

Alma Mater Studiorum Università di Bologna
Archivio istituzionale della ricerca

A Bayesian-optimized design for an interpretable convolutional neural network to decode and analyze the P300 response in autism

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

Borra, D., Magosso, E., Castelo-Branco, M., Simoes, M. (2022). A Bayesian-optimized design for an interpretable convolutional neural network to decode and analyze the P300 response in autism. JOURNAL OF NEURAL ENGINEERING, 19(4), 1-26 [10.1088/1741-2552/ac7908].

Availability:

This version is available at: <https://hdl.handle.net/11585/899621> since: 2025-01-22

Published:

DOI: <http://doi.org/10.1088/1741-2552/ac7908>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

1 **A Bayesian-optimized design for an interpretable convolutional**
2 **neural network to decode and analyze the P300 response in**
3 **autism**

4 Davide Borra¹ (0000-0003-3791-8555)

5 Elisa Magosso^{1,2,3} (0000-0002-4673-2974)

6 Miguel Castelo-Branco⁴ (0000-0003-4364-6373)

7 Marco Simões^{4,5} (0000-0003-3713-2464)

8
9 *¹Department of Electrical, Electronic and Information Engineering “Guglielmo Marconi”*
10 *(DEI), University of Bologna, Cesena Campus, Cesena, Italy*

11 *²Alma Mater Research Institute for Human-Centered Artificial Intelligence, University of*
12 *Bologna, Bologna, Italy*

13 *³Interdepartmental Center for Industrial Research on Health Sciences & Technologies,*
14 *University of Bologna, Bologna, Italy*

15 *⁴CIBIT – Coimbra Institute for Biomedical Imaging and Translational Research, ICNAS –*
16 *Institute of Nuclear Sciences Applied to Health, University of Coimbra, Coimbra, Portugal*

17 *⁵CISUC – Center for Informatics and Systems, University of Coimbra, Coimbra, Portugal*

18
19 Word count: (max:12000)

20
21
22 **Correspondence**

23 Davide Borra

24 email: davide.borra2@unibo.it

1 **ABSTRACT (max: 300 words)**

2 *Objective.* P300 can be analyzed in autism spectrum disorder (ASD) to derive biomarkers and
3 can be decoded in BCIs to reinforce ASD impaired skills. Convolutional neural networks
4 (CNNs) have been proposed for P300 decoding, outperforming traditional algorithms but they
5 i) do not investigate optimal designs in different training conditions; ii) lack in interpretability.
6 To overcome these limitations, an interpretable CNN (ICNN), that we recently proposed for
7 motor decoding, has been modified and adopted here, with its optimal design searched via
8 Bayesian optimization.

9 *Approach.* The ICNN provides a straightforward interpretation of spectral and spatial features
10 learned to decode P300. The Bayesian-optimized (BO) ICNN design was investigated
11 separately for different training strategies (within-subject, within-session, and cross-subject)
12 and BO models were used for the subsequent analyses. Specifically, transfer learning (TL)
13 potentialities were investigated by assessing how pretrained cross-subject BO models
14 performed on a new subject vs. random-initialized models. Furthermore, within-subject BO-
15 derived models were combined with an Explanation Technique (ICNN+ET) to analyze P300
16 spectral and spatial features.

17 *Main results.* The ICNN resulted comparable or even outperformed existing CNNs, at the same
18 time being lighter. Bayesian-optimized ICNN designs differed depending on the training
19 strategy, needing more capacity as the training set variability increased. Furthermore, TL
20 provided higher performance than networks trained from scratch. The ICNN+ET analysis
21 suggested the frequency range [2, 5.8] Hz as the most relevant, and spatial features showed a
22 right-hemispheric parietal asymmetry. The ICNN+ET-derived features, but not ERP-derived
23 features, resulted significantly and highly correlated to ADOS clinical scores.

24 *Significance.* This study substantiates the idea that a CNN can be designed both accurate and
25 interpretable for P300 decoding, with an optimized design depending on the training condition.

1 The novel ICNN-based analysis tool was able to better capture ASD neural signatures than
2 traditional ERP analysis, possibly paving the way for identifying novel biomarkers.

3

1 **KEYWORDS**

2 Convolutional neural network, Electroencephalography, Autism Spectrum Disorder, P300,

3 Interpretability, Bayesian optimization

4

1 1. INTRODUCTION

2 Autism spectrum disorder (ASD) is a set of neurodevelopmental conditions with persistent
3 deficits in social communication and social interaction across multiple contexts, together with
4 restricted, repetitive patterns of behavior, interests, or activities [1]. ASD people show
5 difficulties in social-emotional reciprocity, in developing, maintaining, and understanding
6 relationships, and in non-verbal communicative behaviors used for social interactions, such as
7 joint attention [2–6]. Joint attention emerges during the first year of life and involves the non-
8 verbal coordination of attention of two individuals towards an object or event [7], playing an
9 important role in the development of social and language capabilities [8,9].

10 Approaches based on neuroimaging, e.g., diffusion tensor imaging and functional magnetic
11 resonance imaging, are used to characterize and identify potential neural biomarkers of
12 information-processing deficits in children with autism [10,11]. In addition, Event-Related
13 Potentials (ERPs) computed from the electroencephalogram (EEG) provide a less expensive
14 and portable way to study sensory information processing and are applied to study the neural
15 response in autism following an incoming stimulus [12]. Investigations concern alterations in
16 both early ERP components reflecting pre-attentive sensory processing and/or initial
17 orientation and capture of attention, such as P100 [13] and N100 [14], and later components,
18 such as P300 [15,16]. The P300 response is characterized by a positive deflection that occurs
19 while attending a target stimulus; it peaks between 250-500 ms after the stimulus onset and it
20 is mostly distributed on the scalp around midline electrode sites (Fz, Cz, Pz), increasing its
21 amplitude from frontal to parietal sites [17]. This response can be elicited in two-stimuli
22 oddball paradigms, where an infrequent (target) stimulus is immersed into a sequence of more
23 frequent (standard) stimulus. The amplitude of the P300 response was found to be positively
24 correlated with allocation of attention resources, stimulus recognition and updating of working
25 memory [18,19]. Abnormalities in the sensory information processing following an incoming

1 stimulus were found in ASD, as quantified by a reduced P300 amplitude in auditory and visuo-
2 spatial tasks compared to healthy subjects, reflecting deficiencies in cognitive, attentional, and
3 working memory processes [12,20–27].

4 Besides being potentially useful as an EEG-based ASD biomarker, the P300 response can
5 be used for ASD intervention. Indeed, Brain-Computer Interfaces (BCIs) proved to be useful
6 personalized therapeutic approaches in ASD [30–34]; these interfaces can be designed to train
7 ASD people via an EEG-based neurofeedback aimed to reinforce social interactions and
8 communication skills (e.g., joint attention [28,29]). In this scenario, the P300 response elicited
9 in visuo-spatial tasks represents an important control signal for the BCI system [28,29]. A
10 crucial stage of a BCI is represented by the decoding algorithm, that detects the P300 response
11 from the EEG and translates it into a command. Challenges to perform this step arise from the
12 noise sensitivity, non-linearity and non-stationarity of the EEG, as these characteristics depend
13 on the subject and on the environment [33]. In particular, the non-stationarity causes shifts in
14 the EEG across trials and recording sessions. In addition, inter-subject variability across
15 subjects, due to anatomical and physiopathological differences, hinders the design of a
16 ‘participant-agnostic’ BCI. Therefore, due to intra- and inter-subject variabilities, most BCIs
17 require long calibration times on each recording session to achieve satisfactory decoding
18 performance [31]. Furthermore, due to reduced P300 amplitude in ASD, the decoding of the
19 P300 response is even more challenging and, thus, the decoding performance may be
20 negatively affected.

21 Machine learning algorithms have been widely adopted to learn discriminative patterns from
22 the EEG to perform P300 decoding [32]. Among these algorithms, traditional decoders for
23 P300-based BCIs perform a pre-processing step that includes band-pass filtering within fixed
24 EEG bands (e.g., 0.5-4 Hz and 4-8 Hz [33]; 2-20 Hz and 2-8 Hz [34]; 2-12 Hz [35]), followed
25 by extraction of features in the temporal, frequency and spatial domains; the latter are then

1 evaluated by a classifier such as linear discriminant analysis, support vector machine or multi-
2 layer perceptron [33,35–40]. In addition to these traditional decoders, significant
3 improvements in performance were found using convolutional neural networks (CNNs) [41–
4 43]. CNNs are feed-forward neural networks including the convolutional operator at least in
5 one layer. Inspired by the hierarchical structure of the ventral stream of the visual system,
6 CNNs are composed by stacked layers of neurons, each neuron characterized by a local
7 receptive field. Neurons in deeper layers have larger receptive field and respond to more
8 complex features [44], enabling the learning of hierarchically structured features from the input
9 signal. At variance with traditional machine learning algorithms, in which a separation between
10 feature extraction, selection, and classification occurs, CNNs solve the decoding task in an end-
11 to-end fashion, by automatically learning the more meaningful features for the addressed
12 problem, i.e., discrimination of P300 events from single EEG trials. Therefore, CNNs do not
13 rely on some characteristics extracted *a priori* from the signals (e.g., spectral contents within
14 fixed bands), but automatically learn the relevant discriminative features to distinguish the
15 P300 response from the input EEG trial. From their first application in P300 decoding with the
16 simple design proposed by Cecotti et al. [45], CNNs were improved by including progressively
17 more convolutional and regularization layers (e.g., dropout [46] and batch normalization [47])
18 [48]. Among these CNNs, EEGNet [49] and its variants [43,50] were found to be particularly
19 suitable for P300 decoding, outperforming traditional machine learning solutions as well as
20 other CNN-based approaches, also in case of P300-based BCIs aimed at ASD intervention
21 [41,43]. In addition, by adopting specific training strategies (such as within-subject and cross-
22 session strategy or cross-subject strategy), CNNs were found to be capable of learning robust
23 features across sessions and subjects, encapsulating intra-subject and inter-subject variability,
24 thus providing the potentiality of significantly reducing BCI calibration times [43].

1 However, despite the previous advantages, these algorithms have some limitations. First,
2 they introduce a large number of trainable parameters, i.e., parameters to fit during the training
3 process, and require setting many hyper-parameters, i.e., parameters that define the functional
4 form of the decoder (e.g., convolutional filter size, number of convolutional filters, type of
5 activation function, etc.) that must be set before the training.

6 Generally, hyper-parameters are selected by testing only a few configurations via empirical
7 evaluations [43,48–50], and, thus, CNNs are built with a sub-optimal design in terms of
8 performance. Furthermore, hyper-parameters are kept the same across different training
9 strategies [43,50]. In this way, the network capacity (i.e., its ability of approximating a wide
10 variety of functions) is kept the same even though the network faces problems with increasing
11 difficulty across the different strategies, e.g., when moving from within-subject and within-
12 session to cross-session and cross-subject decoders, the architecture must learn the relevant
13 features from training distributions with increased variability. Searching for an optimal CNN
14 hyper-parameter configuration, which is also specific for a particular training strategy, is even
15 more necessary when designing a P300 decoder for a BCI-based ASD therapeutic approach,
16 where the P300 response is attenuated and more difficult to distinguish than in case of healthy
17 users.

18 Moreover, CNNs are scarcely interpretable in their learned features and are often treated as
19 black boxes. As pointed out in a recent survey [51], there is a growing interest to interpret the
20 feature representations provided by deep neural networks and their relation to clinical outcomes
21 quantifying neuropathology. Thus, research is moving from the sole decoding of brain states
22 towards the analysis of EEG features derived from the learning system, for example studying
23 EEG signatures related to movement [52–56], P300 [43,57,58], or depression [59]. Recently,
24 interpretable CNNs (ICNNs), i.e., CNNs that include layers whose learned parameters are
25 directly interpretable, were proposed to decode motor imagery and execution. Zhao et al. [53]

1 increased feature interpretability in the frequency domain by reparametrizing a convolutional
2 layer to learn Morlet wavelets. Recently, we proposed a lightweight (i.e., with a limited number
3 of trainable parameters) ICNN [52,55] able to increase the interpretability of both spectral and
4 spatial features, by reparametrizing a convolutional layer to learn band-pass filters and by
5 learning spatial filters tied to each band-pass filter. In addition to interpretable layers,
6 explanation techniques (ET) can be used to explain the CNN decision by highlighting which
7 EEG features learned in interpretable domains (e.g., spatial, temporal, spectral) are the most
8 discriminative (i.e., most salient) for the decoding of a specific cognitive state, such as the P300
9 [43,57,58]. Thus, by leveraging on the increased interpretability embedded into an ICNN,
10 combined with an ET (denoted as ICNN+ET in the following), a data-driven non-linear
11 analysis tool could be realized able to gain insights into the physiopathological neural
12 signatures contained in the EEG associated to P300.

13 In this study, we aim to contribute to the decoding and, at the same time, to the analysis of
14 the P300 in ASD using an ICNN, to further increase the feasibility of BCI intervention in
15 autism and to potentially characterize novel data-driven biomarkers related to ASD visuo-
16 spatial sensory processing. To this aim, here we adopted an architecture obtained by combining
17 two existing CNN architectures, i.e., Sinc-ShallowNet and EEGNet. Sinc-ShallowNet [52] is
18 an interpretable and lightweight CNN that we recently designed and proposed for decoding
19 purposes other than P300 (specifically, decoding of motor imagery and execution), and that
20 allows a straightforward interpretation of the learned spectral and spatial features. EEGNet
21 represents the state-of-the-art (SOA) CNN for P300 decoding [41,50]. The proposed
22 combination was conceived to obtain a CNN for P300 decoding (named Sinc-ShallowNet-v2
23 in the following), being both interpretable and with high decoding capabilities, as it inherits
24 the interpretable characteristic from Sinc-ShallowNet and, at the same time, it includes some

1 design aspects of EEGNet. Then, using this architecture, we address the following two issues
2 as main points of novelty of this study:

3 i. Investigate the optimal ICNN design (in terms of performance) and the role of the
4 main ICNN hyper-parameters under different training conditions, by performing
5 automatic hyper-parameter search (AHPS) based on Bayesian optimization (BO),
6 separately for each training condition. The investigated training strategies were
7 within-subject and within-session, within-subject and cross-session, and leave-one-
8 subject-out. Furthermore, we also tested the capability to transfer the knowledge
9 from other participants to a new one, adopting a transfer learning strategy to reduce
10 calibration time. Each of these training strategies may represent a different practical
11 scenario of BCI intervention to improve joint attention in ASD.

12 ii. Design of a novel algorithm based on the combination of the ICNN with saliency
13 representations (ICNN+ET) to highlight the most relevant spectral and spatial
14 features that correspond to the visuo-spatial P300 correlate in autism. These features
15 were then used to define clusters of subjects characterized by shared neural
16 signatures. Then, we investigated whether these features were related to clinical
17 scores measuring the severity of ASD symptoms as derived by developmental and
18 behavioral ASD assessment tools, and whether the ICNN+ET analysis better
19 enhanced useful ASD neural signatures compared to a canonical ERP analysis.

20 Both the transfer learning and the explanation technique were applied to models derived
21 from Bayesian optimization, as the latter, optimizing the decoding capabilities of the models,
22 contributes to increase the performance and reliability of the obtained results.

23

1 2. MATERIALS AND METHODS

2 2.1. Dataset description

3 In this study we used the BCIAUT-P300 dataset, an EEG dataset of a feasibility clinical trial
4 [28,29] publicly released for the IFMBE 2019 scientific challenge
5 (<https://www.kaggle.com/disbeat/bciaut-p300>). EEG signals were recorded from 15 high-
6 functioning ASD participants (22 years old, on average) while testing a P300-based BCI (based
7 on statistical classifiers) aimed to improve their joint attention. During a baseline visit, the
8 following clinical scores, useful also to diagnose ASD, were collected: Autism Diagnostic
9 Observation Schedule (ADOS) [60], Autism Diagnostic Interview-Revised (ADI-R) [61], and
10 Intelligence Quotients (IQs, measured by WAIS-III [62]). These scores are reported in Table
11 1.

12 For each participant, the BCI paradigm was carried out during 7 recording sessions (over 4
13 months) and in each session it was based on an immersive virtual environment presented to the
14 user, consisting of a bedroom with an avatar, common furniture (e.g., shelves, a bed, etc.) and
15 eight objects of interest. These objects were: 1) books on a shelf, 2) a radio on top of a dresser,
16 3) a printer on a shelf, 4) a laptop on a table, 5) a ball on the ground, 6) a corkboard on the
17 wall, 7) a wooden plane hanging from the ceiling, and 8) a picture on the wall (see Figure 1).
18 Participants were instructed to read non-verbal social agent cues from the avatar, i.e., avatar
19 head turning to a particular object, and to pay attention to that target object. The 8 objects
20 randomly flashed in the virtual scene and, thus, a visuo-spatial P300 response was elicited
21 when the attended object flashed.

22 Each recording session was composed by a calibration phase, consisting of $N_b = 20$ blocks,
23 and an online phase, consisting of $N_b = 50$ blocks. Each block was related to a particular object
24 selected as target by the avatar that the user tried to identify (see Figure 1 for a scheme about
25 the organization of blocks, runs, and trials). For each block, K runs were repeated ($K = 10$ in

1 the calibration phase, while K varied in the online phase with $K = 7.095$ on average across
2 online blocks); each run consisted of each of the 8 objects flashing once in a random sequence.
3 This resulted in 1600 trials (20 blocks x 10 runs x 8 EEG trials) per each participant and session
4 during the calibration phase, and in 2838 trials on average during the online phase. During the
5 calibration phase, the BCI statistical classifiers were trained to predict the object the participant
6 was paying attention to; during the online phase, the trained classifiers were applied to identify
7 whether the subject attended the target object correctly and, in that case, to provide positive
8 feedback to the user.

9 The IFMBE 2019 scientific challenge was organized to stimulate researchers to develop
10 decoders maximizing the object detection accuracy based on the P300 response. In the
11 challenge, for each subject and each session, trials recorded during the calibration phase were
12 released to *tune* decoders (i.e., to set their parameters), while trials recorded during the online
13 phase were used to *test* decoders. In the present study, we adopted this same split as defined
14 by the challenge, as this choice may also allow a fair and direct comparison with the results of
15 decoders that participated to the challenge [41]. More specifically, for each subject and each
16 session, we further split the 1600 calibration trials into two sets: 80% of trials (1280 trials),
17 were randomly sampled and used as *training set* to optimize the trainable parameters of the
18 decoders, while the remaining 20% of trials (320 trials) were used as *validation set* to optimize
19 the hyper-parameters of the decoders. This further splitting was based on our winning solution
20 of the challenge [41,50] to select the number of training epochs (i.e., to perform early stopping).
21 Therefore, each subject-specific and session-specific dataset was separated into 3 different sets,
22 as commonly performed in the literature [41,43,45,49,50,52,55,57,63]: training set (1280
23 trials), validation set (320 trials), and test set (2838 trials on average). However, since different
24 training strategies were adopted here, by differently aggregating the training set and validation
25 set across sessions and subjects, the overall number of training and validation trials were

1 different depending on the strategy (see Section 2.5.1). Finally, test trials were used to test the
 2 tuned decoders on a held-out set.

