

Interpretable SincNet-based Deep Learning for Emotion Recognition from EEG brain activity

Juan Manuel Mayor-Torres¹, Mirco Ravanelli², Sara E. Medina-DeVilliers³, Matthew D. Lerner³,
and Giuseppe Riccardi¹

Abstract—Machine learning methods, such as deep learning, show promising results in the medical domain. However, the lack of interpretability of these algorithms may hinder their applicability to medical decision support systems. This paper studies an interpretable deep learning technique, called SincNet. SincNet is a convolutional neural network that efficiently learns customized band-pass filters through trainable sinc-functions. In this study, we use SincNet to analyze the neural activity of individuals with Autism Spectrum Disorder (ASD), who experience characteristic differences in neural oscillatory activity. In particular, we propose a novel SincNet-based neural network for detecting emotions in ASD patients using EEG signals. The learned filters can be easily inspected to detect which part of the EEG spectrum is used for predicting emotions. We found that our system automatically learns the high- α (9-13 Hz) and β (13-30 Hz) band suppression often present in individuals with ASD. This result is consistent with recent neuroscience studies on emotion recognition, which found an association between these band suppressions and the behavioral deficits observed in individuals with ASD. The improved interpretability of SincNet is achieved without sacrificing performance in emotion recognition.

I. INTRODUCTION

The recent development of artificial intelligence fosters unprecedented innovation in the healthcare domain. The availability of big data repositories, the development of robust learning algorithms, and the availability of appropriate computational resources are making technologies such as deep learning, applicable and challenging the state-of-the-art systems. Deep learning (DL) has indeed been used in many medical applications [1], including disease diagnosis [2], and personalized medicine [3], just to name a few. In such critical cases, the lack of interpretability of this technology limits its widespread end-user adoption and may even lead to adverse consequences [4].

Current deep neural networks map low-level data into higher-level concepts using a pipeline of non-linear transformations [5]. Therefore, the final prediction is normally not explained by the network, and end-users (i.e. healthcare professionals) do not know if the outcome is based on solid evidence or due to some statistical biases [6]. Moreover, the inspection of the intermediate representations learned by the network rarely helps to explain the neural predictions. Improving the

interpretability of the current technology is an essential condition for overcoming understandable skepticism, reluctance, and hesitations of the medical personnel.

Interpretable deep learning has been the object of increasing research efforts over the last years [7]. Historically, there was a trade-off between performance and interpretability. Simple models like linear regression are transparent but not competitive in terms of performance with deeper models such as fully-connected, convolutional, and recurrent neural networks. A possible solution may be post-hoc interpretability [7], in which, a complex neural model is built and analyzed afterward. The interpretability can also be achieved with surrogate models (e.g., locally interpretable model explanations - LIME [8]) and with gradient-based methods like in the saliency maps [9]. In contrast, the alternative is to train machine learning models that are interpretable by design. Along this line, a novel model called SincNet has been proposed [10]. SincNet is a convolutional neural network that learns a custom filter-bank using sinc-based convolutional kernels. Following the network training, the filters can be inspected in the frequency domain to identify which parts of the spectrum are used by the neural network to perform a prediction. The filter inspection is straightforward and often insightful, as emerged in some recent studies [11]–[13]. The interpretability insights are directly mapping the neural network parameters to the input signal. SincNet has been originally proposed for processing audio sequences [10], but it has been recently used for EEG-based brain signals as well. In particular, SincNet has been successfully adopted for EEG-based motor-imagery tasks [12], [13]. The full potential of this model in high-level processing of input stimuli (e.g. visual or audio), however, is yet to be explored.

In this paper, we propose SincNet for studying the brain activity of patients with Autism Spectrum Disorder (ASD). In particular, we analyze EEG signals of ASD and non-ASD individuals while performing a Facial Emotion Recognition (FER) task. The EEG recordings feed a deep learning model based on SincNet, which tries to guess patients' emotions from their brain waves. From the study of filters learned by SincNet, we found an interesting consistency with the Power-Spectral-Density analysis of previous EEG studies on ASD. The improved interpretability is achieved without sacrificing the performance of the machine learning models. The proposed system provides a contribution to the interpretability-by-design alternative by applying deep learning techniques to the medical domain.

