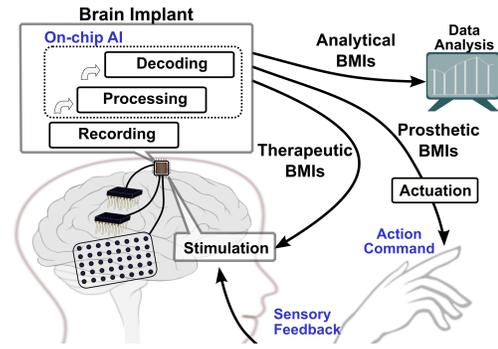


Fig. 1. Schematic view of an implantable BMI. In all three BMI applications, neural signals are initially recorded and processed (e.g., amplified, digitized, filtered) and subsequently decoded (either on- or off-implant). The output can be an actuation or neuromodulatory command for prosthetic and therapeutic applications, respectively.

aim at restoring lost brain functions and treating symptoms (e.g., memory enhancement, seizure or pain suppression).

In a conventional BMI, the recorded neural signals are transmitted to an external device for further processing, while an ‘on-implant’ AI unit performs this function in modern BMIs (Fig. 1). Future-generation BMIs will be miniaturized and ubiquitous prostheses, applicable to daily tasks and chronic use. Thus, designing hardware-efficient and compact implantable BMIs is of crucial importance. Such systems should obtain a high accuracy to be reliable for various prosthetic applications. In this paper, we will first discuss the emerging applications and critical challenges for designing modern implantable BMIs. Next, we will review the state-of-the-art BMI system-on-chips (SoCs) with on-chip AI. Finally, we will discuss novel algorithmic and circuit-level solutions to realize next-generation high-density and intelligent implantable BMIs.

## II. MODERN BMIs WITH INTEGRATED AI

The recording capacity of implantable BMIs has grown rapidly over years, promising higher levels of proficiency and performance. For instance, Neuralink and Paradromics develop implantable BMIs with thousands of channels. Yet, local on-implant processing remains a challenge in many BMI directions. Developing high-performance, energy-efficient, and scalable AI techniques could enable a new generation of implantable BMIs with minimal need for data transmission, enhanced security and privacy, and higher independence. To achieve this goal, critical challenges at the algorithm and circuit levels must be addressed, as discussed below.

### A. BMI Effectiveness Metrics and Design Challenges

The key dimensions for the effectiveness of AI models in the context of a BMI include (I) accuracy, (II) robustness,

TABLE I  
PERFORMANCE SUMMARY AND COMPARISON OF THE STATE-OF-THE-ART BMI SoCs WITH ON-CHIP AI.

Parameter	[3]	[4]	[5]	[6]	[7]	[8]	[9]	[10]
Task	Movement classification	Movement classification	Spike sorting	Spike sorting	Spike sorting	Spike sorting	Spike sorting	Movement classification
Input Signal	Spiking train	Spiking rate	Spike waveform	Spike waveform	Spike waveform	Spike waveform	Spike waveform	ECoG
Classification Model	SNN	ELM	K-means	BOTM	K-means	K-means	K-means	NeuralTree
Features	–	–	Symmlt-2 Wavelet	–	Filtering	Spike derivative	Adaptive filtering	Multi-band filters, LMP
Online Training	N	N	Y	N	N	Y	Y	N
Closed-loop	Y	N	Y	N	N	N	N	Y
# Recording Channels	16	128	10	16	128	N/A	16	256/64
CMOS Process (nm)	180	350	130	40	65	180	22	65
Sampling Rate (kS/s)	1	0.05	N/A	24	25	N/A	20	2
Resolution (bits)	–	–	16	12	8	–	9	10
AI Area/ch ( $mm^2$ )	3.21	0.191*	0.08	0.0175	0.003	1.023	0.014	0.0187 <sup>Ⓒ</sup>
AI Power/ch ( $\mu W$ )	250	0.0032*	56.9	19	0.175	4.35	2.79	4.23 <sup>Ⓒ</sup>
System Area/ch ( $mm^2$ )	N/A	N/A	N/A	N/A	N/A	1.4 <sup>†</sup>	0.038 <sup>‡</sup>	0.013 <sup>§</sup>
System Power/ch ( $\mu W$ )	N/A	N/A	N/A	N/A	N/A	10.25 <sup>†</sup>	4.31 <sup>†</sup>	7.08 <sup>§</sup>
Accuracy (%)	N/A	99.3	N/A	93.4	72-86	93.2	94.12	71.5-75

\* Power and area of the on-chip part (i.e. the hidden layer of ELM), divided by the number of channels.