3 EEG signals were recorded at 250 Hz from C3, Cz, C4, CPz, P3, Pz, P4, and POz locations
 4 ($C = 8$ electrode sites), with the reference placed at the right ear and the ground at AFz. These
 5 signals were acquired notch filtered at 50 Hz and filtered between 2 and 30 Hz [61]. In this
 6 study, as in the winning solution [50] that we proposed for the challenge, each trial contained
 7 signals from -0.1 s to 1 s relative to the stimulus onset and signals were downsampled to 128
 8 Hz, so that each trial contained $T = 140$ time steps.

9

10 2.2 Problem definition

11 Based on the previous description, the collection of trials associated to the s -th subject
 12 acquired during the r -th recording session can be formalized as the collection $D^{(s,r)} =$
 13 $\{(X_0^{(s,r)}, y_0^{(s,r)}), \dots, (X_i^{(s,r)}, y_i^{(s,r)}), \dots, (X_{M^{(s,r)}-1}^{(s,r)}, y_{M^{(s,r)}-1}^{(s,r)})\}$. $M^{(s,r)}$ denotes the total number of
 14 trials for that subject and that session, $X_i^{(s,r)} \in \mathbb{R}^{C \times T}$ represents the pre-processed EEG signals
 15 of the i -th trial ($0 \leq i \leq M^{(s,r)} - 1$), and $y_i^{(s,r)}$ represents the label of the i -th trial, i.e., $y_i^{(s,r)} \in$
 16 $L = \{l_0, l_1\} = \{\text{non-P300}, \text{P300}\}$, where the label P300 was assigned to those trials where the
 17 flashing object coincided with the object the avatar was looking at.

18 In this study, each decoder consisted in a parametrized classifier f (representing the ICNN
 19 and having a different functional form depending on the hyper-parameters), which solves a
 20 binary classification task: $f(X_i^{(s,r)}; \theta): \mathbb{R}^{C \times T} \rightarrow L$, where θ represents the array of trainable
 21 parameters. Thus, the ICNN input was represented by $X_i^{(s,r)}$ that can be viewed as a 2D matrix
 22 of shape $(C, T) = (8, 140)$ with electrodes along the height and time steps along the width, and
 23 the output consisted of two neurons corresponding to either one or the other class. The
 24 validation set was used to perform the automatic search of ICNN hyper-parameters via

1 Bayesian optimization, and the ICNN learned automatically from the training set the relevant
2 features to assign the correct label to unseen input trials belonging to the test set. The
3 knowledge learned by the ICNN during the training process was stored in its trainable
4 parameters θ . These parameters, thanks to their increase interpretability, can be exploited to
5 gain insights into the neural signatures related to a specific class (e.g., P300 class) in a data-
6 driven way, without relying on handcrafted features based on expected EEG responses.

7

8 **2.3. An update of Sinc-ShallowNet: Sinc-ShallowNet-v2**

9 Sinc-ShallowNet [52] is an ICNN that we developed to decode motor imagery and execution
10 from single EEG trials. This ICNN is composed of two blocks, each consisting of several
11 stacked layers: an interpretable spectral and spatial (ISS) feature extractor followed by a fully-
12 connected (FC) block that performs classification. The ISS block was designed to increase the
13 interpretability of the learned parameters, at the same time keeping limited the model size, i.e.,
14 the number of trainable parameters, and included a temporal convolutional layer (with a
15 reparameterization of the kernels), a depthwise spatial convolution and an averaging pooling
16 layer (see also below for a more detailed description of the ISS block). Crucially, in Sinc-
17 ShallowNet the output of the ISS block was provided directly as input to the FC block. Here,
18 we proposed an updated version for P300 decoding, named Sinc-ShallowNet-v2, by integrating
19 in the design also structure elements of EEGNet, a CNN proved to be particularly suitable to
20 decode the P300 response from EEG [41,49,50] but not designed to be interpretable.
21 Specifically, the updated Sinc-ShallowNet-v2 embraces three main blocks by including an
22 additional block, the fixed-scale temporal (FST) feature extractor (inspired by EEGNet),
23 between the ISS block and the FC classification block. This is an important modification, since
24 Sinc-ShallowNet-v2 further processes the output of the ISS block by learning features in the
25 temporal domain, and this may help to better capture intra- and inter-subject variability of P300

1 in time. A schematization of Sinc-ShallowNet-v2 is reported in Figure 2a; a more detailed
 2 description reporting the hyper-parameters, output shape and number of trainable parameters
 3 per layer, is reported in Table 2. In the following, the blocks are described in detail.

4 *i. Block 1: Interpretable spectral and spatial (ISS) feature extractor*

5 The first block was inspired by the ISS block of Sinc-ShallowNet [52] and was devoted to
 6 separately learn spectral and spatial features from the input EEG trials in an easy interpretable
 7 way. The very first layer of the ISS block was a temporal sinc-convolutional layer [52,55,64],
 8 learning K_0^{ISS} filters with filter size $F_0^{ISS} = (1,65)$, unitary stride and zero-padding $P_0^{ISS} =$
 9 $(0,32)$ to preserve the number of input temporal samples. This temporal convolutional layer
 10 was devoted to filter each electrode signal in time. By using the sinc-convolutional layer to
 11 perform such processing step instead of a conventional convolutional layer, each convolutional
 12 filter is forced to describe a band-pass filter in the temporal domain. Denoting with k_j the j -th
 13 convolutional kernel, in a conventional convolutional layer each value of the filter (i.e.,
 14 $k_j[0, n], n \in [0,64]$) has to be learned during the optimization process; conversely, in a sinc-
 15 convolutional layer, each value of the filter is defined by a parametrized function, forcing the
 16 filters to belong to a specific subset of temporal filters (here band-pass filters). Therefore, in a
 17 sinc-convolutional layer a re-parametrization of each convolutional kernel occurs:

$$18 \quad k_j'[0, n; \{f_{0,j}, f_{1,j}\}] = 2f_{1,j} \text{sinc}(2\pi f_{1,j}n) - 2f_{0,j} \text{sinc}(2\pi f_{0,j}n). \quad (1)$$

19 In Equation 1, $\{f_{0,j}, f_{1,j}\}$ are the trainable parameters related to the j -th kernel, including
 20 only the inferior ($f_{0,j}$) and superior ($f_{1,j}$) cutoff frequencies of the band-pass filter. In this way,
 21 the number of trainable parameters reduces from 65 ($= F_0^{ISS}[0] \cdot F_0^{ISS}[1]$) to 2, for each
 22 temporal filter. Lastly, to alleviate the effects of the inevitable truncation of k_j' on the
 23 characteristics of each filter, k_j' is multiplied by a Humming window:

$$24 \quad \begin{cases} k_{w,j}'[0, n; \{f_{0,j}, f_{1,j}\}] = k_j'[0, n; \{f_{0,j}, f_{1,j}\}] \cdot w[n] \\ w[n] = 0.54 - 0.46 \cos\left(\frac{2\pi n}{F_0^{ISS}[1]-1}\right) \end{cases}. \quad (2)$$

1 Accordingly, the temporal sinc-convolution computes the convolution between the input
 2 and $k_{w,j}'[0, n; \{f_{0,j}, f_{1,j}\}]$, learning only the following 2 parameters for each kernel:

$$3 \quad \theta_{spect,j} = \{f_{0,j}, f_{1,j}\} \in \theta, 0 \leq j \leq K_0^{ISS} - 1. \quad (3)$$

4 Thus, the output of this first layer consists of stacked feature maps containing band-pass
 5 filtered versions of the input EEG trial within specific frequency ranges that were explicitly
 6 learned during training.

7 The use of a temporal sinc-convolutional layer, besides substantially reducing the
 8 parameters to fit, promotes the learning of more meaningful and well-defined temporal filters,
 9 and provides a straightforward interpretation of the learned spectral features [64], being the
 10 cutoff frequencies of the learned band-pass filters. Conversely, in case of a conventional
 11 temporal convolutional layer, the learned spectral features are not immediately accessible and
 12 interpretable.

13 Downstream the temporal sinc-convolutional layer, a batch normalization layer [47] was
 14 included. Then, a spatial depthwise convolutional layer was introduced: for each band-pass
 15 filtered map, D_1^{ISS} spatial filters were learned having size $(C, 1)$, unitary stride and no zero-
 16 padding, i.e., D_1^{ISS} spatial combinations of electrodes were learned for each band-pass filtered
 17 map. Therefore, a total number of $K_1^{ISS} = K_0^{ISS} \cdot D_1^{ISS}$ spatial filters were learned and
 18 constrained to have a norm upper bounded by c_1^{ISS} (kernel max-norm constraint). This type of
 19 convolution does not exploit dense connections across feature maps as in traditional
 20 convolutional layers, reducing the number of trainable parameters. In addition, the combination
 21 of temporal sinc-convolution with spatial depthwise convolution leads to an interpretable
 22 spectral-spatial feature learning, as each group of D_1^{ISS} spatial filters is strictly tied (via
 23 depthwise convolution) to a specific band-pass filter, i.e., to a specific frequency range:

$$24 \quad \theta_{spat,j} = \{\theta_{j0}, \dots, \theta_{jk}, \dots, \theta_{jD_1^{ISS}-1}\} \in \theta, 0 \leq j \leq K_0^{ISS} - 1, \quad (4)$$

25 indicating with θ_{jk} the k-th spatial filter ($0 \leq k \leq D_1^{ISS} - 1$) tied to the j-th band-pass filter.

1 The whole parameters learned in these two first convolutional layers can be denoted as the
 2 set:

$$\begin{aligned}
 3 \quad \theta_{ISS} &= \left\{ \theta_{ISS,0}, \dots, \theta_{ISS,j}, \dots, \theta_{ISS,K_0^{ISS}-1} \right\} = \\
 4 \quad &\left\{ (\theta_{spect,0}, \theta_{spat,0}), \dots, (\theta_{spect,j}, \theta_{spat,j}), \dots, (\theta_{spect,K_0^{ISS}-1}, \theta_{spat,K_0^{ISS}-1}) \right\}, \quad (5)
 \end{aligned}$$

5 where each pair $(\theta_{spect,j}, \theta_{spat,j})$, $0 \leq j \leq K_0^{ISS} - 1$ contains the cutoff frequencies of the j -
 6 th band-pass filter and the associated D_1^{ISS} spatial filters exploited to decode the input EEG
 7 trial.

8 Output activations of the interpretable spectral-spatial feature extractor were then
 9 normalized via batch normalization [47] and activated via an Exponential Linear Unit (ELU)
 10 non-linearity [65], i.e., $f(x) = x, x > 0$ and $f(x) = \exp(x) - 1, x \leq 0$. Average pooling was
 11 introduced to reduce the temporal dimension of the activations by a factor of 4 ($F_p^{ISS} = S_p^{ISS} =$
 12 $(1, 4)$) from $T = 140$ to $T//4 = 35$ (indicating with $//$ the floor division operator). Lastly, a
 13 dropout layer [46] was added with dropout rate p^{ISS} .

14 *ii. Block 2: Fixed-scale temporal (FST) feature extractor*

15 This block was inspired by EEGNet, and in Sinc-ShallowNet-v2 was used to learn how to
 16 summarize the feature maps provided by the ISS block in the time domain; it contained a
 17 separable convolutional layer [66], defined by a temporal depthwise convolution with $D_0^{FST} =$
 18 1 , filter size F_0^{FST} , unitary stride and zero padding $P_0^{FST} = (0, F_0^{FST}[1]//2)$, followed by a
 19 pointwise convolution with K_0^{FST} filters. In the first layer of this composition, depending on
 20 F_0^{FST} , temporal features were learned on the input ISS feature maps within a temporal window
 21 of specific size, i.e., at a fixed temporal scale, and without using dense connections across
 22 feature maps. The temporal window in which features are learned corresponds to
 23 $F_0^{FST}[1]/(sf//4)$ s, indicating with sf the sampling frequency of the input EEG signals. In
 24 the second layer of the separable convolutional, feature maps at the output of the previous layer

1 are recombined, and, this is performed learning compressed, equal or overcomplete
2 representations, depending on whether $K_0^{FST} < K_1^{ISS}$, $K_0^{FST} = K_1^{ISS}$ or $K_0^{FST} > K_1^{ISS}$,
3 respectively. Then, activations were normalized via batch normalization [47] and activated via
4 an ELU non-linearity [65]. An average pooling layer was introduced to further reduce the
5 temporal dimension of the activations by a factor of 8 ($F_p^{FST} = S_p^{FST} = (1, 8)$), i.e., from
6 $T//4 = 35$ to $T//32 = 4$. Lastly, a dropout layer [46] was added with dropout rate p^{FST} .

7 *iii. Block 3: Fully-connected (FC) block*

8 This block traduces the activations provided by the FST block into conditional probabilities,
9 finalizing the decoding task. At first, the input feature maps of the FC block were unrolled
10 along one single dimension using a flatten layer, producing a single feature array. Then, this
11 array was provided as input to a single fully-connected layer with $N^{FC} = 2$ neurons associated
12 to the non-P300 and P300 classes. Here, trainable parameters were constrained to have a norm
13 upper bounded by c_0^{FC} . Lastly, these 2 neurons were activated using a softmax activation
14 function to obtain the output conditional probabilities $p(l_k | X_i^{(s,r)})$, $k = 0, 1$.

15

16 **2.4. Performance metric and comparison of Sinc-ShallowNet-v2 with EEGNet and Sinc-** 17 **ShallowNet**

18 In this study, we adopted the object-level accuracy as performance metric, i.e., accuracy in
19 decoding the flashing object the participant was paying attention to (among the 8 possible
20 objects, see Section 2.1 for additional details). This performance metric was the same as
21 adopted in the IFBME 2019 competition to evaluate decoders [41]. To this aim, the trial-level
22 EEG decoding provided by the CNN (binary classification, i.e., non-P300 vs. P300) was used
23 to produce an object-level decoding (8-class classification, i.e., discrimination of the attended
24 object among 8 objects). Considering a specific recording block, associated to a specific object
25 the participant was paying attention to, the following processing was performed. Indicating

1 with l_1 the P300 condition, we considered the probabilities $p(l_1 | X_i^{(s,r)})$ predicted for the EEG
2 trials $X_i^{(s,r)}$ within that block when each of the 8 objects flashed (including the attended one).
3 These probabilities were averaged across runs separately for each object, obtaining the
4 average probability $\bar{p}_o, 0 \leq o \leq 7$ that the participant was paying attention to the o -th object
5 Then, the object with the highest probability was predicted as the one attended in that block.

6 We compared the performance of Sinc-ShallowNet-v2 against our approach derived from
7 EEGNet [50] that won the IFMBE 2019 challenge. In the competition, that approach
8 significantly outperformed traditional machine learning solutions as well as other deep neural
9 networks, including both CNNs [48] and recurrent neural networks [41], thus, it can be
10 considered the current SOA algorithm for the addressed decoding problem. Furthermore, Sinc-
11 ShallowNet-v2 performance was compared against Sinc-ShallowNet [52], to assess whether
12 the inclusion of the FST block in the updated version led to potential benefit in P300 decoding
13 compared to the previous architecture (lacking this block).

14

15 **2.5. Trainable parameter optimization**

16 **2.5.1. Training strategies**

17 The structure of the adopted dataset BCIAUT-P300 (including several sessions per subject
18 and a large number of trials per session), allows to train and evaluate P300 decoders in a large
19 spectrum of training conditions, each reflecting a different scenario in which a P300-based BCI
20 intervention might be applied. In this study, we assessed the optimal design of the proposed
21 ICNN for P300 decoding, separately for each of different training strategies, characterized by
22 increasing variability in the dataset, i.e., simulating BCI paradigms applied progressively to
23 more sessions and subjects.

24 As introduced in Section 2.1, the trials of each session and subject were divided into training
25 validation sets (80% and 20% of trials of the calibration phase, randomly sampled) and test set

1 (trials of the online phase). Then, depending on whether training and validation sets were
2 considered separately for each subject or session, or were differently aggregated across
3 sessions and subjects, four different training strategies were designed and investigated. It is
4 worth mentioning that, despite different training strategies were employed, the definition of the
5 test set was the same across training strategies to perform a fair comparison across them, that
6 is, the test set always included trials belonging to the 50 online blocks. These strategies are
7 described in the following. Furthermore, Figure 2b provides a schematization of the adopted
8 strategies, and Table 3 a summary of the number of trials belonging to the training, validation
9 and test sets in each strategy.

10 *i. Within-subject and within-session (WS-WS) strategy.*

11 In this strategy, subject- and session-specific training and validation sets were used to tune
12 subject- and session-specific CNNs. The test set was defined as the subject- and session-
13 specific test set, i.e., trials of the online phase belonging to the subject and session the CNN
14 was trained for (see Table 3A). Furthermore, to simulate a practical scenario where limited
15 trials are available in a single recording session, we trained and validated CNNs using a
16 progressively increasing number of calibration blocks (see Section 2.1), and thus, using
17 variable-sized training and validation sets. This condition was implemented by sampling
18 training and validation examples from the first n calibration blocks, where n was progressively
19 increased from 2 to 20 with a step of 2 blocks (see Table 3B), where 20 is the overall number
20 of calibration blocks in each session.

21 *ii. Within-subject and cross-session (WS-CS) strategy.*

22 In this strategy, subject-specific training and validation sets were used to tune subject-
23 specific CNNs. For the s -th subject, training and validation sets were merged across all
24 recording sessions of that specific subject (see Table 3A). incorporating within-subject
25 variability. The architecture was then tested separately on each subject- and session-specific

1 test set, i.e., on trials of the online phase of each session belonging to the subject the CNN was
2 trained for.

3 *iii. Leave-one-subject-out (LOSO) strategy.*

4 In this strategy, cross-subject training and validation sets were used to tune cross-subject,
5 cross-session and subject-agnostic CNNs, i.e., the training and validation sets included only
6 examples sampled from other subject distributions. In particular, for each subject s , named
7 “held-back subject”, training and validation sets were merged across all sessions from all the
8 other subjects (different from the s -th subject, see Table 3A), incorporating between-subject
9 and within-subject variability. The architecture was then tested separately on each subject- and
10 session-specific test set, i.e., on trials of the online phase of each session belonging to the held-
11 back subject (s -th subject).