¹University of Trento, Department of Information Engineering and Computer Science (DISI), Via Sommarive 5, Povo, Trento 38123, ²Mila - Quebec Artificial Intelligence Institute, Montreal, QC, Canada, ³StonyBrook University, NY, USA, Department of Psychology corresponding author: juan.mayortorres@unitn.it

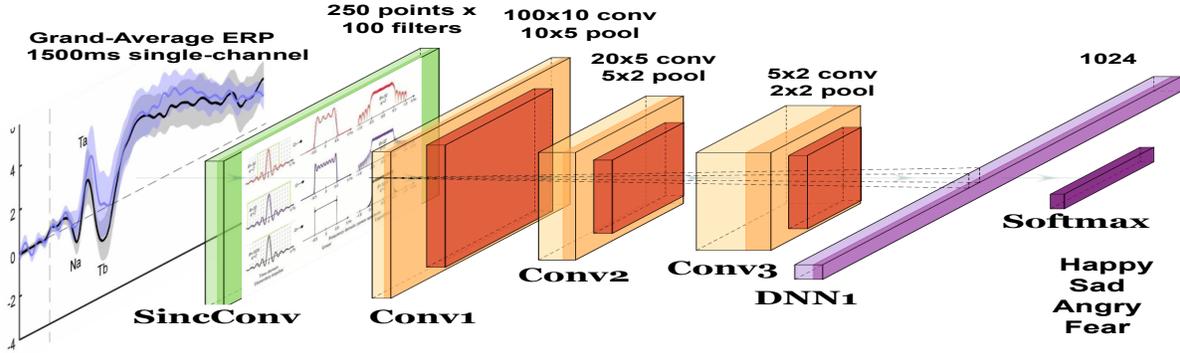


Fig. 1: The proposed architecture for EEG-based Face Emotion Recognition. The neural pipeline is composed of a *SincConv* layer, three standard conv-pool blocks (*Conv1*, *Conv2*, and *Conv3*), and a fully-connected network (*DNN1*) connected to a softmax classifier.

II. FACIAL EMOTION RECOGNITION WITH SINCNET

The Facial Emotion Recognition (FER) task measures the ability to identify basic emotions in facial expressions. It is particularly important to assess this ability in individuals with Autism Spectrum Disorder (ASD) because they frequently demonstrate impairments in the ability to accurately categorize and label the emotional facial expressions of others. In this section, we first describe the FER task. Then, we describe the proposed architecture for EEG-based emotion recognition in ASD and typically developing (TD) individuals.

A. FER Task Description

Eighty-eight 14-to-17-year old adolescent participants (48 without ASD, while 40 with ASD) completed a FER task. The FER task, called Diagnostic Analysis of Nonverbal Behavior (DANVA-2) [14], consists of presenting 48 photographs in random order. Each photograph represents one of following emotions: *happy*, *sad*, *angry*, and *fear*. The picture was first shown to the users for 2 second and, within this interval, their EEG signal was recorded. Subsequently, each photograph was shown again (for up 4 seconds) to give the participant enough time to select the emotion they have seen in the picture using a button box. The EEG activity of each participant was fed into the proposed machine learning system, which is designed to predict the actual emotion from the given brain recordings. The full pipeline for EEG-based emotion recognition is shown in Figure 1.

B. EEG pre-processing

The EEG signals were captured with a 32-channel ActiCHAMP device from Brain Products and digitized with a sampling rate of 500 Hz (16-bit resolution). The data were recorded continuously using the BrainVision Recorder software and processed using the BrainVision Analyzer 2.0 for offline data reduction. We converted all the input channels averaging them into a single sequence or a grand-average Event- Related Potential (ERP) representation for each participant and for each emotion. This helps reducing noise and statistical variability of EEG trials [15]. We employed

a standard EEG pre-processing pipeline based on band-pass filtering, amplitude normalization, and Zero Component Analysis (ZCA) whitening [16]. We also apply the ADJUST algorithm [17] and the PREP pipeline [18] to clean the EEG signal before feeding the neural network.

C. Interpretable SincNet

The resulting grand-average ERP representation feeds the proposed SincNet-based system. SincNet is a convolutional neural network whose first layer, called *SincConv*, is designed to learn tunable Finite Impulse Response (FIR) filters. In standard CNNs, all the elements (taps) of each filter are learned from the data. In SincConv, instead, we parametrize the kernels in order to implement rectangular bandpass filters, whose cut-off frequencies are the only two parameters learned from data. This can be achieved with the following parametrization:

$$y[n] = x[n] * g[n, f_1, f_2], \quad (1)$$

where $x[n]$ is the input EEG signal and $y[n]$ is the output of the SincConv layer. The convolution is performed with the kernel g , which is defined in this way:

$$g[n, f_1, f_2] = 2f_2 \text{sinc}(2\pi f_2 n) - 2f_1 \text{sinc}(2\pi f_1 n), \quad (2)$$

where the frequencies f_1 and f_2 are the learned parameters. This technique not only saves a lot of parameters but naturally leads to a more interpretable model. Moreover, the filters depend on human-readable parameters with a clear physical meaning. At the end of the training, the filters are inspected to identify which parts of the spectrum are covered by filters. This helps users better understand what the network has learned [11], [13]. In addition to the SincConv layer, the proposed architecture employs three 2-D convolutional blocks based on standard convolution, batch normalization, pooling, ReLU activations, and dropout. Finally, we plug a fully connected layer followed by a softmax classifier.