† Power and area of the analog front-end and spike detector/sorter, divided by the number of channels.

‡ Power and area of the analog front-end, ADC, and spike detector/sorter, divided by the number of channels.

Ⓒ Power and area of the feature extractor and decoder, divided by the number of active (i.e., selected) channels during inference (64).

§ Area of the front-end, feature extractor, decoder, and 16-channel stimulator, divided by the number of recording and stimulation channels (256+16).

¶ Power of the mixed-signal front-end, feature extractor, and decoder, divided by the number of active channels (64).

(III) online adaptability, (IV) interpretability, (V) computational speed, (VI) scalability, and (VII) hardware efficiency.

The decoding performance—often measured as classification or regression accuracy—needs to be high for clinically viable BMIs. AI models are susceptible to performance loss due to signal variability (e.g., noise). Therefore, the model needs to be robust to handle variations and guarantee reliability over time. In addition, online adaptability enables the model to adjust to non-stationary signal changes in chronic settings. Another emerging AI trend is to develop interpretable models that establish an explainable relationship between brain activity and the associated physical phenomena. Furthermore, BMI systems are desired to make predictions in real-time, with minimal time required for training. Thus, improving the training and inference speed of edge AI models is critical.

Enhancing the recording capacity of modern BMIs will significantly increase the data dimensionality. Thus, an emerging challenge is the processing and decoding of neural data in high-channel-count BMIs with inevitably limited hardware resources and under strict power requirements near neural tissue. Thus, the hardware scalability of the model, i.e., its capacity to handle high-dimensional data without a significant increase in hardware cost (power, chip area) is critical for implantable BMIs. Increasing the neural data dimensionality may also increase the risk of model overfitting.

### B. State-of-the-art BMI SoCs

AI-based methods have been utilized for spike sorting, detection of neurological symptoms (e.g., epileptic seizures [11]), and motor intention decoding in a number of neural interface SoCs. Here, we focus on AI methods used for spike sorting and/or motor decoding in prosthetic BMIs.

The AI models used in the BMI domain generally aim at spike sorting and movement classification. A classic approach to solve a classification/clustering problem is to allocate each data point to the neighboring class with minimal distance. In such methods, the distance/proximity metric could be the Manhattan ( $l_1$ -norm) distance [7], [5] or cosine similarity [12]. K-means [12], [8] and template matching [6] are the distance-based methods widely used for spike sorting (Fig. 2(a)). For

instance, wavelet features were extracted and used to cluster spikes based on the  $l_1$ -norm distance in [5]. The chip consumed  $56.9\mu W/ch$  and occupied a silicon area of  $0.08mm^2/ch$ . In [6], local extrema were detected from neural signal and their adjacent samples were selected and classified using the *Bayes optimal template matching (BOTM)*. This spike sorter consumed  $19.0\mu W/ch$  and  $0.0175mm^2/ch$ .

An integer-coefficient filter was used for feature extraction from spikes in [7], followed by feature selection and clustering with a simplified K-means algorithm. Thanks to the dimensionality reduction and ultra-low-voltage SRAM usage, the fabricated spike sorter consumed  $0.175\mu W/ch$  and  $0.003mm^2/ch$ , albeit at the cost of a degraded accuracy (76-86%). Similarly, adaptive FIR filtering was used for feature extraction in [9], followed by feature selection. Then,  $l_1$ -norm distance in the feature space was employed for spike sorting. The design was fabricated in a 22nm CMOS process ( $2.79\mu W/ch$  and  $0.014mm^2/ch$ ). More recently, an analog implementation of the *first and second derivative extrema (FSDE)* for feature extraction as well as K-means were reported for spike sorting ( $4.35\mu W/ch$ ,  $1.023mm^2/ch$ ) [8].