12 *iv. Transfer learning on single sessions (TL-WS).*

13 Transfer learning focuses on transferring the knowledge across domains/tasks. It is inspired
14 by the human capability to use the knowledge learned in a source domain/task to improve the
15 performance and/or reduce the training time in a related domain/task [67]. In this strategy, as
16 for the WS-WS strategy (2.5.1-i), training and validation trials were subject- and session-
17 specific. In particular, the definition of training and validation sets was the same as the one
18 adopted in the WS-WS strategy when analyzing variable-sized training and validation sets
19 (depending on the number of included calibration blocks), while keeping unchanged the
20 subject- and session-specific test set, i.e., the trials of the online phase of that subject and
21 session (see Table 3B). However, differently from the WS-WS strategy where the trainable
22 parameters were initialized randomly, in the TL-WS strategy, for the s -th subject and r -th
23 recording session the CNN was initialized using the CNN trained with LOSO strategy when
24 the s -th subject was held-back. That is, the knowledge learned during the LOSO strategy from
25 many subjects except the held-back one was transferred on the held back subject. The use of

1 TL-WS strategy could be useful when a new user approaches the BCI system in a new
2 recording session and a calibration stage, as short as possible, is needed to tune an accurate
3 decoder. In this view, the knowledge embedded in LOSO models, i.e., incorporating between-
4 subject and within-subject variabilities, could represent a better initialization point in the space
5 of parameter θ than the random one, potentially leading to an improvement in performance
6 and/or to a reduction of training and validation examples needed to achieve high performance.
7 The potentialities of transfer learning were tested by comparing the performance of the TL-WS
8 models with the performance of the WS-WS models, while increasing the size of the training
9 and validation sets, thus, assessing the benefits of the use of models pre-trained on other
10 subjects compared to models trained from scratch.

11

12 **2.5.2. Training settings**

13 CNN optimization consisted of the minimization of the cross-entropy between the empirical
14 probability distribution defined by the training labels, and the probability distribution defined
15 by the model. This corresponds to minimize the Kullback-Leibler divergence between the two
16 probability distributions at trial-level, and thus also at block-level, i.e., object-level. Adaptive
17 moment estimation (Adam) [68] was used as optimizer with learning rate lr , mini-batch size
18 $bs = 64$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$ for computing the running averages of the gradient and its
19 square, and $\varepsilon = 10^{-8}$ to improve numerical stability. To address class imbalance, parameter
20 updates were weighted depending on the class occurrence of the training examples: indicating
21 with M_0 and M_1 , the number of trials for non-P300 and P300 conditions in the training set and
22 given that $M_0 > M_1$, class weights were defined as 1 and $M_0/M_1 = 7$, respectively for non-
23 P300 and P300 classes. The maximum number of epochs was set to 1000 and early stopping
24 was performed by interrupting the training loop when the validation loss did not decrease for
25 50 consecutive epochs. Besides early stopping, which acts as regularizer, Sinc-ShallowNet-v2

1 implicitly included in its structure also many layers devoted to increase the generalization, such
2 as batch normalization [47] (with a momentum term $m = 0.99$ and $\varepsilon = 1e-3$ for numerical
3 stability), dropout [46] and kernel max norm constraint.

4

5 **2.6. Hyper-parameter optimization**

6 Hyper-parameter optimization is devoted to find the optimal hyper-parameter configuration
7 of a learning system associated with the best performance measured on a separate validation
8 set. In the following, an overview of hyper-parameter tuning via Bayesian optimization (BO)
9 [69] is provided; then, the procedure adopted to perform BO hyper-parameter search and to
10 report the obtained results is described.

11 **2.6.1. Hyper-parameter search via Bayesian optimization**

12 Denoting with h the array containing the hyper-parameters of interest ($h \in H$, where H is
13 the hyper-parameter search space), the aim of hyper-parameter optimization is to find $h^* =$
14 $\mathit{arg} \min_{h \in H} k(h)$, that minimizes the objective function $k(h)$ (the cross-entropy loss, in this study)
15 evaluated on the validation set. However, the evaluation of the objective function $k(h)$ requires
16 a new training stage and a new evaluation stage (on the validation set) for each hyper-parameter
17 configuration h . Thus, depending on the number of hyper-parameters to optimize and on the
18 complexity of the model, which are generally both high in deep learning-based approaches,
19 this process can be expensive. Common hyper-parameter search algorithms, such as grid search
20 and random search, do not take advantage of the results of the previous iterations to select the
21 next hyper-parameters to evaluate (the latter are, thus, uninformed by past evaluations), often
22 wasting computational time on hyper-parameters with poor performance. Conversely, BO
23 methods keep track of results related to past evaluations to update a probabilistic model,
24 mapping hyper-parameters to the probability to obtain a specific score from the objective
25 function. That is, BO can be formalized as a sequential model-based optimization, performing

1 several iterations each one considering a specific hyper-parameter configuration and updating
2 the probability model. This model is a surrogate of the objective function to be minimized and
3 is easier to optimize than the original objective function. BO methods suggest the hyper-
4 parameters in an informed way after each iteration step, based on a “selection function”. By
5 investigating hyper-parameters that seem promising based on past results, BO methods can
6 find better configurations than other approaches within fewer iterations compared to
7 uninformed algorithms (e.g., grid or random search) [69].

8 **2.6.2. Hyper-parameter search settings and analysis**

9 Optimal hyper-parameters were searched via BO (using the Python library Optuna [70],
10 version 2.3.0), using tree-structured Parzen estimator as surrogate function and expected
11 improvement as selection function, and sequential model-based optimization was performed
12 for 100 iterations (which is the default value [70]). The hyper-parameters of Sinc-ShallowNet-
13 v2 subjected to BO, defining the array $h \in H$, were K_0^{ISS} , D_1^{ISS} , c_1^{ISS} , $p^{ISS} = p^{FST} = p$,
14 K_0^{FST} , F_0^{FST} , c_0^{FC} and lr . Within each BO iteration, the hyper-parameters were sampled from
15 the distributions specified in Table 4. Note that, based on the adopted distributions, since K_0^{FST}
16 can be equal to or less than (but not greater than) $K_0^{ISS} \cdot D_1^{ISS}$, only compressed or equal (but
17 not overcomplete) representations in FST layer were examined. BO was applied to models
18 trained with WS-WS, WS-CS and LOSO strategies to investigate the optimal design of Sinc-
19 ShallowNet-v2 depending on the training strategy (see Figure 2b). BO was not applied to TL-
20 WS, as models trained in this strategy strictly depended on LOSO models. Indeed, in TL-WS
21 the ICNN was initialized with LOSO parameters and then trained on the held-back subject.
22 Therefore, the hyper-parameter configuration adopted in TL-WS models needed to be the same
23 as in the LOSO strategy, i.e., hyper-parameters needed to be kept fixed as the ones of Bayesian-
24 optimized LOSO models. Therefore, also the WS-WS models used for comparison with TL-
25 WS models, while training and evaluating decoders with a progressively increasing size of

1 training and validation sets, were assigned to the Bayesian-optimized hyper-parameters of the
2 LOSO models, for a fair comparison. It is worth remarking that while analyzing the
3 potentialities of transfer learning with TL-WS and WS-WS models, where hyper-parameters
4 were inherited from LOSO models, the validation set was used only to perform early stopping.

5 For each training condition (WS-WS, WS-CS and LOSO), and for a given q -th hyper-
6 parameter $\in h$ ($0 \leq q \leq 7$) the optimal values obtained across the Bayesian-optimized ICNNs
7 were extracted (we had 15·7 ICNNs in WS-WS and 15 ICNNs in WS-CS and LOSO). Then,
8 the probability (P) to obtain a specific hyper-parameter value via BO across ICNNs was
9 computed. This allowed a visualization and a comparison of the optimal model design across
10 training conditions. Moreover, for each training condition the importance score z_q was derived
11 by using the fANOVA hyper-parameter importance evaluation algorithm proposed in [71] that
12 fits a random forest regression model predicting the objective value given a parameter
13 configuration. Importance scores sum up to 1 over the investigated hyper-parameters (i.e.,
14 $\sum_q z_q = 1$) and the higher the score of a specific hyper-parameter, the higher its importance.

15 To provide a fair performance comparison with the other two CNN-based approaches, the
16 EEGNet adaptation used in [50] and Sinc-ShallowNet [52] (see Section 2.4), the optimal hyper-
17 parameters of EEGNet and Sinc-ShallowNet were searched using BO for each training
18 condition too. In particular, the hyper-parameters of EEGNet and Sinc-ShallowNet were
19 sampled using the same distributions defined as for Sinc-ShallowNet-v2 (see Table 4); of
20 course, Sinc-ShallowNet lacked the hyper-parameters of the FST block (K_0^{FST} and F_0^{FST}) as
21 this block was absent in this architecture.

22

23 **2.7. Explanation technique: spectral and spatial features analysis**

24 The interpretation of the features learned by Sinc-ShallowNet-v2 incorporating subject-
25 specific knowledge ($\theta = \theta^{(s)}$, in WS-CS) could provide insights about neural signatures

1 related to each specific ASD subject in a data-driven way, enabling between-subject variability
2 investigation in ASD. For this analysis, we retrained the ICNN in WS-CS using the most
3 frequent Bayesian-optimized configuration obtained in WS-WS across all subjects and
4 sessions. By doing so, only one fixed hyper-parameter configuration was used to retrain WS-
5 CS models across all subjects, so that the same number of temporal and spatial filters in the
6 ISS block – whose learned features are the objects of the following analysis – was used across
7 subjects. In this way, we were able to evidence differences/similarities among ASD subjects in
8 terms of the learned features, excluding that these differences may arise from differences in
9 ICNN configurations. The fixed configuration was based on the most frequent WS-WS optimal
10 configuration, as the Bayesian-optimized WS-WS models resulted the lightest and fastest to
11 train (see Figure 4 and Table 5). It is worth noticing that the WS-CS retrained models achieved
12 comparable performance with the Bayesian-optimized WS-CS models (see Section 1 of
13 Supplementary Materials).

14 As described in Section 2.3-i, the adopted architecture was designed to provide interpretable
15 parameters in the array $\theta_{ISS}^{(s)}$. The ICNN processes input trials by filtering out P300-unrelated
16 spectral and spatial information while preserving only the most significant ones for P300
17 decoding. However, features may have a different importance on the discrimination, i.e., a
18 band-pass filtering in a peculiar frequency range and a subset of electrodes may be more
19 relevant to distinguish the P300 response. Therefore, the processing of θ_{ISS} should also include
20 an explanation technique, devoted to highlight contributions of the more important features in
21 ASD related to P300, realizing the combination ICNN+ET.

22 In the following, the proposed algorithm based on ICNN+ET for the investigation of
23 subject-specific spectral and spatial P300 features is described and a schematic representation
24 is reported in Figure 3.

25 **2.7.1. Computation of the relevance of spectral features (spectral relevance)**

1 At first, given a single input EEG trial of the s -th subject containing the P300 response
2 $(X_i^{(s,r)})$, we computed the saliency [72] for the P300 class of each spatio-temporal sample ($C \cdot$
3 T samples) within each feature map provided by the temporal sinc-convolutional layer. This
4 ET consists in computing the gradient of the output of the neuron associated to the P300
5 condition (immediately before the softmax function) with respect to the feature maps provided
6 by the first convolutional layer. The output of the ET is one relevance map for each feature
7 map of the first layer (i.e., $\forall j$, with $0 \leq j \leq K_0^{ISS} - 1$), and can be interpreted as a
8 transformation $e(X_i^{(s,r)}): \mathbb{R}^{C \times T} \rightarrow \mathbb{R}^{K_0^{ISS} \times C \times T}$ whose output quantifies how much the spatio-
9 temporal samples in each filtered version of the input affect the P300 prediction. Relevance
10 maps were averaged across P300 trials, and the absolute value of these maps was computed.
11 Therefore, an average relevance map $\bar{G}_j^{(s)} \in \mathbb{R}^{C \times T}$, $0 \leq j \leq K_0^{ISS} - 1$ was obtained, and
12 finally, the relevance score of each feature map was computed as:

$$13 \quad g_j^{(s)} = \max_{c,t} \bar{G}_j^{(s)} / \max_{c,t,j} \bar{G}_j^{(s)}, \quad 0 \leq j \leq K_0^{ISS} - 1, \quad 0 \leq c \leq C - 1, \quad 0 \leq t \leq T - 1, \quad (6)$$

14 where $g_j^{(s)} \in [0,1]$ is a scalar quantity that summarizes the importance of each band-pass filter
15 for discriminating the P300 response from input EEG trials.

16 The spectral relevance scores obtained in Equation 6 were used to weight the associated
17 passbands of the temporal filters, defined by $\theta_{spect,j}^{(s)}$ (see Equation 3), computing $r_j^{(s)}(f)$ as
18 follows:

$$19 \quad \begin{cases} p_j^{(s)}(f) = \begin{cases} 1, & \text{if } f_{0,j}^{(s)} \leq f \leq f_{1,j}^{(s)} \\ 0, & \text{elsewhere} \end{cases} & 0 \leq j \leq K_0^{ISS} - 1 \\ r_j^{(s)}(f) = g_j^{(s)} \cdot p_j^{(s)}(f) & 0 \leq j \leq K_0^{ISS} - 1 \end{cases}. \quad (7)$$

20 In Equation 7, $p_j^{(s)}(f)$ indicates the probability that a specific frequency f was included in the
21 passband of the j -th band-pass filter, i.e., marking with a probability of 0 and 1 frequencies
22 outside and inside the passband, respectively.

1 Then, the spectral relevance $r^{(s)}(f)$, quantifying the relevance of each frequency bin, was
 2 obtained as:

$$3 \quad r^{(s)}(f) = \max_j r_j^{(s)}(f). \quad (8)$$

4 **2.7.2. Computation of the relevance of spatial features related to the most relevant** 5 **spectral features (spatial relevance)**

6 By averaging $r^{(s)}(f)$ across subjects, the frequency ranges retaining a relevance greater
 7 than or equal to 0.75 across subjects ($\forall s$) were identified. As described in Section 3.3, only
 8 one frequency range was identified, and it is indicated in the following with $[f_0^*, f_1^*]$. This
 9 interval is characterized by a high relevance with high agreement ($\geq 75\%$) across subjects, as
 10 the spectral relevance was normalized for each subject. For each subject the following analysis
 11 was performed. First, we considered the subset of the band-pass filters, denoted as $S^{(s)}$, that
 12 contained in their passband the frequency bins belonging to $[f_0^*, f_1^*]$ and we selected the spatial
 13 filters (i.e., the learned features $\theta_{jk}^{(s)}$, see Equation 4) associated to this subset of band-pass
 14 filters ($j \in S^{(s)}, 0 \leq k \leq D_1^{ISS} - 1$). As we were interested in the investigation of the relevance
 15 at the level of single electrode, spatial filters were considered in their absolute value, as done
 16 in [45,52]. Subsequently, the absolute spatial features were averaged together electrode per
 17 electrode ($\forall c$) and normalized to the maximum across electrodes, obtaining the spatial
 18 relevance:

$$19 \quad \bar{\theta}_{spat}^{(s)} = \frac{1}{N_{spat}^{(s)}} \sum_{j \in S^{(s)}, k} abs(\theta_{jk}^{(s)}) / \max_c \frac{1}{N_{spat}^{(s)}} \sum_{j \in S^{(s)}, k} abs(\theta_{jk}^{(s)}), \quad 0 \leq c \leq C - 1, \quad (9)$$

20 where $N_{spat}^{(s)}$ represents the overall number of spatial filters associated to the subset $S^{(s)}$ of the
 21 band-pass filters (D_1^{ISS} spatial filters for each band-pass filter in the subset).

22 Based on Equation 9, each value of $\bar{\theta}_{spat}^{(s)} \in \mathbb{R}^C$ quantifies the relevance (normalized
 23 between $[0,1]$) of a specific electrode signal filtered in the frequency range most important for

1 decoding the P300 response. Abnormalities in P300 (e.g., ASD abnormalities) may modulate
2 $\bar{\theta}_{spat}^{(s)}$.

3 **2.7.3. Clustering**

4 The spatial relevance $\bar{\theta}_{spat}^{(s)}$ was clustered, automatically finding clusters of subjects
5 characterized by different patterns of $\bar{\theta}_{spat}^{(s)}$. This was made possible by adopting the WS-CS
6 strategy, as it allows the learning of robust subject-specific neural signatures across recording
7 sessions. Here, Hierarchical Density-Based Spatial Clustering of Applications with Noise
8 (HDBSCAN) [73] was used as clustering algorithm (using the Python library hdbscan [74],
9 version 0.8.27), using the Euclidean distance between observations as distance metric. With
10 HDBSCAN, clustering is performed without specifying the number of clusters as input
11 parameter, including also an outlier detection algorithm [73]. Therefore, the adopted clustering
12 algorithm provides the optimal number of clusters N_{clust} populated by the more reliable
13 observations and marks as outliers a subset O of observations. As observations belonging to
14 the same cluster share a common structure, these were averaged together, obtaining $\bar{\theta}_{c,spat}$ for
15 the c -th cluster. Therefore, cluster-level representations were obtained by averaging across a
16 small subset of subjects sharing a peculiar spatial pattern as learned by the model, reflecting
17 neural signatures belonging to different clusters of ASD subjects.

18 **2.8. Statistical analysis**

19 **2.8.1. Performance**

20 As described in Section 2.4, object-level accuracy was used as performance metric, and it
21 was computed separately for each session-specific test set. In addition, as performed in [41],
22 as the performance resulted robust across recording sessions, the performance metric was
23 averaged across all sessions for each subject. Then, the following tests were performed.

- 24 i. A Friedmann test was performed to compare the performance obtained by Sinc-
25 ShallowNet-v2 trained with the three training strategies (WS-WS, WS-CS, LOSO

1 strategies). Then, as significant differences were found (see Section 3.2), post-hoc
2 pairwise comparisons were performed testing all combinations; a total of 3 tests were
3 performed to check for differences between different strategies.

4 ii. For each training strategy, pairwise comparisons were performed between Sinc-
5 ShallowNet-v2 and the other two CNNs that inspired Sinc-ShallowNet-v2 design: the
6 winning solution of the IFMBE 2019 challenge based on EEGNet [50], representing
7 the SOA for P300 decoding, and Sinc-ShallowNet [52]. A total of 6 tests were
8 performed to check for differences between the proposed ICNN and the other two
9 CNN-based solutions.

10 iii. Pairwise comparisons were performed between Sinc-ShallowNet-v2 trained with TL-
11 WS strategy and trained with WS-WS strategy for each number of calibration blocks
12 used for tuning the models. This was done to test the benefits of pre-training the models
13 on other subjects compared to training models from scratch. A total of 10 tests were
14 performed to check for differences.

15 All previous pairwise comparisons were performed using Wilcoxon signed-rank tests and
16 false discovery rate correction at $\alpha = 0.05$ using the Benjamini–Hochberg procedure [75] to
17 correct for multiple tests.