III. EXPERIMENTAL SETUP

In this section, we describe our SincNet-based emotion recognition architecture and its training procedure.

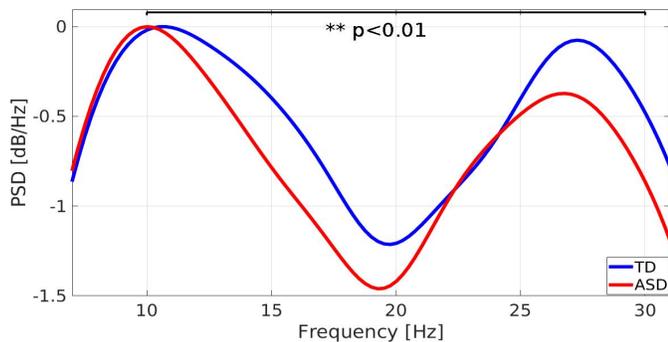


Fig. 2: Power Spectral Density (PSD) learned by SincNet filters on TD and ASD individuals. Significant differences are observed in high- α (9-13 Hz) and β (13-30 Hz) bands.

A. Architecture details

The adopted SincConv layer employed 100 filters with a kernel size of 250 points (which represent the impulse response of the filter). As in the original paper, [10], a Hamming window was used to mitigate ripples in the filters. The 2-D convolutional blocks were based on 32, 64, and 128 channels with kernel sizes of (100 x 10), (20 x 5), (5 x 2), respectively. Max-pooling used kernel sizes of (10 x 5), (5 x 2), (2 x 2). Batch normalization was added between the convolution and the ReLU activations. Next, we employed a single fully connected layer composed of 1024 ReLU neurons. Dropout was used in both convolutional and fully connected layers with a rate of 0.5. The final prediction over the four emotions was performed with a softmax classifier.

B. Network Training

We trained and evaluated the SincNet-based pipeline using a Leave-One-Trial-Out (LOTO) approach [16]. Thus, for each participant we used 47 out of the 48 trials for training, leaving a different test trial out every time. Therefore, we completed a total of 48 training/validation experiments. This modality was needed due to the lack of in-domain data for this specific task. The neural network was initialized with the standard Glorot’s initialization scheme [19]. We used categorical cross-entropy as a loss function. The gradient was computed with the backpropagation algorithm, while parameters are updated using the Adam optimizer [20]. We used a learning rate of 0.001 with a weight-decay penalty of $1e - 05$. We trained the neural network for 400 epochs with a batch size of 30. For more information on the neural architecture and training modality, please refer to the open-source code repository of this project¹.

IV. RESULTS

In the following section, we report the experimental evidence that emerged from the FER task using SincNet.

¹<https://github.com/meiyor/SincNet-for-Autism-EEG-based-Emotion-Recognition>

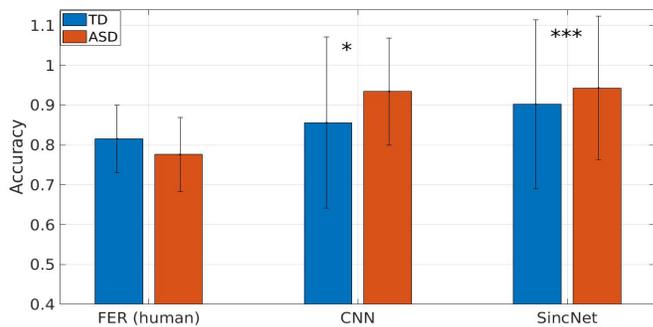


Fig. 3: Accuracy comparison between FER (humans), CNN, and SincNet. Human accuracy is significantly lower than the one of deep learning systems.

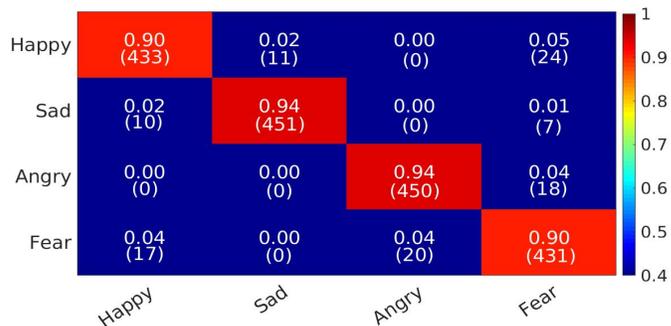


Fig. 4: Confusion Matrix for the SincNet pipeline on ASD.