Hardware-efficient classification models such as *window discrimination (WD)* have also been reported for BMIs, where a hyperrectangle discrimination window is assigned to each individual class (Fig. 2(b)) [13]. A discrimination window is composed of two decision boundaries for each data dimension, simply implemented by a few digital comparators. Similarly, *decision tree (DT)-based* models classify the data with a set of successive comparisons and achieve excellent energy efficiencies [11], Fig. 2(c). Rather than using a single feature per node as in conventional axis-parallel DTs, an *oblique decision tree (OT)* uses a linear combination of features per node, thus forming a more accurate, oblique boundary for movement classification [14] or spike sorting [15] (Fig. 2(d)).

*Neural networks-based* models have also been used in BMI applications [3], [4], [16]. In [3], a brain-inspired *spiking neural network (SNN)* was implemented to classify stimulation-evoked brain activity. The SNN classifier contained integrate-and-fire interconnected neurons that mimic

the micro-scale property of the neuronal network. The chip consumed  $250\mu\text{W}/\text{ch}$  and  $3.213\text{mm}^2/\text{ch}$ . In [4], an *extreme learning machine (ELM)* was employed to classify finger movements in monkeys. The ELM is a single hidden-layer feedforward network that performs a random projection of the inputs through a nonlinear hidden layer and generates the classification results through a linear output layer. The nonlinear hidden layer was implemented on-chip, while the linear output layer was implemented off-chip via a commercial microcontroller. The chip achieved an excellent power consumption and silicon area of  $0.0032\mu\text{W}/\text{ch}$  and  $0.191\text{mm}^2/\text{ch}$ , respectively. Alternatively, a *binarized neural network (BNN)* with binary weights and activation function was reported as a hardware-efficient model in [17]. Moreover, combing a DT structure with a highly-pruned neural network led to a lightweight *NeuralTree* model for finger movement classification from human ECoG [10]. The model comprised sparsely-connected internal nodes with fewer computations and memory needs compared to conventional OTs. Fabricated in a 65nm process, the on-chip decoder consumed  $4.23\mu\text{W}/\text{ch}$  and  $0.0187\text{mm}^2/\text{ch}$ .

In addition to the classifiers discussed above, a number of AI methods seek efficient, minimal data representation to reduce hardware complexity. Namely, *salient feature selection* selects a small number of salient features that achieve the highest class discrimination from other classes for on-implant spike sorting [13], thus improving the hardware efficiency. Table I summarizes the details of the state-of-the-art BMIs with on-chip AI. Designs with spike sorting only, or spiking rate as input, may need additional blocks for movement classification and spiking rate extraction, respectively.

### III. AI ALGORITHM AND HARDWARE SOLUTIONS FOR NEXT-GENERATION BMIS

Neuronal spiking rate is commonly used as the input to intracortical BMIs, requiring a complex spike sorting phase. Recent studies, however, show that simple threshold-crossing rate (without sorting) can lead to successful decoding of motor intentions. Although the decoding accuracy with spike sorting is generally higher, the low complexity threshold-crossing method or extraction of *spiking band power* (i.e., the power within 0.3–1kHz) can be beneficial in certain applications [18].

To date, various hardware-algorithm co-design solutions have been introduced to shrink the AI models and improve hardware efficiency, such as network pruning and weight quantization [15], [14]. Similarly, SNN models represent the data with a binary stream of spiking events and facilitate hardware implementation of NNs.

A recent AI trend is to implement the training algorithm in tandem with the inference model on chip [9], [8], [19]. Such online (i.e., adaptive) ML approaches could enable autonomous BMIs with minimal need for recalibration, at the cost of extra hardware resources required for on-chip training.