18 **2.8.2. ASD clinical scores**

19 Finally, we investigated whether the relevance of spatial features related to the most relevant
20 spectral features of the P300 response, as learned by the ICNN, was related to ASD clinical
21 scores (ADOS, ADI-R, and IQs, see Section 2.1). As the cluster analysis highlighted a strong
22 contribution of parietal sites across clusters (see Section 3.3), $\bar{\theta}_{spat}^{(s)}$ was averaged across P3,
23 Pz and P4 sites (thus, resulting in a scalar value for each subject s). Then, the Pearson
24 correlation coefficient between this ICNN+ET-derived measure and each ASD clinical score
25 was computed. Lastly, we performed a similar correlation analysis by using a measure derived

1 from subject-specific evoked potentials, to investigate whether useful ASD neural signatures
2 were enhanced by the proposed ICNN+ET analysis or they emerged already from evoked
3 potentials. To this aim, for each subject, EEG trials of the test set containing the P300 response
4 were band-pass filtered in the range $[f_0^*, f_1^*]$ and were averaged together. Then, the resulting
5 evoked potential was averaged across P3, Pz and P4 sites and across the time samples where
6 the P300 response was stronger, i.e., between 400-600 ms, see Borra et al. [43]. Thus, after this
7 averaging procedure, a scalar value was obtained for each subject s . Then, the Pearson
8 correlation coefficient between this ERP-derived measure and each ASD clinical score was
9 computed.

1 3. RESULTS

2 First, this section describes the results of AHPS via BO and the obtained decoding
3 performance; then, it shows the results obtained via the proposed ICNN+ET algorithm to
4 analyze the neural signatures corresponding to attentional P300 in ASD.

6 3.1. Hyper-parameter search via Bayesian optimization

7 The hyper-parameter configurations and their importance scores were separately obtained
8 for each tested training strategy (WS-WS, WS-CS, LOSO). Figure 4 shows the probability to
9 obtain a specific hyper-parameter value (P) among the admitted ones (Figure 4a) and the hyper-
10 parameter importance scores z (Figure 4b) for each training condition.

11 The more frequent design of the ISS block included a lower number of band-pass filters in
12 WS-WS than WS-CS and LOSO (i.e., $K_0^{ISS} = 8$ vs. $K_0^{ISS} = 16$) and the same number of
13 spatial filters for each band-pass filter across training conditions (i.e., $D_1^{ISS} = 4$). Among the
14 regularization techniques employed, spatial kernel max-norm constraint was applied more
15 frequently in the LOSO strategy, while was progressively applied less frequently (i.e., $c_1^{ISS} =$
16 *None*) moving to WS-CS and WS-WS strategies. The FST block was mostly designed learning
17 an equal representation over the input ISS feature maps (i.e., $K_0^{FST} = K_0^{ISS} \cdot D_1^{ISS}$) for WS-
18 WS and LOSO strategy, while equal representation and compressed representation (with
19 $K_0^{FST} = 2$) were almost evenly frequent in WS-CS strategy. Interestingly, the most frequent
20 temporal window in which features were learned in the FST block was approx. of 600 ms for
21 all training strategies ($= F_0^{FST} [1] / (sf // 4)$, where $F_0^{FST} = (1, 21)$ and $sf = 128$ Hz), and
22 this could be related to the temporal dynamics of the P300 response. Regarding the FC block,
23 weights were always constrained in case of LOSO strategy mainly with low upper bound (c_0^{FC}
24 $= 0.25$ or 0.5), while the probability distribution of this hyper-parameter appeared almost
25 uniform for the other strategies.

1 Dropout (both in the ISS and in the FST blocks) was mostly absent in LOSO strategy while
2 it was included in the other strategies with a more frequent dropout probability of $p = 0.25$.
3 Lastly, higher learning rates were adopted in WS-WS strategy, i.e., $lr \in [5 \cdot 1e - 3, 1e - 2)$,
4 compared to the other strategies where $lr \in [5 \cdot 1e - 4, 1e - 3)$.

5 As shown in Figure 4b, the learning rate resulted the most important hyper-parameter to
6 optimize in case of WS-WS strategy, while the other hyper-parameters resulted almost equally
7 important. Regarding the WS-CS strategy, the number of band-pass filters (K_0^{ISS}) and of
8 spatial filters for each band-pass filter (D_1^{ISS}) were the most important hyper-parameters, while
9 regarding the LOSO strategy, the number of band-pass filters (K_0^{ISS}) and the dropout
10 probability (p) were the most important ones.

11 Table 5 lists the model size (expressed as the number of trainable parameters) and training
12 time (expressed in s/epoch) of the models. As expected, WS-WS models were the lightest
13 (overall and within each block) and fastest to train, while LOSO models were the heaviest
14 (overall and within each block) and slowest to train. Interestingly, the proportion of trainable
15 parameters across blocks between WS-WS and WS-CS models did not change substantially
16 (ISS: 28%, 28%; FST: 62%, 63%; FC: 10%, 9%, respectively for WS-WS and WS-CS);
17 compared to WS-WS and WS-CS strategies, LOSO models presented a higher concentration
18 of the trainable parameters (i.e., higher capacity) in the FST block (i.e., ISS: 12%; FST: 80%;
19 FC: 8%).

20

21 **3.2. Decoding performance**

22 **3.2.1. Comparison with EEGNet and Sinc-ShallowNet in different training strategies**

23 The performance metrics averaged across subject-specific and session-specific test sets (see
24 Section 2.8.1) are reported here. Figure 5 displays the performance metrics of the Bayesian-
25 optimized EEGNet, the Bayesian-optimized Sinc-ShallowNet, and the Bayesian-optimized

1 Sinc-ShallowNet-v2. Sinc-ShallowNet-v2 scored significant different performance across the
2 three training strategies ($p < 0.001$, Friedmann test). As expected, post-hoc comparisons
3 highlighted that the network performed significantly worse in LOSO strategy than both WS-
4 WS and WS-CS ($p < 0.001$ for both comparisons) strategies. Lastly, the network in WS-CS
5 performed significantly better than WS-WS strategy ($p = 2.3 \cdot 1e - 3$). Sinc-ShallowNet-v2
6 scored object-level accuracies of $88.5 \pm 2.3\%$, $91.5 \pm 1.8\%$, $76.0 \pm 3.2\%$, compared to $85.5 \pm 2.5\%$,
7 $91.7 \pm 1.6\%$, $76.3 \pm 3.2\%$, of EEGNet, and $76.1 \pm 3.9\%$, $84.0 \pm 2.9\%$, $67.6 \pm 3.6\%$, of Sinc-
8 ShallowNet, respectively in WS-WS, WS-CS and LOSO strategies. Remarkably, despite the
9 reduction in model capacity by re-parametrizing part of the architecture to design an
10 interpretable spectral and spatial feature extractor, Sinc-ShallowNet-v2 achieved comparable
11 performance compared to EEGNet for P300 decoding in WS-CS and LOSO strategies ($p =$
12 0.84), while significantly outperformed it in the WS-WS strategy ($p = 0.01$). Lastly, it is
13 worth noticing that the updated version of Sinc-ShallowNet-v2 significantly outperformed
14 Sinc-ShallowNet in all training strategies.

15

16 **3.2.2. Transfer learning**

17 Transfer learning was adopted to test the capability of models to transfer the knowledge
18 from other subjects to a new one, for a more practical usage of the decoder in a BCI for ASD
19 treatment (i.e., reduction of BCI calibration times). To this aim, Figure 6 shows the accuracy
20 obtained while transferring the knowledge from the other participants to the held-out
21 participant (TL-WS strategy, see Section 2.5.1-iv), using a progressively increasing number of
22 calibration blocks of the held-out participant (from 2 to 20, where 20 corresponds to the entire
23 subject-specific calibration set, see Section 2.5.1-i, iv). The WS-WS performance are reported
24 together with TL-WS performance, to highlight the potential benefit of transfer learning
25 compared to a random initialization (proper of the WS-WS strategy). Remarkably, transfer

1 learning was found to be beneficial with significant improvements ($p < 0.01$) for all the tested
2 number of calibration blocks; using just 2 calibration blocks per subject, TL-WS reached an
3 object-level performance as high as $80.3 \pm 3.0\%$ with an average improvement as high as 37.1%
4 compared to WS-WS.

5

6 **3.3. Subject-specific ASD neural signatures related to P300**

7 Figure 7 (top panel) shows the spectral relevance (average \pm standard error of the mean
8 across subjects). The frequency range retaining most of the relevance (≥ 0.75) resulted
9 $[f_0^*, f_1^*] = [2, 5.8]$ Hz. Then, by clustering the relevance of spatial features related to this
10 frequency range, two clusters were obtained and their $\bar{\theta}_{c,spat}$ are reported in the bottom panel
11 of Figure 7 (left). Each of these cluster-level representations evidenced a peculiar strategy at
12 the level of scalp that the system automatically learned for P300 decoding. Cluster 0 showed a
13 strong and wide contribution of parietal sites (peaking at Pz) symmetrical across hemispheres,
14 while cluster 1 – that resulted the most populated cluster (with 10 observations) – showed a
15 strong right-lateralized contribution at parietal sites (P4). Lastly, in addition to these clusters,
16 two outliers were detected, with a different spatial relevance.

17 The results of the statistical analysis conducted on the relevance of spatial features ($\bar{\theta}_{spat}^{(s)}$)
18 at parietal sites is reported in Figure 8. Among all clinical scores, ADOS scores (both A-
19 Communication, B-Social Interaction and A+B-Communication-Social Interaction scores)
20 were the ones that mostly correlated to the ICNN+ET-derived measures, with high ($r=0.770$,
21 $r=0.657$ and $r=0.801$, respectively) and significant ($p < 0.01$) positive correlations.
22 Conversely, by using ERP-derived measures no significant correlations with clinical scores
23 were found.

24

1 4. DISCUSSION

2 4.1. Hyper-parameter search

3 Depending on the training strategy adopted, AHPS selected different hyper-parameter
4 configurations for Sinc-ShallowNet-v2. A peculiar trend of the probability distributions related
5 to K_0^{ISS} , D_1^{ISS} , K_0^{FST} , F_0^{FST} (see Figure 4a) was observed while moving from the WS-WS to
6 WS-CS and LOSO strategies. In particular, the probability distributions of these hyper-
7 parameters progressively moved from distributions more focused on small values (WS-WS) to
8 distributions more focused on higher values (LOSO) among the admitted ones (see Table 4 for
9 the admitted values), with the WS-CS strategy exhibiting an intermediate behavior in between
10 the previous strategies. As an example, the probability distributions related to K_0^{ISS} moved
11 from a distribution more focused on lower values including also the lowest admitted value
12 ($K_0^{ISS} = 4$), in the WS-WS strategy, to a distribution entirely focused on the highest admitted
13 value ($P(K_0^{ISS} = 16) = 1$), in the LOSO strategy. That is, Sinc-ShallowNet-v2, moving from
14 WS-WS to WS-CS and LOSO strategies, required progressively more filtered versions of the
15 input EEG trial (K_0^{ISS}), more electrode combinations tied to each band-pass filtered
16 representation (D_1^{ISS}), more temporal representations to learn (K_0^{FST}) and a wider window size
17 in which learn these representations (F_0^{FST}).

18 Among regularization techniques, the needing of dropout was stronger in WS-WS strategy,
19 using mostly $p = 0.5$ and $p = 0.25$, with a progressively reduction in WS-CS and LOSO
20 strategies, up to $p = 0$ corresponding to no dropout applied (denoted by $p = None$). That is,
21 the smaller and less variable (in terms of intra-subject and inter-subject variabilities) the
22 dataset, the higher the needing of dropout to provide a better generalization. Similar to other
23 regularization techniques, dropout acts reducing the algorithm capacity. Thus, from the
24 previous considerations about K_0^{ISS} , D_1^{ISS} , K_0^{FST} , F_0^{FST} and from the consideration on p , as
25 the training strategy involves a more challenging decoding task, the architecture needs more

1 capacity to solve the task with high performance. This is particularly relevant in strategies such
2 as WS-CS and LOSO, where single-trial EEG decoding is performed using signals collected
3 across several recording sessions (training set with high intra-subject variability) and across
4 several subjects and sessions (training set with high intra-subject and inter-subject variability),
5 respectively. Interestingly, this result was confirmed by looking to the hyper-parameter
6 importance scores (see red and light blue bars in Figure 4b), where K_0^{ISS} and D_1^{ISS} , and K_0^{ISS}
7 and p were the more important hyper-parameters to optimize for the WS-CS and LOSO
8 strategies, respectively. Furthermore, the need of an increased capacity as observed in the
9 structural hyper-parameters of Sinc-ShallowNet-v2 is reflected onto its model size (see Table
10 5), resulting in 1207, 1655, and 4638 (on average) trainable parameters for WS-WS, WS-CS,
11 and LOSO, respectively. In addition, while progressively increasing the variability in the
12 training examples, the capacity of the model increased differently across the network. While
13 increasing the intra-subject variability (WS-CS) the number of trainable parameters increased
14 by the same proportion across ISS, FST, FC blocks compared to WS-WS, suggesting that
15 increasing the capacity equally across the network could help addressing an increased intra-
16 subject variability in the training examples. Conversely, while also increasing the inter-subject
17 variability (LOSO), the network needed more abstract temporal features to learn in deeper
18 layers compared to the models trained in the other strategies, i.e., 80% (vs. ~60%) of the
19 parameters were in the FST block (see Table 5). This suggests that increasing the capacity at
20 the deeper layers of the architecture could help addressing an increased inter-subject
21 variability.

22 Lastly, the probability distribution of the learning rate was characterized by being more
23 biased towards higher learning rates in the WS-WS strategy and towards lower learning rates
24 in the LOSO strategy, with the WS-CS representing an intermediate condition. That is, as the
25 decoding task became less challenging, i.e., moving from the LOSO to WS-WS strategies, a

1 higher learning rate resulted beneficial. As for the previous hyper-parameters, this result was
2 confirmed by looking to the importance scores for the WS-WS strategy (Figure 4b), where the
3 learning rate was the most important hyper-parameter to search.

4 5 **4.2. Decoding performance**

6 In a preliminary analysis, we investigated the utility of AHPS, by comparing the
7 performance of architectures optimized via BO for each session or each subject instead of using
8 a fixed configuration. Results showed that the use of a fixed configuration provided
9 significantly lower accuracies (see Section 2 of Supplementary Materials). This suggests that
10 when a new subject or a subject in a new session approaches the BCI, BO should be performed
11 again to achieve higher performance.

12 Object-level accuracy varied using different training strategies. Despite the larger within-
13 subject variability in the input distributions (due to the involvement of many recording
14 sessions), CNNs trained with WS-CS were able to significantly outperform CNNs trained with
15 WS-WS, although the improvement is modest. A possible explanation for this result may be
16 that some subjects exhibited high variability even intra-session, leading to smaller performance
17 in WS-WS than in WS-CS (subjects 1,13,14 in Supplementary Figure 1), on average. Indeed,
18 in these cases WS-CS took great advantage of the larger training set across all sessions,
19 allowing object-level accuracy to be increased up to 8.6% compared to WS-WS. However,
20 WS-WS strategy reached high performance (i.e., >90%) with small variance in most of the
21 subjects (subjects 2, 4, 6, 7, 8, 10, 11, 15 in Supplementary Figure 1), suggesting consistency
22 in within-subject responses both intra- and inter-session. In these cases, WS-CS can still benefit
23 from the 7-time increased training set, slightly improving accuracy compared to WS-WS.
24 Lastly, further increasing the decoding difficulty, by including the cross-subject variability into

1 the training distributions, the challenge represented by the LOSO strategy was reflected into
2 the lowest performance across all the training strategies.

3 4 **4.2.1. Comparison with EEGNet and Sinc-ShallowNet in different training strategies**

5 Sinc-ShallowNet-v2 significantly outperformed the solution based on EEGNet [50] in the
6 WS-WS strategy, while it was comparable in the other strategies (WS-CS and LOSO). Due the
7 introduction of the temporal sinc-convolution in Sinc-ShallowNet-v2, the number of trainable
8 parameters to perform band-pass filtering were only 1.5% of the ones used in EEGNet that
9 exploits conventional temporal convolution. This reduction of trainable parameters provided a
10 beneficial effect when using the training strategy involving the most compact training set (WS-
11 WS), in which a limited algorithm capacity resulted optimal, while no effect was observed in
12 the conditions with larger and more variable training sets (WS-CS and LOSO), requiring a
13 higher algorithm capacity). Therefore, the introduction of the interpretable spectral and spatial
14 feature extractor, not only resulted beneficial for feature interpretability, but also did not
15 negatively affect the performance. Indeed, surprisingly the performance of Sinc-ShallowNet-
16 v2 was comparable or even significantly higher, while providing at the same time a
17 straightforward access to the spectral and spatial features. Furthermore Sinc-ShallowNet-v2
18 significantly outperformed the previous version Sinc-ShallowNet [52] in all training
19 conditions. These results suggest that the inclusion of the FST block in Sinc-ShallowNet-v2,
20 by learning temporal features on the output of ISS block, enables the extraction of more
21 relevant P300-related features than the ISS block only, easing the discrimination of the P300
22 response. In particular, the FST block may better cope with the large variability in time of
23 P300; the lower performance of Sinc-ShallowNet is then ascribable to the lack of this block
24 and to a reduced ability to catch intra-session, inter-session and inter-subject variability.

1 **4.2.2. Transfer learning**

2 Transfer learning could be used (TL-WS strategy) to reduce the calibration time on a new
3 subject recorded in a new recording session (i.e., requiring WS-WS decoding). Indeed, results
4 reported in Figure 6 highlight a significant benefit in transferring the knowledge from other
5 participants to a new one recorded in a new session, for each calibration set adopted.
6 Remarkably, the performance improvement was particularly relevant (on average 37.1%) in
7 the lowest data regime (i.e., 2 calibration blocks corresponding to 160 trials), suggesting that
8 the proposed interpretable approach was also capable of transferring the knowledge from other
9 participants enabling its usage with extremely compact-sized calibration sets. Lastly, transfer
10 learning not only improved the performance obtained in a WS-WS strategy using a small
11 portion of calibration blocks, which could be useful in practice to reduce calibration time, but
12 also improved the performance when using all 20 the calibration blocks (i.e., 1600 trials).

13

14 **4.3. Subject-specific ASD neural signatures related to P300**

15 The results obtained by analyzing the neural signatures related to P300 response via the
16 proposed ICNN+ET algorithm suggest that the more relevant (and shared across subjects)
17 features to distinguish P300 were obtained by filtering EEG signals including the frequency
18 range [2, 5.8] Hz in the passband. This is in line with other studies performing P300
19 classification both on healthy and ASD people, where features from EEG signals were
20 computed mainly in [2,4] Hz, [4,8] Hz, [2,8] Hz and [2,12] Hz frequency ranges [33–35]. In
21 these previous studies, EEG signals recorded from ASD people were pre-processed based on
22 other P300 decoding approaches validated on healthy subjects, adopting fixed cutoff
23 frequencies. Conversely, here the learning system was left free to explore all frequency
24 contents, adapting them to EEG signals of ASD people in an unbiased way. Indeed, some
25 learned filters included also higher cutoff frequencies too, but only the filters containing lower

1 frequency content resulted more important to provide the correct discrimination of the P300
2 response. Furthermore, the spectral relevance peaked approximatively in [2,4] Hz, matching
3 the findings of the wavelet analysis conducted by Demiralp et al. [36] where single EEG trials
4 were successfully decoded from features designed in delta ([2,4] Hz) frequency range. These
5 results suggest that future studies on traditional machine learning applications for P300
6 decoding in ASD, could design filter banks focusing especially on the interval [2, 5.8] Hz, i.e.,
7 in the delta and theta ranges.