#### A. Toward Kinetic Trajectory Decoding in Implantable BMIS

While AI models have been implemented for discrete movement classification in BMIs, no SoC for continuous trajectory decoding has been reported so far (Table I). This could be due to the higher complexity of accurate regression

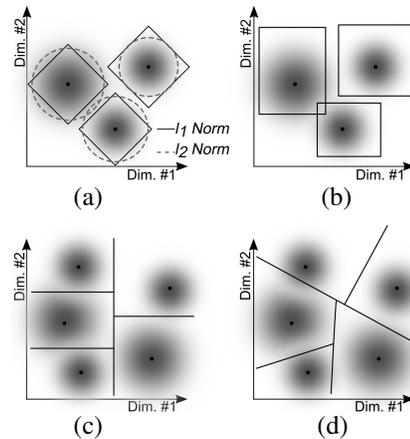


Fig. 2. Hardware-efficient classification models; (a) Distance-based classification, (b) Window discrimination, (c) Axis-parallel DT (d) Oblique DT. models compared to binary or multi-class classifiers. Given the crucial role of continuous trajectory decoding to control prosthetic BMIs, here we introduce an algorithmic solution that transforms the regression problem into a classification task, with key advantages over conventional regressors.

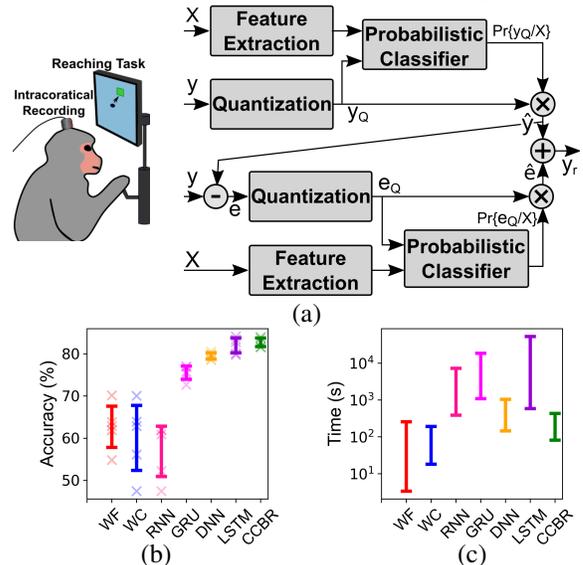


Fig. 3. (a) Experimental setup and block diagram of the proposed CCBR for trajectory decoding.  $X$ ,  $y$ , and  $e$  indicate the neural data, movement trajectory, and error signal. (b) Prediction accuracy (in  $R^2$ ) for CCBR, Wiener filter (WF), Wiener cascade (WC), RNN, gated recurrent unit (GRU), and LSTM. Error bars represent mean  $\pm$  STD across validation folds. (c) Runtime comparison for tuning different regression models vs. CCBR.

Traditionally, accurate regression models (e.g., LSTM) are used to decode movement trajectory from neural signals [20]. However, for realizing implantable BMIs, the accuracy, training speed, and hardware complexity need to be considered simultaneously. Here, we transform the regression problem into a classification task by quantizing the kinematic signal (e.g., velocity) and solving a multi-class classification problem. We used PCA for dimensionality reduction, and a simple probabilistic SVM as the classifier. The model reconstructs the kinematic signal by multiplying the kinematic quanta (i.e., class labels) by the corresponding probabilities. In the second phase of training, we run a similar classifier to predict the reconstruction error. The proposed *cascaded classification-*

based regressor (CCBR) generates the trajectory output by adding the predicted movement and error signals (Fig. 3(a)). This model can decode the first-order along with higher-order errors and converge to maximum performance with no need for hyperparameter tuning, thus results in fast training.

*Low Sensitivity to Hyperparameters:* Reducing the sensitivity of the decoder to hyperparameters is crucial to improve the reliability and lower the need for retraining in an implantable BMI [21]. We tested the performance of CCBR on a monkey intracortical dataset recorded from dorsal premotor cortex during a reaching task [21], by varying various hyperparameters, including the number of principal components (PCs), SVM’s regularization factor (C), and quantization level (QL) of the movement signal ( $y$ ). We observed that selecting the proper number of PCs, which indicates the input dimensionality, is sufficient to maximize the accuracy. Thus, the model performance is highly robust to hyperparameters (C and QL).

*Prediction Accuracy and Training Speed:* As shown in Fig. 3(b), compared to various decoders reported in [21], CCBR achieved the highest accuracy (82.7%) and one of the lowest variances (1%), proving its high reliability for neural decoding tasks. Although complex AI models (e.g., LSTM) may result in high performance, they often require a considerable training time. Thanks to the robustness of our proposed technique, we can replace the traditional methods for hyperparameter tuning (e.g., grid and random search) by error prediction and significantly reduce the model training time (Fig. 3(c)).

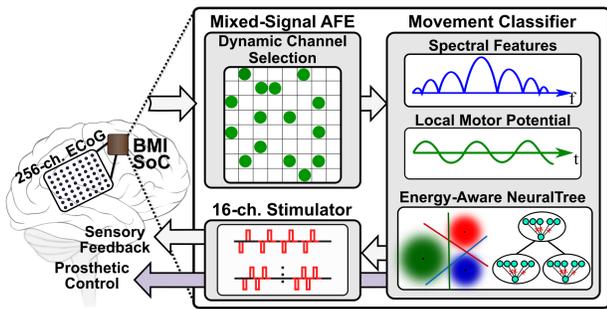


Fig. 4. Block diagram of a high-density BMI with on-chip AI.

### B. High-Density BMI Design with On-Chip AI

While the existing BMI SoCs have integrated a limited number of channels (and often lack the integration of AI with the analog front-end), high-density sensing is critical to improve the motor decoding accuracy for next-generation BMIs. One such design [10] implements a 256-channel highly-scalable SoC for low-noise ECoG/LFP sensing, multi-class movement classification, and closed-loop stimulation. The key idea is the use of a highly-pruned energy-aware AI model compatible with a dynamic channel selection scheme, and a highly-multiplexed and low-power mixed-signal front-end for successive processing of selected channels upon training. Figure 4 depicts the architecture of this high-density BMI. The fully-integrated SoC occupies a silicon area of  $0.013\text{mm}^2/\text{ch}$  in 65nm process, with state-of-the-art power consumption and energy efficiency of  $7.08\mu\text{W}/\text{ch}$  and  $0.227\mu\text{J}/\text{class}$  for the entire BMI system, respectively.

## IV. CONCLUSION

Next-Generation BMIs will be implantable prostheses with on-chip AI. We discussed the algorithmic and hardware challenges for realizing implantable BMIs and reviewed the state-of-the-art AI models used for spike sorting and movement classification. Finally, we introduced CCBR as a potential solution for trajectory decoding in Next-Generation BMIs.