8 The spatial relevance showed a strong right-lateralization for most of the observations (see
9 cluster 1 in Figure 7). This indicates that the visuo-spatial sensory processing during the BCI
10 task involved mainly the right hemisphere, differing from the classic P300 scalp distribution.
11 In the literature, there is evidence of hemispheric asymmetries underlying social perception
12 [76], e.g., right-lateralization while processing facial expressions related to emotions [77–79].
13 Furthermore, a right-lateralization was found also in the P300 response in Amaral et al. [80]
14 and was associated to the high-level characteristics in the realism of the animated paradigms
15 provided via the virtual environment (e.g., reflexive attention generated by social gaze
16 orientation). The BCI intervention investigated in this study was based on a visuo-spatial
17 oddball task which exploits a complex animated paradigm (see Section 2.1) and here this right-
18 lateralization emerged in the electrode discriminatory power related to P300, as learned by the
19 ICNN. Thus, our results further substantiate the idea that neural signatures related to social
20 perception (e.g., perception of gaze, faces and related gestures) are characterized by a peculiar
21 right-hemispheric asymmetry. Two clusters of spatial relevance for the P300 discrimination
22 were obtained (see Section 3.3), potentially highlighting a modulation of the right-hemispheric
23 asymmetry: from an asymmetric involvement of mainly P4 (cluster 1, the most populated – 10
24 subjects) to the symmetric involvement of mainly Pz, P3 and P4 (cluster 0 – 3 subjects).

1 From the correlation analysis between the ASD clinical scores and spatial relevance ($\bar{\theta}_{spat}^{(s)}$) of
2 parietal electrodes (P3, Pz, P4), a strong positive and significant correlation with ADOS scores
3 was observed (see Figure 8). This result is particularly interesting, as ADOS together with
4 ADI-R are two important assessment tools used to characterize and diagnose ASD [81]. In
5 addition, recent results [82] suggested that ADOS-related scores are more useful than the
6 combination of ADOS and ADI-R scores to diagnose ASD, which is in line with our finding.
7 Notably, at variance with the results obtained with the ICNN+ET-derived measure, no
8 significant correlations were found between ASD clinical scores and the P300 ERP-derived
9 measure.

10 Overall, we believe that the proposed method presents some points of novelty and strength,
11 compared to previous CNN-based approaches [41,43,45,48–50,57,58], with perspective not
12 only for practical engineering applications but also for theoretical neuroscience knowledge.

13 First, from a decoding perspective, we optimized our CNN-based decoders in their structure,
14 by tuning their hyper-parameters to achieve higher decoding performance, and this was done
15 separately for different training strategies. To the best of our knowledge, this is the first time
16 that such analysis is performed on P300 CNN-based decoders, using an automatic strategy
17 (Bayesian optimization) and in a training-specific way. Indeed, previous studies mainly
18 proposed CNN structures with hyper-parameters manually selected, without discussing the
19 specific choice [41,45,48–50,57], or performing sensitivity analysis only on few hyper-
20 parameters and for a single training strategy [43,58]. A significant result of our analysis is that
21 for each investigated training strategy (WS-WS, WS-CS, LOSO), different optimal hyper-
22 parameters were found, depending on the variability incorporated into the training examples,
23 and that a different architecture with a proper capacity should be adopted depending on the
24 specific real-life scenario in which the P300-BCI is used. Thus, the results obtained in this
25 study may help researchers in the design of CNN decoders for P300-based BCI applications,

1 by driving them in the choice of the more appropriate structure that takes into account the
2 adopted training strategy and the level of variability inside the training examples.

3 Second, we introduced an interpretable layer in our CNN decoders, using convolutional
4 kernels defined as function of only two and directly interpretable parameters per kernel. The
5 performance evaluation of our decoders compared to EEGNet (which is the state-of-the-art for
6 P300 decoding), indicates that this modification, while decreasing the number of trainable
7 parameters, did not reduce the decoding accuracy. Thus, a CNN can be made intrinsically more
8 interpretable without losing decoding capabilities. ICNNs are receiving growing interest in
9 EEG applications, moving beyond the use of CNNs as black-box decoders only, towards a
10 neurophysiological interpretation of the learned features. In this regard, it is highly important
11 to first verify the decoding performance of ICNNs, to ensure reliability of the learned features
12 as discriminative of the process under investigation, otherwise the neurophysiological
13 interpretation of these features would be unreliable. So far, ICNNs have been adopted in
14 literature to decode and analyze motor states [52,53] and not the P300 response and not even
15 in neurological disorder conditions such as ASD. Thus, the present study serves also for
16 validating ICNNs for P300 decoding in ASD subjects.

17 Third and related to previous point, by taking advantage of the interpretability embedded
18 into the proposed ICNN, we have formalized an analysis workflow that combines our ICNN
19 with an explanation technique (ET) to derive novel ASD biomarkers related to P300 by
20 exploiting the knowledge learned by the ICNN. The ICNN+ET combination was able to
21 capture and enhance, better than a canonical ERP analysis, meaningful spectral and spatial
22 characteristics underlying visuo-spatial sensory and attentional processing that are related to
23 autism. The enhancement of P300-related features is in line with the results of our recent study
24 [57] on healthy subjects, where representations derived from a CNN+ET framework better
25 highlighted P300 subcomponents (i.e., P3a and P3b) from single trial level, compared to ERP

1 canonical analysis. The remarkable aspect is that not only the ICNN-derived ASD markers
2 obtained here have a straightforward neurophysiological interpretation, but they also appeared
3 more sensitive than markers derived from traditional ERP-analysis.

4 Of course, the present study has also some limitations that can be the subject of future
5 improvements, and can foster further investigations along this line of research.

6 First, the optimal models obtained via Bayesian optimization may be task-related i.e., may
7 result optimal only for EEG decoding of P300 response elicited through the specific kind of
8 stimuli applied during the adopted oddball paradigm (visual stimuli presented inside a virtual
9 environment). It would be of great value to test these same models on other BCI-based P300
10 paradigms, i.e., using a different type of visual stimuli or a different sensory modality (acoustic
11 stimuli), or even on BCI paradigms based on event-related potentials other than P300. This
12 would test the ability of the proposed models to generalize across different BCI paradigms.
13 This kind of analysis is very important to promote a more efficient use of decoders, limiting
14 the proliferation of novel decoders for every new paradigm and task, but it is only rarely
15 performed (see [52] as a useful example of this analysis).

16 Another limitation concerns the limited number of subjects (15 subjects) included in the
17 examined dataset. Thus, the biomarkers derived from the ICNN+ET analysis that were found
18 to be correlated with ASD social impairment scores, although promising, require a more
19 extensive validation on a larger set of participants. Furthermore, in future research, the
20 workflow ICNN+ET, applied here to study P300 in autism, could be focused on analyzing
21 other ERP components in autism, and/or in other neurological and neurodevelopmental
22 disorders that involve ERP abnormalities (e.g., schizophrenia, depression, etc.) [83].

23

1 5. CONCLUSION

2 In conclusion, in this study we investigated the Bayesian-optimized design of an ICNN in
3 different training strategies when decoding the P300 response for BCI intervention in ASD.
4 While performing AHPS, models were found needing more capacity and lower learning rates
5 to decode EEG signals when more variability was included in the training distribution (i.e.,
6 including progressively intra-subject and inter-subject variabilities), i.e., depending on the
7 practical scenario in which the decoder is used in the BCI system. Despite its interpretable and
8 lightweight nature, the proposed ICNN performed as well as EEGNet, even significantly
9 outperforming it in the strategy with the lowest variability in the training examples (within-
10 subject and within-session). Furthermore, Sinc-ShallowNet-v2 significantly outperformed the
11 previous Sinc-ShallowNet design. Lastly, transferring the knowledge from other users to a new
12 one proved to lead to a substantial reduction of calibration times. All these results contribute
13 to the development of optimal decoders for P300-based BCI for ASD interventions, by
14 specifically improving CNN designs, performance, and calibration times. Furthermore, in this
15 study we leveraged the interpretable nature embedded into the ICNN to design an ICNN+ET
16 algorithm for the analysis of the visuo-spatial P300 in ASD in the frequency and spatial
17 domains. The analysis on spectral features matched known P300-related correlates, and while
18 analyzing spatial features, a right-hemispheric asymmetry was found, in line with the literature
19 of social perception. The modulation of this asymmetry, as provided by the ICNN+ET analysis,
20 was found to be correlated to ADOS scores, while no significant correlations with any clinical
21 score were found by using a simpler ERP analysis. This suggests that the ICNN+ET algorithm
22 was capable of better characterize and enhance useful ASD-related features than a canonical
23 analysis. In the future, the analysis on optimal ICNN designs for P300 decoding may be
24 extended to other oddball recording paradigms involving different stimuli properties.
25 Furthermore, the proposed ICNN+ET combination could be applied on more subjects for

- 1 further validation, generalized to other ERP components and neurological disorders to study
- 2 alterations in ERP components in a data-driven way, and also possibly extended on other
- 3 recording modalities of neural activity, e.g., magnetoencephalography.

1 **ACKNOWLEDGMENTS**

2 The authors declare no conflict of interest.

3 We gratefully acknowledge the support of NVIDIA Corporation with the donation of the
4 TITAN V used for this research. The provider was not involved in the study design, collection,
5 analysis, interpretation of data, the writing of this article or the decision to submit it for
6 publication.

7 The dataset adopted was collected in the scope of the European project BRAINTRAIN-FP7-
8 HEALTH-2013-INNOVATION-1-602186.

9

1 REFERENCES

- 2 [1] American Psychiatric Association 2013 *Diagnostic and Statistical Manual of Mental*
3 *Disorders* (American Psychiatric Association)
- 4 [2] Baron-Cohen S 1989 Perceptual role taking and protodeclarative pointing in autism
5 *British Journal of Developmental Psychology* **7** 113–27
- 6 [3] Baron-Cohen S, Baldwin D A and Crowson M 1997 Do Children with Autism Use the
7 Speaker’s Direction of Gaze Strategy to Crack the Code of Language? *Child*
8 *Development* **68** 48
- 9 [4] Swettenham J, Baron-Cohen S, Charman T, Cox A, Baird G, Drew A, Rees L and
10 Wheelwright S 1998 The Frequency and Distribution of Spontaneous Attention Shifts
11 between Social and Nonsocial Stimuli in Autistic, Typically Developing, and Nonautistic
12 Developmentally Delayed Infants *Journal of Child Psychology and Psychiatry* **39** 747–
13 53
- 14 [5] Klin A, Jones W, Schultz R, Volkmar F and Cohen D 2002 Visual Fixation Patterns
15 During Viewing of Naturalistic Social Situations as Predictors of Social Competence in
16 Individuals With Autism *Arch Gen Psychiatry* **59** 809
- 17 [6] Dawson G, Toth K, Abbott R, Osterling J, Munson J, Estes A and Liaw J 2004 Early
18 Social Attention Impairments in Autism: Social Orienting, Joint Attention, and Attention
19 to Distress. *Developmental Psychology* **40** 271–83
- 20 [7] Bakeman R and Adamson L B 1984 Coordinating Attention to People and Objects in
21 Mother-Infant and Peer-Infant Interaction *Child Development* **55** 1278
- 22 [8] Charman T 1998 Specifying the Nature and Course of the Joint Attention Impairment in
23 Autism in the Preschool Years: Implications for Diagnosis and Intervention *Autism* **2** 61–
24 79
- 25 [9] Charman T 2003 Why is joint attention a pivotal skill in autism? ed U Frith and E L Hill
26 *Phil. Trans. R. Soc. Lond. B* **358** 315–24
- 27 [10] Travers B G, Adluru N, Ennis C, Tromp D P M, Destiche D, Doran S, Bigler E D,
28 Lange N, Lainhart J E and Alexander A L 2012 Diffusion Tensor Imaging in Autism
29 Spectrum Disorder: A Review: Diffusion tensor imaging *Autism Res* **5** 289–313
- 30 [11] Minshew N J and Keller T A 2010 The nature of brain dysfunction in autism:
31 functional brain imaging studies *Current Opinion in Neurology* **23** 124–30
- 32 [12] Townsend J, Westerfield M, Leaver E, Makeig S, Jung T-P, Pierce K and
33 Courchesne E 2001 Event-related brain response abnormalities in autism: evidence for
34 impaired cerebello-frontal spatial attention networks *Cognitive Brain Research* **11** 127–
35 45
- 36 [13] Luck S J, Hillyard S A, Mouloua M, Woldorff M G, Clark V P and Hawkins H L
37 1994 Effects of spatial cuing on luminance detectability: Psychophysical and
38 electrophysiological evidence for early selection. *Journal of Experimental Psychology:*
39 *Human Perception and Performance* **20** 887–904

- 1 [14] Heinze H J, Luck S J, Mangun G R and Hillyard S A 1990 Visual event-related
2 potentials index focused attention within bilateral stimulus arrays. I. Evidence for early
3 selection *Electroencephalography and Clinical Neurophysiology* **75** 511–27
- 4 [15] Herrmann C S and Knight R T 2001 Mechanisms of human attention: event-related
5 potentials and oscillations *Neuroscience & Biobehavioral Reviews* **25** 465–76
- 6 [16] Sokhadze E, Baruth J, Tasman A, Sears L, Mathai G, El-Baz A and Casanova M F
7 2009 Event-related Potential Study of Novelty Processing Abnormalities in Autism *Appl*
8 *Psychophysiol Biofeedback* **34** 37–51
- 9 [17] Polich J 2007 Updating P300: An integrative theory of P3a and P3b *Clinical*
10 *Neurophysiology* **118** 2128–48
- 11 [18] Azizian A and Polich J 2007 Evidence for Attentional Gradient in the Serial Position
12 Memory Curve from Event-related Potentials *Journal of Cognitive Neuroscience* **19**
13 2071–81
- 14 [19] Mecklinger A and Pfeifer E 1996 Event-related potentials reveal topographical and
15 temporal distinct neuronal activation patterns for spatial and object working memory
16 *Cognitive Brain Research* **4** 211–24
- 17 [20] Cui T, Wang P P, Liu S and Zhang X 2017 P300 amplitude and latency in autism
18 spectrum disorder: a meta-analysis *Eur Child Adolesc Psychiatry* **26** 177–90
- 19 [21] Ciesielski K T, Courchesne E and Elmasian R 1990 Effects of focused selective
20 attention tasks on event-related potentials in autistic and normal individuals
21 *Electroencephalography and Clinical Neurophysiology* **75** 207–20
- 22 [22] Courchesne E, Kilman B A, Galambos R and Lincoln A J 1984 Autism: Processing
23 of novel auditory information assessed by event-related brain potentials
24 *Electroencephalography and Clinical Neurophysiology/Evoked Potentials Section* **59**
25 238–48
- 26 [23] Courchesne E, Lincoln A J, Kilman B A and Galambos R 1985 Event-related brain
27 potential correlates of the processing of novel visual and auditory information in autism *J*
28 *Autism Dev Disord* **15** 55–76
- 29 [24] Courchesne E, Lincoln A J, Yeung-Courchesne R, Elmasian R and Grillon C 1989
30 Pathophysiologic findings in nonretarded autism and receptive developmental language
31 disorder *J Autism Dev Disord* **19** 1–17
- 32 [25] Dawson G, Finley C, Phillips S, Galpert L and Lewy A 1988 Reduced P3 amplitude
33 of the event-related brain potential: Its relationship to language ability in autism *J*
34 *Autism Dev Disord* **18** 493–504
- 35 [26] Verbaten M N, Roelofs J W, van Engeland H, Kenemans J K and Slangen J L 1991
36 Abnormal visual event-related potentials of autistic children *J Autism Dev Disord* **21**
37 449–70
- 38 [27] Kemner C, van der Gaag R J, Verbaten M and van Engeland H 1999 ERP differences
39 among subtypes of pervasive developmental disorders *Biological Psychiatry* **46** 781–9

- 1 [28] Amaral C, Mouga S, Simões M, Pereira H C, Bernardino I, Quental H, Playle R,
2 McNamara R, Oliveira G and Castelo-Branco M 2018 A Feasibility Clinical Trial to
3 Improve Social Attention in Autistic Spectrum Disorder (ASD) Using a Brain Computer
4 Interface *Frontiers in Neuroscience* **12** 477
- 5 [29] Amaral C P, Simões M A, Mouga S, Andrade J and Castelo-Branco M 2017 A novel
6 Brain Computer Interface for classification of social joint attention in autism and
7 comparison of 3 experimental setups: A feasibility study *Journal of Neuroscience*
8 *Methods* **290** 105–15
- 9 [30] Yger F, Berar M and Lotte F 2017 Riemannian Approaches in Brain-Computer
10 Interfaces: A Review *IEEE Trans. Neural Syst. Rehabil. Eng.* **25** 1753–62
- 11 [31] Saha S and Baumert M 2020 Intra- and Inter-subject Variability in EEG-Based
12 Sensorimotor Brain Computer Interface: A Review *Front. Comput. Neurosci.* **13** 87
- 13 [32] Lotte F, Bougrain L, Cichocki A, Clerc M, Congedo M, Rakotomamonjy A and Yger
14 F 2018 A review of classification algorithms for EEG-based brain–computer interfaces: a
15 10 year update *Journal of Neural Engineering* **15** 031005
- 16 [33] de Arancibia L, Sánchez-González P, Gómez E J, Hernando M E and Oropesa I 2020
17 Linear vs Nonlinear Classification of Social Joint Attention in Autism Using VR P300-
18 Based Brain Computer Interfaces *XV Mediterranean Conference on Medical and*
19 *Biological Engineering and Computing – MEDICON 2019* ed J Henriques, N Neves and
20 P de Carvalho (Cham: Springer International Publishing) pp 1869–74
- 21 [34] Krzemiński D, Michelmann S, Treder M and Santamaria L 2020 Classification of
22 P300 Component Using a Riemannian Ensemble Approach *XV Mediterranean*
23 *Conference on Medical and Biological Engineering and Computing – MEDICON 2019*
24 ed J Henriques, N Neves and P de Carvalho (Cham: Springer International Publishing) pp
25 1885–9
- 26 [35] Chatterjee B, Palaniappan R and Gupta C N 2020 Performance Evaluation of
27 Manifold Algorithms on a P300 Paradigm Based Online BCI Dataset *XV Mediterranean*
28 *Conference on Medical and Biological Engineering and Computing – MEDICON 2019*
29 IFMBE Proceedings vol 76, ed J Henriques, N Neves and P de Carvalho (Cham: Springer
30 International Publishing) pp 1894–8
- 31 [36] Demiralp T, Ademoglu A, Istefanopulos Y, Başar-Eroglu C and Başar E 2001
32 Wavelet analysis of oddball P300 *International Journal of Psychophysiology* **39** 221–7
- 33 [37] Bostanov V and Kotchoubey B 2006 The t-CWT: A new ERP detection and
34 quantification method based on the continuous wavelet transform and Student’s t-
35 statistics *Clinical Neurophysiology* **117** 2627–44
- 36 [38] Bittencourt-Villalpando M and Maurits N M 2018 Stimuli and Feature Extraction
37 Algorithms for Brain-Computer Interfaces: A Systematic Comparison *IEEE Trans.*
38 *Neural Syst. Rehabil. Eng.* **26** 1669–79
- 39 [39] Bittencourt-Villalpando M and Maurits N M 2020 Linear SVM Algorithm
40 Optimization for an EEG-Based Brain-Computer Interface Used by High Functioning
41 Autism Spectrum Disorder Participants *XV Mediterranean Conference on Medical and*

- 1 *Biological Engineering and Computing – MEDICON 2019* ed J Henriques, N Neves and
2 P de Carvalho (Cham: Springer International Publishing) pp 1875–84
- 3 [40] Miladinović A, Ajčević M, Battaglini P P, Silveri G, Ciacchi G, Morra G,
4 Jarmolowska J and Accardo A 2020 Slow Cortical Potential BCI Classification Using
5 Sparse Variational Bayesian Logistic Regression with Automatic Relevance
6 Determination *XV Mediterranean Conference on Medical and Biological Engineering
7 and Computing – MEDICON 2019* IFMBE Proceedings vol 76, ed J Henriques, N Neves
8 and P de Carvalho (Cham: Springer International Publishing) pp 1853–60
- 9 [41] Simões M, Borra D, Santamaría-Vázquez E, GBT-UPM, Bittencourt-Villalpando M,
10 Krzemiński D, Miladinović A, Neural_Engineering_Group, Schmid T, Zhao H, Amaral
11 C, Direito B, Henriques J, Carvalho P and Castelo-Branco M 2020 BCIAUT-P300: A
12 Multi-Session and Multi-Subject Benchmark Dataset on Autism for P300-Based Brain-
13 Computer-Interfaces *Front. Neurosci.* **14** 568104
- 14 [42] Craik A, He Y and Contreras-Vidal J L 2019 Deep learning for
15 electroencephalogram (EEG) classification tasks: a review *J. Neural Eng.* **16** 031001
- 16 [43] Borra D, Fantozzi S and Magosso E 2021 A Lightweight Multi-Scale Convolutional
17 Neural Network for P300 Decoding: Analysis of Training Strategies and Uncovering of
18 Network Decision *Frontiers in Human Neuroscience*
- 19 [44] Lindsay G 2020 Convolutional Neural Networks as a Model of the Visual System:
20 Past, Present, and Future *Journal of Cognitive Neuroscience* 1–15
- 21 [45] Cecotti H and Graser A 2011 Convolutional Neural Networks for P300 Detection
22 with Application to Brain-Computer Interfaces *IEEE Transactions on Pattern Analysis
23 and Machine Intelligence* **33** 433–45
- 24 [46] Srivastava N, Hinton G, Krizhevsky A, Sutskever I and Salakhutdinov R 2014
25 Dropout: a simple way to prevent neural networks from overfitting *The Journal of
26 Machine Learning Research* **15** 1929–58
- 27 [47] Ioffe S and Szegedy C 2015 Batch Normalization: Accelerating Deep Network
28 Training by Reducing Internal Covariate Shift *Proceedings of the 32nd International
29 Conference on Machine Learning* Proceedings of Machine Learning Research vol 37, ed
30 F Bach and D Blei (Lille, France: PMLR) pp 448–56
- 31 [48] Manor R and Geva A B 2015 Convolutional Neural Network for Multi-Category
32 Rapid Serial Visual Presentation BCI *Frontiers in Computational Neuroscience* **9** 146
- 33 [49] Lawhern V J, Solon A J, Waytowich N R, Gordon S M, Hung C P and Lance B J
34 2018 EEGNet: a compact convolutional neural network for EEG-based brain–computer
35 interfaces *Journal of Neural Engineering* **15** 056013
- 36 [50] Borra D, Fantozzi S and Magosso E 2020 Convolutional Neural Network for a P300
37 Brain-Computer Interface to Improve Social Attention in Autistic Spectrum Disorder *XV
38 Mediterranean Conference on Medical and Biological Engineering and Computing –
39 MEDICON 2019* ed J Henriques, N Neves and P de Carvalho (Cham: Springer
40 International Publishing) pp 1837–43

- 1 [51] Zhang X, Yao L, Wang X, Monaghan J, McAlpine D and Zhang Y 2021 A survey on
2 deep learning-based non-invasive brain signals: recent advances and new frontiers *J.*
3 *Neural Eng.* **18** 031002
- 4 [52] Borra D, Fantozzi S and Magosso E 2020 Interpretable and lightweight convolutional
5 neural network for EEG decoding: Application to movement execution and imagination
6 *Neural Networks* S0893608020302021
- 7 [53] Zhao D, Tang F, Si B and Feng X 2019 Learning joint space–time–frequency
8 features for EEG decoding on small labeled data *Neural Networks* **114** 67–77
- 9 [54] Liang S, Hang W, Yin M, Shen H, Wang Q, Qin J, Choi K-S and Zhang Y 2022
10 Deep EEG feature learning via stacking common spatial pattern and support matrix
11 machine *Biomedical Signal Processing and Control* **74** 103531
- 12 [55] Borra D, Fantozzi S and Magosso E 2020 EEG Motor Execution Decoding via
13 Interpretable Sinc-Convolutional Neural Networks *XV Mediterranean Conference on*
14 *Medical and Biological Engineering and Computing – MEDICON 2019* ed J Henriques,
15 N Neves and P de Carvalho (Cham: Springer International Publishing) pp 1113–22
- 16 [56] Vahid A, Mückschel M, Stober S, Stock A-K and Beste C 2020 Applying deep
17 learning to single-trial EEG data provides evidence for complementary theories on action
18 control *Commun Biol* **3** 112
- 19 [57] Borra D and Magosso E 2021 Deep learning-based EEG analysis: investigating P3
20 ERP components *Journal of Integrative Neuroscience* **In press**
- 21 [58] Farahat A, Reichert C, Sweeney-Reed C and Hinrichs H 2019 Convolutional neural
22 networks for decoding of covert attention focus and saliency maps for EEG feature
23 visualization *Journal of Neural Engineering*
- 24 [59] Wu W, Zhang Y, Jiang J, Lucas M V, Fonzo G A, Rolle C E, Cooper C, Chin-Fatt C,
25 Krepel N, Cornelissen C A, Wright R, Toll R T, Trivedi H M, Monuszko K, Caudle T L,
26 Sarhadi K, Jha M K, Trombello J M, Deckersbach T, Adams P, McGrath P J, Weissman
27 M M, Fava M, Pizzagalli D A, Arns M, Trivedi M H and Etkin A 2020 An
28 electroencephalographic signature predicts antidepressant response in major depression
29 *Nat Biotechnol* **38** 439–47
- 30 [60] Lord C, Risi S, Lambrecht L, Cook, Jr. E H, Leventhal B L, DiLavore P C, Pickles A
31 and Rutter M 2000 The Autism Diagnostic Observation Schedule—Generic: A Standard
32 Measure of Social and Communication Deficits Associated with the Spectrum of Autism
33 *Journal of Autism and Developmental Disorders* **30** 205–23
- 34 [61] Lord C, Rutter M and Le Couteur A 1994 Autism Diagnostic Interview-Revised: A
35 revised version of a diagnostic interview for caregivers of individuals with possible
36 pervasive developmental disorders *J Autism Dev Disord* **24** 659–85
- 37 [62] Spek A A, Scholte E M and van Berckelaer-Onnes I A 2008 Brief Report: The Use
38 of WAIS-III in Adults with HFA and Asperger Syndrome *J Autism Dev Disord* **38** 782–7

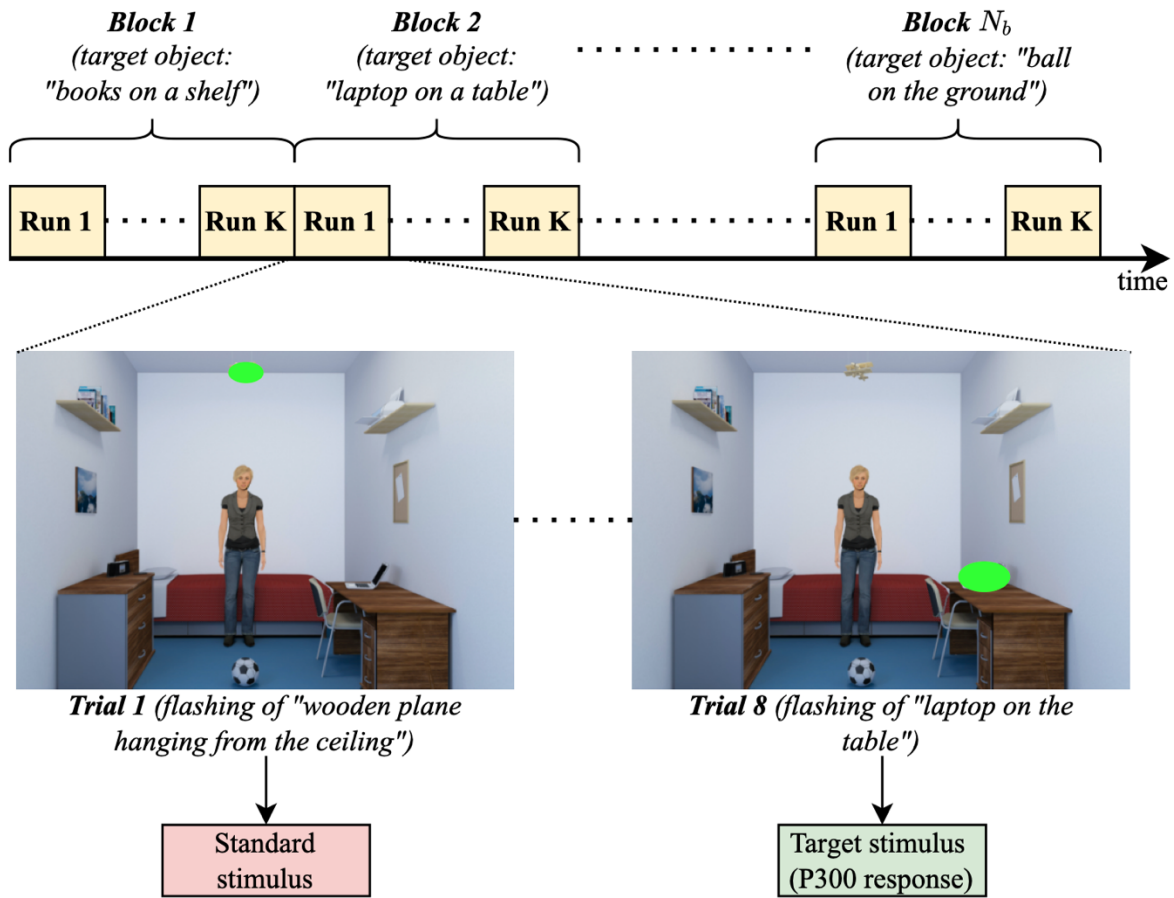
- 1 [63] Schirrmester R T, Springenberg J T, Fiederer L D J, Glasstetter M, Eggenberger K,
2 Tangermann M, Hutter F, Burgard W and Ball T 2017 Deep learning with convolutional
3 neural networks for EEG decoding and visualization *Human brain mapping* **38** 5391–420
- 4 [64] Ravanelli M and Bengio Y 2018 Speaker Recognition from Raw Waveform with
5 SincNet *2018 IEEE Spoken Language Technology Workshop (SLT)* pp 1021–8
- 6 [65] Clevert D-A, Unterthiner T and Hochreiter S 2015 Fast and accurate deep network
7 learning by exponential linear units (elus) *arXiv preprint*
- 8 [66] Chollet F 2016 Xception: Deep Learning with Depthwise Separable Convolutions
9 *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 1800–7
- 10 [67] Zhuang F, Qi Z, Duan K, Xi D, Zhu Y, Zhu H, Xiong H and He Q 2020 A
11 Comprehensive Survey on Transfer Learning *arXiv:1911.02685 [cs, stat]*
- 12 [68] Kingma D P and Ba J 2017 Adam: A Method for Stochastic Optimization
13 *arXiv:1412.6980 [cs]*
- 14 [69] Bergstra J, Bardenet R, Bengio Y and Kégl B 2011 Algorithms for hyper-parameter
15 optimization *Advances in neural information processing systems* **24**
- 16 [70] Akiba T, Sano S, Yanase T, Ohta T and Koyama M 2019 Optuna: A Next-generation
17 Hyperparameter Optimization Framework *arXiv:1907.10902 [cs, stat]*
- 18 [71] Frank Hutter, Holger Hoos, and Kevin Leyton-Brown 2014 An Efficient Approach
19 for Assessing Hyperparameter Importance *Proceedings of the 31st International*
20 *Conference on Machine Learning* ed Eric P. Xing and Tony Jebara (PMLR) pp 754–62
- 21 [72] Simonyan K, Vedaldi A and Zisserman A 2014 Deep Inside Convolutional
22 Networks: Visualising Image Classification Models and Saliency Maps *arXiv:1312.6034*
23 *[cs]*
- 24 [73] Campello R J G B, Moulavi D, Zimek A and Sander J 2015 Hierarchical Density
25 Estimates for Data Clustering, Visualization, and Outlier Detection *ACM Trans. Knowl.*
26 *Discov. Data* **10** 1–51
- 27 [74] McInnes L, Healy J and Astels S 2017 hdbSCAN: Hierarchical density based clustering
28 *JOSS* **2** 205
- 29 [75] Benjamini Y and Hochberg Y 1995 Controlling the False Discovery Rate: A
30 Practical and Powerful Approach to Multiple Testing *Journal of the Royal Statistical*
31 *Society. Series B (Methodological)* **57** 289–300
- 32 [76] Brancucci A, Lucci G, Mazzatenta A and Tommasi L 2009 Asymmetries of the
33 human social brain in the visual, auditory and chemical modalities *Phil. Trans. R. Soc. B*
34 **364** 895–914
- 35 [77] Hartikainen K M 2021 Emotion-Attention Interaction in the Right Hemisphere *Brain*
36 *Sciences* **11** 1006

- 1 [78] Demaree H A, Everhart D E, Youngstrom E A and Harrison D W 2005 Brain
2 Lateralization of Emotional Processing: Historical Roots and a Future Incorporating
3 “Dominance” *Behavioral and Cognitive Neuroscience Reviews* **4** 3–20
- 4 [79] Simões M, Monteiro R, Andrade J, Mouga S, França F, Oliveira G, Carvalho P and
5 Castelo-Branco M 2018 A Novel Biomarker of Compensatory Recruitment of Face
6 Emotional Imagery Networks in Autism Spectrum Disorder *Front. Neurosci.* **12** 791
- 7 [80] Amaral C P, Simões M A and Castelo-Branco M S 2015 Neural Signals Evoked by
8 Stimuli of Increasing Social Scene Complexity Are Detectable at the Single-Trial Level
9 and Right Lateralized ed S Ben Hamed *PLoS ONE* **10** e0121970
- 10 [81] Hiremath C S, Sagar K J V, Yamini B K, Girimaji A S, Kumar R, Sravanti S L,
11 Padmanabha H, Vykunta Raju K N, Kishore M T, Jacob P, Saini J, Bharath R D,
12 Seshadri S P and Kumar M 2021 Emerging behavioral and neuroimaging biomarkers for
13 early and accurate characterization of autism spectrum disorders: a systematic review
14 *Transl Psychiatry* **11** 42
- 15 [82] Kamp-Becker I, Tauscher J, Wolff N, Küpper C, Poustka L, Roepke S, Roessner V,
16 Heider D and Stroth S 2021 Is the Combination of ADOS and ADI-R Necessary to
17 Classify ASD? Rethinking the “Gold Standard” in Diagnosing ASD *Front. Psychiatry* **12**
18 727308
- 19 [83] Picton T W 1992 The P300 Wave of the Human Event-Related Potential *Journal of*
20 *Clinical Neurophysiology* **9** 456–79