## REFERENCES

- [1] S. N. Flesher *et al.*, “A brain-computer interface that evokes tactile sensations improves robotic arm control,” *Science*, vol. 372, no. 6544, pp. 831–836, 2021.
- [2] F. R. Willett *et al.*, “High-performance brain-to-text communication via handwriting,” *Nature*, vol. 593, no. 7858, pp. 249–254, 2021.
- [3] F. Boi *et al.*, “A bidirectional brain-machine interface featuring a neuromorphic hardware decoder,” *Front. Neurosci.*, vol. 10, p. 563, 2016.
- [4] Y. Chen *et al.*, “A 128-channel extreme learning machine-based neural decoder for brain machine interfaces,” *IEEE Trans. Biomed. Circuits Syst.*, vol. 10, no. 3, pp. 679–692, June 2016.
- [5] G. Gagnon-Turcotte *et al.*, “A wireless electro-optic headstage with a  $0.13\text{-}\mu\text{m}$  cmos custom integrated DWT neural signal decoder for closed-loop optogenetics,” *IEEE Trans. Biomed. Circuits Syst.*, vol. 13, no. 5, pp. 1036–1051, Oct 2019.
- [6] H. Xu *et al.*, “Unsupervised and real-time spike sorting chip for neural signal processing in hippocampal prosthesis,” *J. Neurosci. Meth.*, vol. 311, pp. 111–121, 2019.
- [7] A. T. Do *et al.*, “An area-efficient 128-channel spike sorting processor for real-time neural recording with  $0.175\mu\text{W}/\text{channel}$  in 65-nm cmos,” *IEEE Trans. VLSI Syst.*, vol. 27, no. 1, pp. 126–137, 2018.
- [8] H. Hao *et al.*, “A  $10.8\mu\text{W}$  neural signal recorder and processor with unsupervised analog classifier for spike sorting,” *IEEE Trans. Biomed. Circuits Syst.*, vol. 15, no. 2, pp. 351–364, 2021.
- [9] S. M. A. Zeinolabedin *et al.*, “A 16-channel fully configurable neural soc with  $1.52\text{W}/\text{Ch}$  signal acquisition,  $2.79\text{W}/\text{Ch}$  real-time spike classifier, and  $1.79\text{TOPS}/\text{W}$  deep neural network accelerator in 22nm FDSOI,” *IEEE Trans. Biomed. Circuits Syst.*, 2022.
- [10] U. Shin *et al.*, “A 256-channel  $0.227\mu\text{J}/\text{class}$  versatile brain activity classification and closed-loop neuromodulation soc with  $0.004\text{mm}^2\text{-}1.51\mu\text{W}/\text{channel}$  fast-settling highly multiplexed mixed-signal front-end,” in *IEEE Inter. Solid-State Circuits Conf.*, vol. 65, 2022, pp. 338–340.
- [11] B. Zhu *et al.*, “Closed-loop neural prostheses with on-chip intelligence: A review and a low-latency machine learning model for brain state detection,” *IEEE Trans. Biomed. Circuits Syst.*, 2021.
- [12] C. Seong *et al.*, “A multi-channel spike sorting processor with accurate clustering algorithm using convolutional autoencoder,” *IEEE Trans. Biomed. Circuits Syst.*, 2021.
- [13] M. Shaeri and A. M. Sodagar, “A framework for on-implant spike sorting based on salient feature selection,” *Nature Comm.*, vol. 11, no. 3278, pp. 1–9, June 2020.
- [14] B. Zhu *et al.*, “Resot: Resource-efficient oblique trees for neural signal classification,” *IEEE Trans. Biomed. Circuits Syst.*, vol. 14, no. 4, pp. 692–704, 2020.
- [15] Y. Yang *et al.*, “A hardware-efficient scalable spike sorting neural signal processor module for implantable high-channel-count brain machine interfaces,” *IEEE Trans. Biomed. Circuits Syst.*, vol. 11, no. 4, pp. 743–754, August 2017.
- [16] J. Yoo and M. Shoaran, “Neural interface systems with on-device computing: machine learning and neuromorphic architectures,” *Curr. Opin. Biotech.*, vol. 72, pp. 95–101, 2021.
- [17] D. Valencia and A. Alimohammad, “Neural spike sorting using binarized neural networks,” *IEEE Trans. Neur. Syst. Rehab. Eng.*, vol. 29, pp. 206–214, 2020.
- [18] S. R. Nason *et al.*, “A low-power band of neuronal spiking activity dominated by local single units improves the performance of brain-machine interfaces,” *Nat. Biom. Eng.*, vol. 4, no. 10, pp. 973–983, 2020.
- [19] S. Shaikh *et al.*, “Lightweight reinforcement algorithms for autonomous, scalable intra-cortical brain machine interfaces,” *bioRxiv*, 2020.
- [20] L. Yao *et al.*, “Fast and accurate decoding of finger movements from ECoG through riemannian features and modern machine learning techniques,” *J. Neur. Eng.*, vol. 19, no. 1, pp. 1–14, Feb 2022.
- [21] J. I. Glaser *et al.*, “Machine learning for neural decoding,” *Eneuro*, vol. 7, no. 4, 2020.