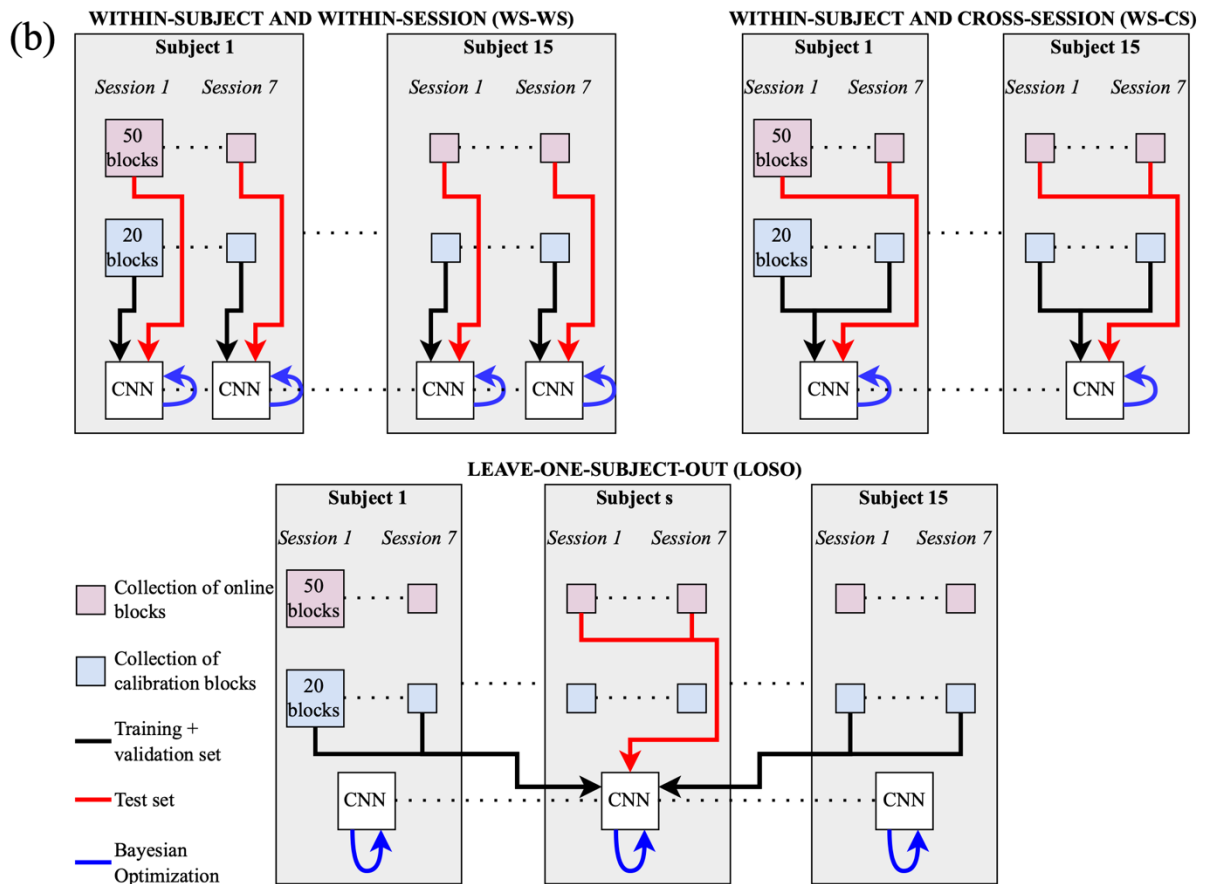
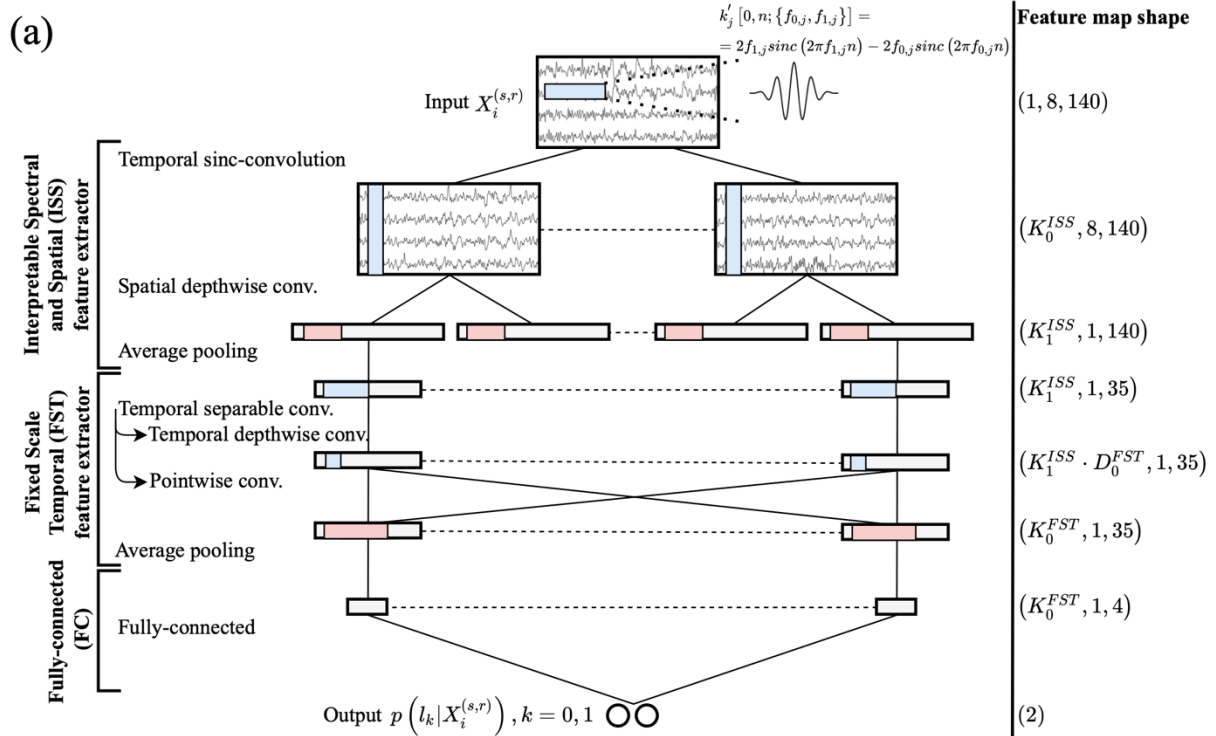
21

22

1 **FIGURES**

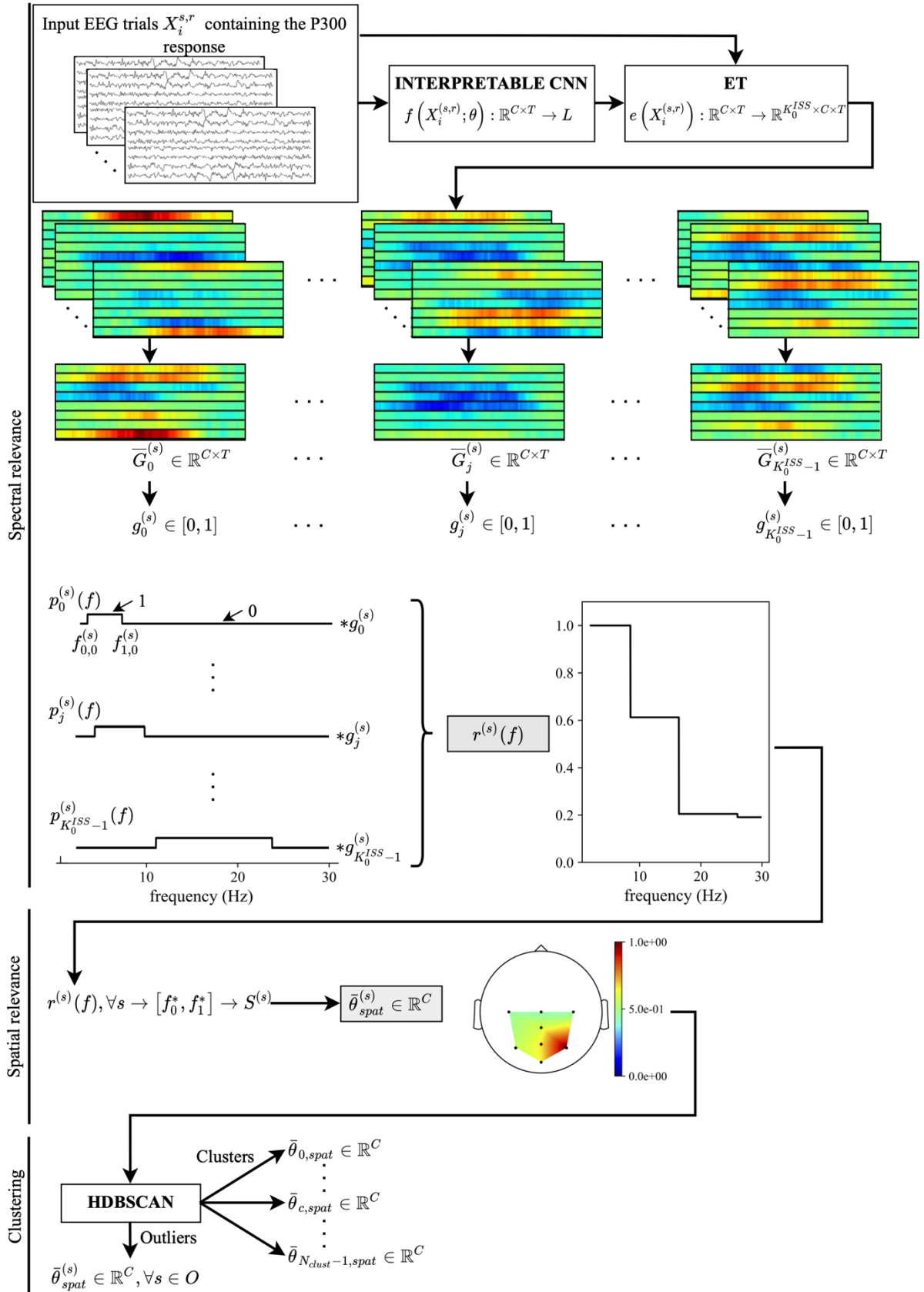


2
3 **Figure 1** – Structure of the BCI paradigm in its division into blocks, runs, and trials. Each
4 recording block was used to identify one of the eight target objects visualized in the virtual
5 environment, with a total of $N_b = 20$ (calibration phase) and of $N_b = 50$ (online phase) blocks.
6 Within each block, K runs were recorded, where $K = 10$ (calibration phase) and $K = 7.095$
7 (online phase) on average across subjects and sessions. In each run, the 8 objects randomly
8 flashed (schematized by the green ellipses in the figure) and 8 EEG trials per run where
9 recorded. Therefore, overall, within each recording session $N_b \cdot K \cdot 8$ EEG trials were recorded,
10 corresponding 1600 trials in the calibration phase and 2838 trials on average in the online
11 phase, respectively.

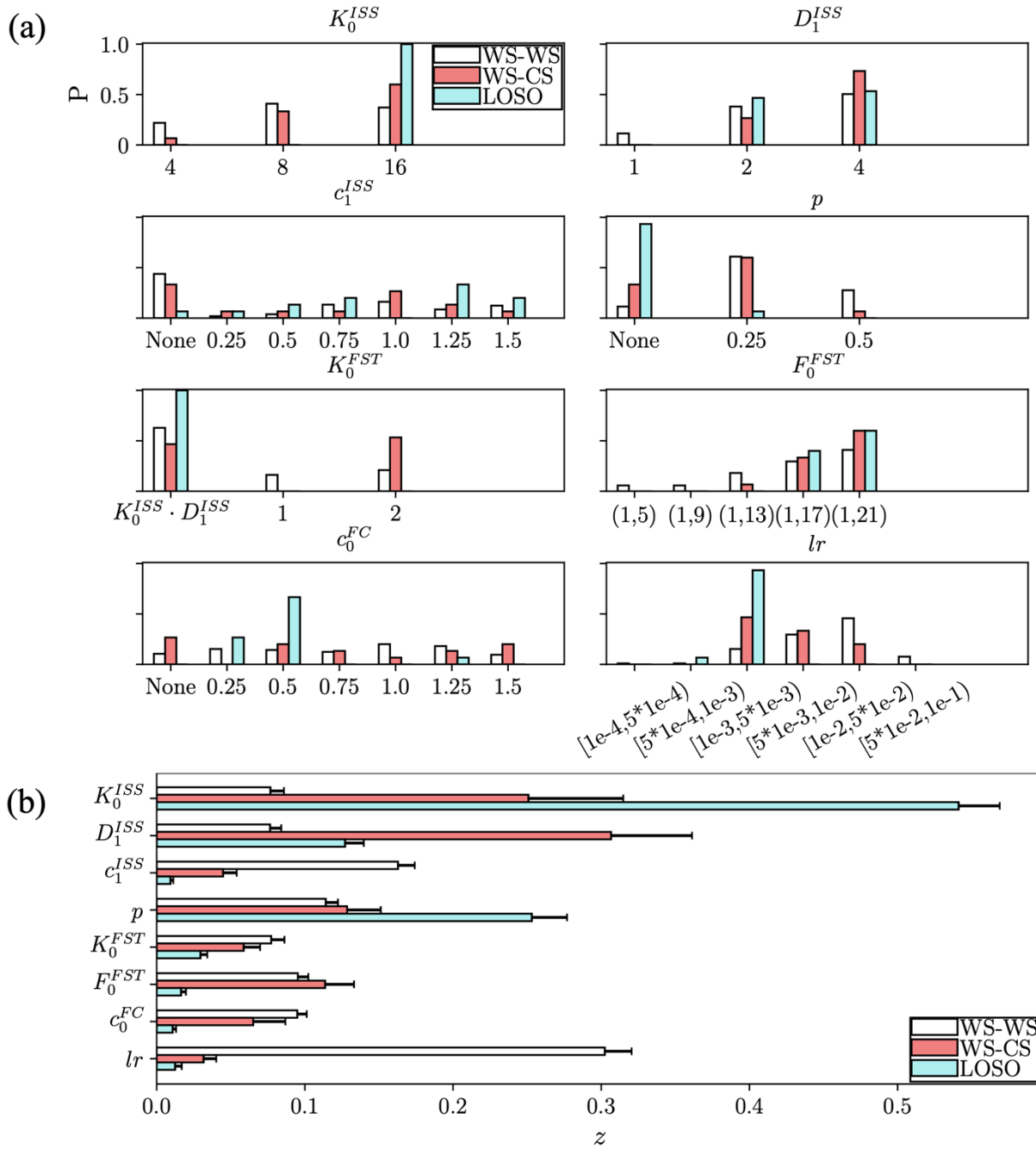


1
2 **Figure 2** – Sinc-ShallowNet-v2 (a) and training strategies investigated while performing
3 Bayesian optimization (b). In Figure 2a, blocks and main layers are listed on the left side. To
4 keep the figure as clear and simple as possible, only the main trainable layers (i.e.,
5 convolutional and fully-connected layers) in addition to pooling layers (to highlight the
6 temporal dimension reduction) were displayed. Boxes represent the output feature maps of

1 each layer, and colored rectangles represent convolutional (blue) and pooling (red) kernels.
2 Tuples reported on the right side represent the shape of the feature maps. For all outputs except
3 the last, tuples are composed by three numbers representing the number of the feature maps,
4 the number of spatial samples and the number of temporal samples within each map,
5 respectively. The input layer provides an output of shape $(1, C, T) = (1, 8, 140)$, as it just
6 replicates the original input EEG trial with shape $(8, 140)$, providing a single feature map as
7 output. The temporal dimension changed from $T = 140$ to $T//4 = 35$ and to $T//32 = 4$ along
8 the entire CNN due to the average pooling operations (where $//$ indicates the floor division
9 operator). See Section 2.2, 2.3 and Table 2 for further details. Figure 2b shows a schematic
10 representation of how training, validation, and test examples (by means of black and red
11 arrows) were sampled from calibration (blue boxes) and online (purple boxes) blocks recorded
12 in the BCI paradigm when training decoders in Bayesian optimization in within-subject and
13 within-session (WS-WS), within-subject and cross-session (WS-CS), and leave-one-subject-
14 out (LOSO) training conditions. For brevity, in the LOSO strategy the aggregation across
15 blocks is reported only for the s -th subject. See Section 2.1 and Section 2.5.1 for further details
16 about the definition of recording blocks and the training strategies, respectively.
17

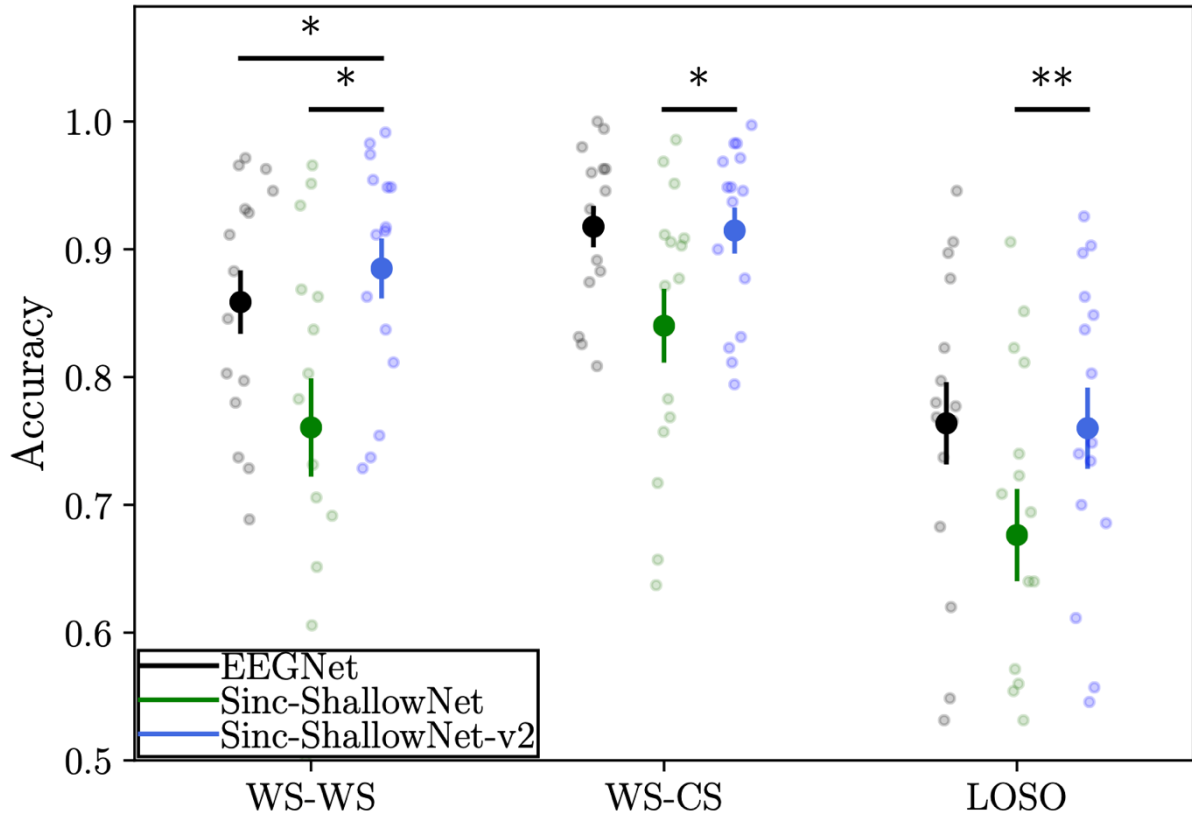


1
2 **Figure 3** – Schematic representation of the algorithm based on ICNN+ET adopted to gain
3 insights about the neural signatures of the visuo-spatial P300 correlate in autism in its three
4 main steps: spectral relevance computation, spatial relevance computation, and clustering.
5

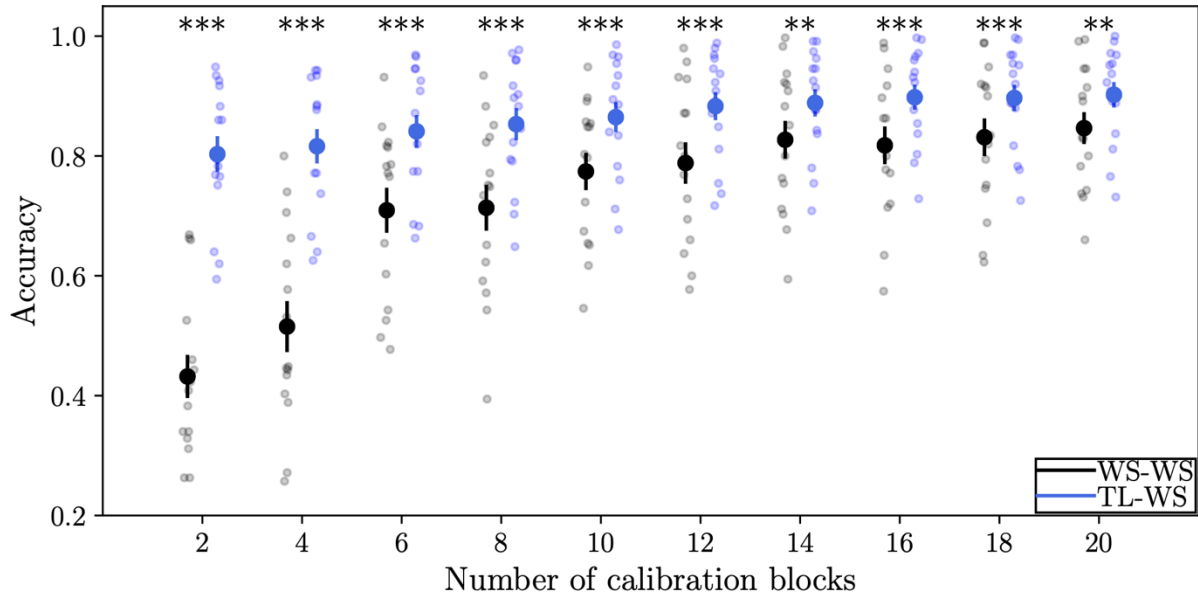


1
2 **Figure 4** – Hyper-parameter distributions (a) and importance scores (b) when training Sinc-
3 ShallowNet-v2 with within-subject and within-session (WS-WS), within-subject and cross-
4 session (WS-CS) and leave-one-subject-out (LOSO) strategies. Hyper-parameter distributions
5 were reported representing the probability (P) that a specific hyper-parameter value resulted as
6 optimal via BO, and it is reported separately for the different training strategies. Hyper-
7 parameter importance scores were reported in their mean values across decoders trained with
8 each training strategy, separately. Bar heights represent the mean value, while the error bars
9 represent the standard error of the mean.

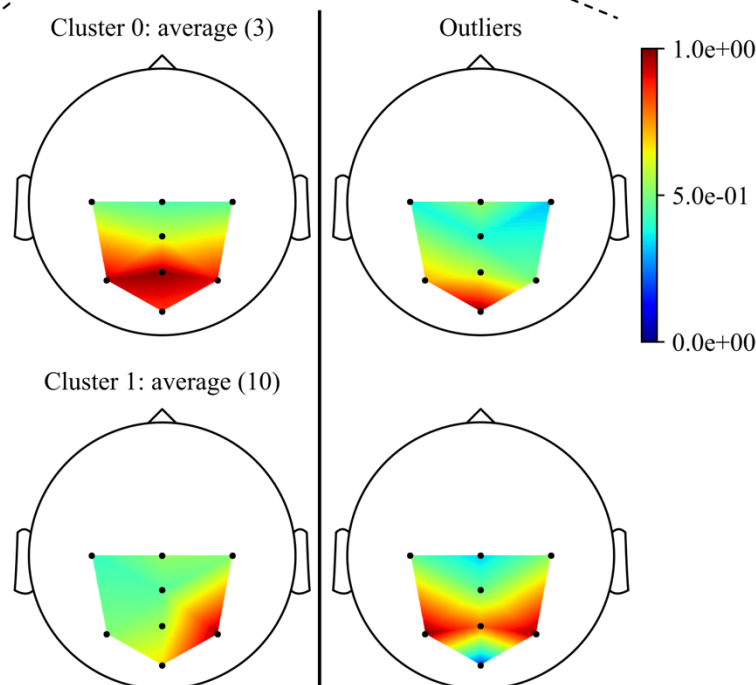
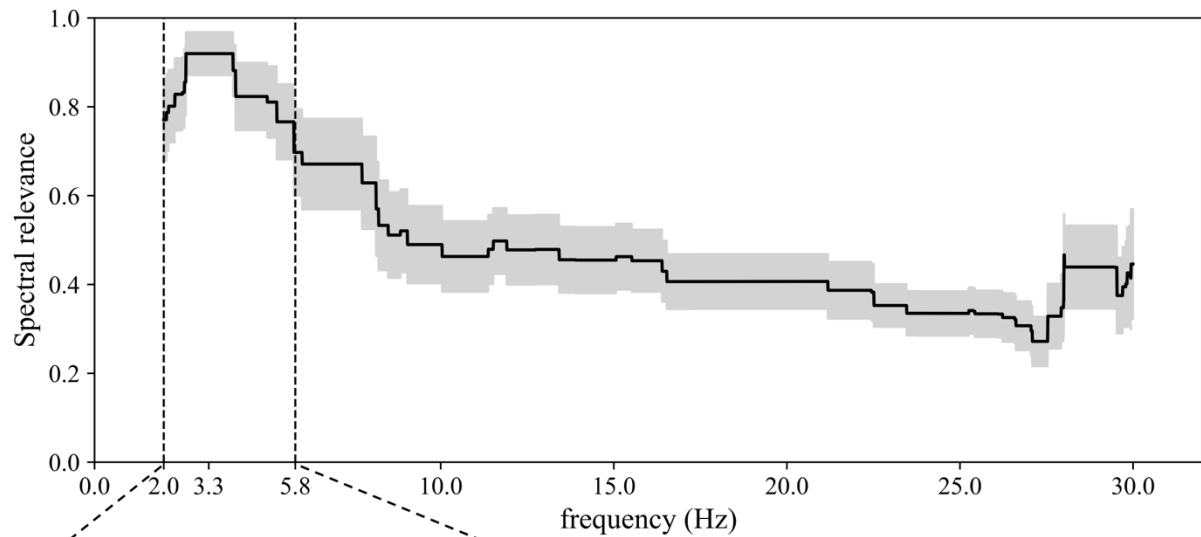
10



1
2 **Figure 5** – Accuracies of the Bayesian-optimized EEGNet (black), Bayesian-optimized Sinc-
3 Sinc-ShallowNet (green) and Bayesian-optimized Sinc-ShallowNet-v2 (blue) scored in the WS-W,
4 WS-CS and LOSO strategies. Bigger dots represent the mean value across subjects, while the
5 error bars represent the standard error of the mean. Smaller dots represent the accuracy scored
6 for each subject (i.e., 15 data points). See part A of Table 3 for details about the number of
7 examples defining the training, validation and test sets. Wilcoxon signed-rank tests were
8 performed to compare the performance of the Bayesian-optimized Sinc-ShallowNet-v2 with
9 the Bayesian-optimized EEGNet and with the Bayesian-optimized Sinc-ShallowNet, within
10 each training strategy. P-values were corrected for multiple comparisons (6 in total) via the
11 Benjamini–Hochberg procedure and significant comparisons are marked once applied the
12 correction (* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$). See Section 2.8.1 for further details about the
13 statistical tests.
14

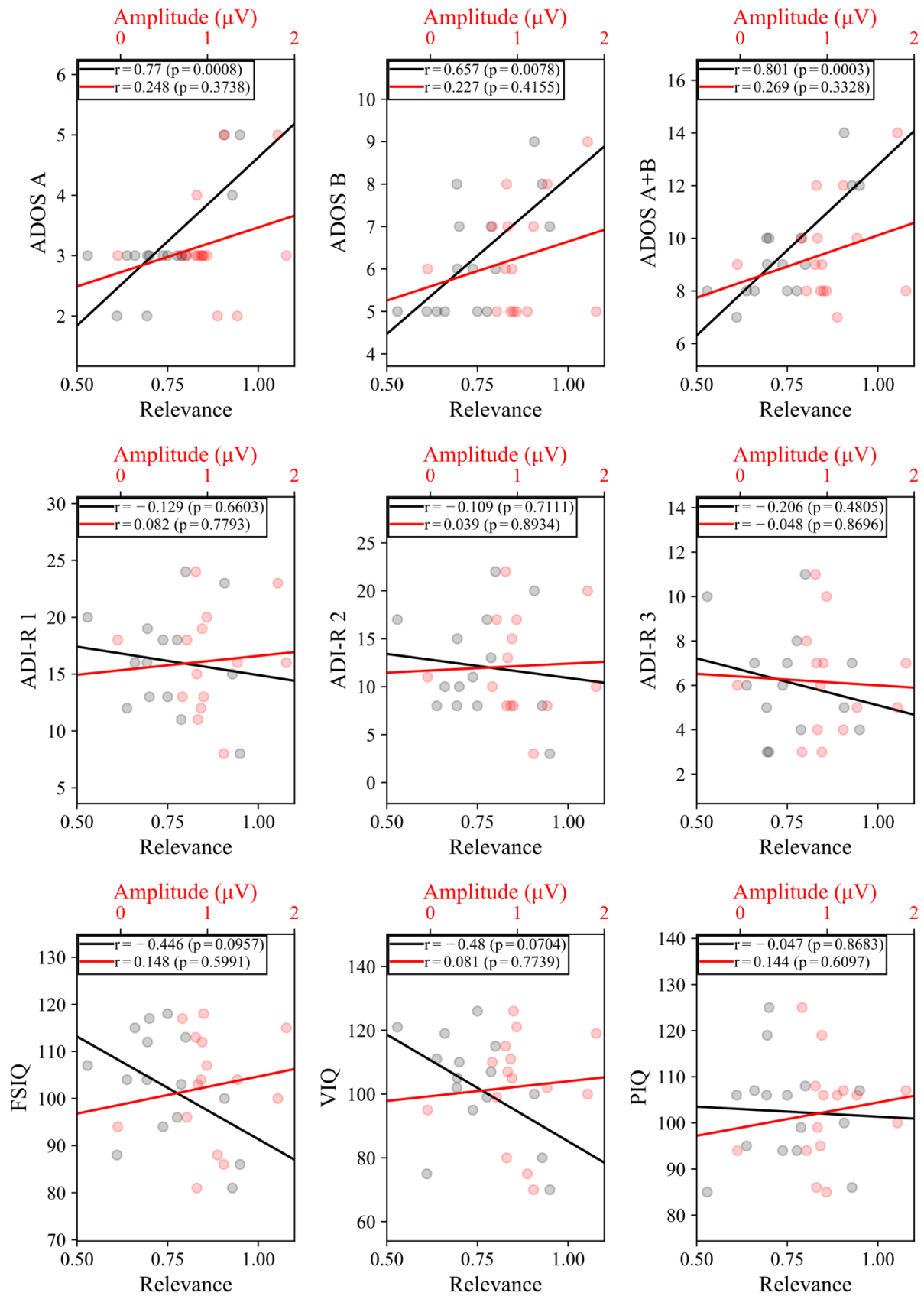


1
2 **Figure 6** – Results of transfer learning as a function of the number of calibration blocks of the
3 subject. Blue distributions show the accuracies obtained with the TL-W models while
4 transferring the knowledge from multiple subjects on a new subject (i.e., using the knowledge
5 embedded in a LOSO model to initialize the WS-W model on the held-out subject). Black
6 distributions show the accuracies obtained by randomly initializing the CNN (i.e., WS-W
7 models trained from scratch, without exploiting knowledge from other subjects). The
8 performance metric is reported for different numbers of calibration blocks used for training the
9 TL-W and WS models. Bigger dots represent the mean value across subjects, while the error
10 bars represent the standard error of the mean. Smaller dots represent the accuracy scored for
11 each subject (i.e., 15 data points). See part B of Table 3 for details about the number of
12 examples defining the training, validation, and test sets. Wilcoxon signed-rank tests were
13 performed to compare the performance of TL-W models and WS-W models. P-values were
14 corrected for multiple comparisons via the Benjamini–Hochberg procedure and significant
15 comparisons are marked once applied the correction (* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$). See
16 Section 2.8.1 for further details about the statistical tests.
17



1
2
3
4
5
6
7
8

Figure 7 – Spectral relevance (top) and spatial relevance (bottom). The spectral relevance is reported in its mean value (black line) and standard error of the mean (grey shaded area) across subjects. The more relevant frequency range (associated to a relevance ≥ 0.75) is delimited by the two dashed vertical lines. The spatial relevance is displayed for the two clusters (averaged within each cluster) and to the outliers detected. The number of subjects in each cluster is reported within brackets.



1
2 **Figure 8** – Correlation analysis between each ASD clinical score and the ICNN+ET-derived
3 (black) or ERP-derived measures (red). Subject-specific observations are reported with dots
4 together with regression lines. Pearson's correlation coefficients together with p-values are
5 reported inside each plot.

1 **TABLES**

2 **Table 1** – Measured ASD clinical scores, i.e., ADOS, ADI-R, and IQs (mean ± standard
3 deviation across subjects).

Measure	Value
<i>ADOS A: Communication</i>	3.20±0.86
<i>ADOS B: Social Interaction</i>	6.27±1.29
<i>ADOS A+B: Communication-Social Interaction</i>	9.47±1.86
<i>ADI-R 1: Social Interaction</i>	16.14±4.39
<i>ADI-R 2: Communication and Language</i>	12.14±5.19
<i>ADI-R 3: Restricted and Repetitive Behavior</i>	6.14±2.33
<i>FSIQ: Full Scale IQ</i>	102.53±11.24
<i>VIQ: Verbal IQ</i>	102.33±16.06
<i>PIQ: Performance IQ</i>	102.46±10.59

4

5

1 **Table 2** – Sinc-ShallowNet-v2. Each layer is provided with its name, main hyper-parameters
2 and number of trainable parameters. See Sections 2.2 and 2.3 for the meaning of the symbols.
3 In all layers, where not specified, stride (S) and padding (P) were set to (1,1) and (0,0),
4 respectively. The structural hyper-parameters (i.e., hyper-parameters of the architecture
5 structure) reported in bold were searched via Bayesian optimization.

Block	Layer name	Hyper-parameters	Number of trainable parameters
	Input	$K_0 = 1$	0
ISS	Sinc-Conv2D	$\mathbf{K}_0^{ISS}, F_0^{ISS} = (1,65),$ $P_0^{ISS} = (0, F_0^{ISS}[1]//2)$	$2 \cdot K_0^{ISS} \cdot K_0$
	BatchNorm2D		$2 \cdot K_0^{ISS}$
	Depthwise-Conv2D	$\mathbf{D}_1^{ISS}, K_1^{ISS} = K_0^{ISS} \cdot D_1^{ISS},$ $F_1^{ISS} = (C, 1), \mathbf{c}_1^{ISS}$	$F_1^{ISS}[0] \cdot F_1^{ISS}[1] \cdot K_1^{ISS}$
	BatchNorm2D		$2 \cdot K_1^{ISS}$
	ELU		0
	AvgPool2D	$F_p^{ISS} = S_p^{ISS} = (1,4)$	0
	Dropout	$\mathbf{p}^{ISS} = \mathbf{p}^{FST} = \mathbf{p}$	0
FST	Separable-Conv2D	$\mathbf{K}_0^{FST}, \mathbf{F}_0^{FST},$ $D_0^{FST} = 1, P_0^{FST} = (0, F_0^{FST}[1]//2)$	$F_0^{FST}[0] \cdot F_0^{FST}[1] \cdot K_1^{ISS} + K_1^{ISS} \cdot K_0^{FST}$
	BatchNorm2D		$2 \cdot K_0^{FST}$
	ELU		0
	AvgPool2D	$F_p^{FST} = S_p^{FST} = (1,8)$	0
	Dropout	$\mathbf{p}^{FST} = \mathbf{p}^{ISS} = \mathbf{p}$	0
FC	Flatten		0
	Fully-Connected	$N^{FC} = 2, \mathbf{c}_0^{FC}$	$N^{FC} \cdot (T//32 \cdot K_0^{FST} + 1)$
	Softmax		0

6
7

1 **Table 3** –Part A: The number of trials in the training set and validation set used to tune the
2 models (under Bayesian optimization) and number of trials in the test set used to test the
3 models’ performance, for the different training strategies. Part B: The number of trials in the
4 training set, validation set and test set used to tune and to test the WS-WS models and the TL-
5 WS models, to evaluate the beneficial effect of transfer learning. Note that in the computational
6 experiments of part B, the validation set was used only for early stopping, as the other hyper-
7 parameters were inherited from the LOSO models tuned in experiments of part A. The trials in
8 the test set were 2838 (on average) in each training strategy, as each model was tested
9 separately on each session-specific test set, and then the performance was averaged across all
10 sessions for each specific subject.
11

	Training strategy	No. of trials in the training set	No. of trials in the validation set	No. of trials in the test set
A	<i>Within-subject and within-session (WS-WS)</i>	1280	320	2838
	<i>Within-subject and cross-session (WS-CS)</i>	8960	2240	2838
	<i>Leave-one-subject-out (LOSO)</i>	125440	31360	2838
B	<i>Within-subject and within-session (WS-WS)</i>	From 128 to 1280 (step of 128)	From 32 to 320 (step of 32)	2838
	<i>Transfer learning on single session (TL-WS)</i>	From 128 to 1280 (step of 128)	From 32 to 320 (step of 32)	2838

12
13

1 **Table 4** – Searched hyper-parameters of Sinc-ShallowNet-v2: distributions and admitted
 2 values sampled during Bayesian optimization. Curly brackets denote discrete admitted values,
 3 while square brackets denote interval of admitted values.

Hyper-parameter	Distribution	Values
K_0^{ISS}	uniform	{4, 8, 16}
D_1^{ISS}	uniform	{1, 2, 4}
c_1^{ISS}	uniform	{None, 0.25, 0.5, 0.75, 1, 1.25, 1.5}
$p^{ISS} = p^{FST} = p$	uniform	{None, 0.25, 0.5}
K_0^{FST}	uniform	$\{K_0^{ISS} \cdot D_1^{ISS}, 1, 2\}$
F_0^{FST}	uniform	{(1,5), (1,9), (1,13), (1,17), (1,21)}
c_0^{FC}	uniform	{None, 0.25, 0.5, 0.75, 1, 1.25, 1.5}
lr	log-uniform	[1e-4, 1e-1]

4
5

1 **Table 5** – Model size (expressed as the number of trainable parameters) and training time
 2 (expressed in s/epoch) of Sinc-ShallowNet-v2 for each training condition. The total number of
 3 trainable parameters (mean \pm standard error of the mean) is reported together with the number
 4 of parameters specific for each block, indicating within brackets the percentage of parameters
 5 exploited in each block.

	Block	Within-subject and within-session (WS-WS)	Within-subject and cross-session (WS-CS)	Leave-one- subject-out (LOSO)
<i>Model size (# tr. parameters)</i>	-	1207 \pm 141	1655 \pm 412	4638 \pm 559
	<i>ISS</i>	336 \pm 21 (28%)	466 \pm 47 (28%)	555 \pm 43 (12%)
	<i>FST</i>	749 \pm 121 (62%)	1042 \pm 368 (63%)	3689 \pm 483 (80%)
	<i>FC</i>	122 \pm 13 (10%)	147 \pm 40 (9%)	394 \pm 34 (8%)
<i>Training time (s/epoch)</i>	-	0.980 \pm 0.014	6.380 \pm 0.319	23.1 \pm 1.1

6