A. Filter Analysis

After training, we inspected the filters learned by SincNet for ASD and non-ASD (TD) participants. To analyze which EEG frequency bands were used for the emotion prediction, we performed a Fourier Transform of the learned filters. We then averaged their frequency responses and computed the cumulative Power Spectral Density (PSD) reported in Fig. 2. We observed significant differences in the filters learned from ASD and non-ASD participants. In particular, an attenuation in the high- α (9-13 Hz) and β (13-30 Hz) emerged in the cumulative PSD spectrum of ASD participants. We observed significant differences between non-ASD and ASD groups on high- α (9-13 Hz) $F(1, 87) = 3.331, p = 0.000267$ and β (13-30 Hz) $F(1, 87) = 2.05, p = 0.00102$ bands after Bonferroni-Holm correction. Our findings are consistent with previous studies on ASD, indicating that high- α and β attenuations are related to FER deficits on individuals with ASD [21], [22]. Although this frequency-band attenuation has been observed in some EEG studies including individuals with ASD, the causes of this are not fully understood in the scientific community. Some authors hypothesized that this is associated with multiple behavioral deficits of individuals with ASD [23], [24]. However, it is worth notice that SincNet learns that these bands are not useful to predict emotions in individuals with ASD. Notably, these predictions were learned automatically from raw EEG data only, without providing any additional information to the network.

B. Performance Analysis

Figure 3 compares the accuracy achieved by humans in FER with the one reached by the deep learning systems based on CNNs and SincNet. As for the human performance or FER, we found a small difference between ASD and non-ASD accuracies (79% vs 81%, respectively), thus, showing some consistent deficits in individuals with ASD performing FER tasks [23], [24]. Interestingly, deep learning systems outperformed the FER human accuracy. This suggests that there were some cases where a participant wrongly labeled the photograph, but the SincNet-based system was able to detect the correct emotion from participants' EEG brain activity. This difference is more evident in ASD participants, where the SincNet improvement is greater (79% vs 92%) than FER. The SincNet confusion matrix for the ASD group is shown in Figure 4. The proposed SincNet pipeline turned out to slightly outperforming a CNN-based system (90% vs 85% for the non-ASD group and 92% vs 91% for the ASD group). The CNN architecture was obtained by replacing the SincConv layer with a standard convolution. This improvement is consistent with what was observed for audio [10] and motor-based EEG signal in previous studies [13]. In sum, our results confirm that SincNet improves the interpretability of the model without sacrificing performance.

V. CONCLUSIONS

This paper has proposed the application of an interpretable deep-learning architecture, SincNet, in a medical domain. We applied this model to study EEG activity patterns of ASDs and non-ASDs in a FER task. Our results indicate that SincNet transparently learns the high- α and β suppressions observed in ASD individuals when perceiving and recognizing emotional faces. SincNet improves the interpretability of the neural model without affecting its performance, thus offering a convenient way to avoid the performance versus interpretability dilemma.

ACKNOWLEDGMENT

The authors would like to thank the University of Trento High-Performance Computing and Stony Brook Research Computing and Cyber-infrastructure which was made possible by a \$1.4M National Science Foundation grant (1531492). This research was supported by NIMH grant R01MH110585, grants from the AAF Fund for Communication.

REFERENCES

- [1] Juan M Mayor Torres, Tessa Clarkson, Evgeny A Stepanov, Christian C Luhmann, Matthew D Lerner, and Giuseppe Riccardi. Enhanced error decoding from error-related potentials using convolutional neural networks. In *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, pages 360–363. IEEE, 2018.
- [2] Shu Lih Oh, Yuki Hagiwara, U. Raghavendra, Yuvaraj Rajamanickam, N. Arunkumar, M. Murugappan, and U. Acharya. A deep learning approach for parkinson's disease diagnosis from eeg signals. *Neural Computing and Applications*, pages 1–7, 2018.
- [3] G. Z. Papadakis, A. Karantanas, Tsiknakis H, Tsatsakis M, Spandidos A, and Marias D. A. Deep learning opens new horizons in personalized medicine (review). *Biomedical Reports*, 10(4):215–217, 2019.
- [4] Muhammad Aurangzeb Ahmad, Carly Eckert, and Ankur Teredesai. Interpretable machine learning in healthcare. In *Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*, page 559–560, 2018.
- [5] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- [6] Sebastian Lapuschkin, Stephan Wäldchen, Alexander Binder, Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. Unmasking clever hans predictors and assessing what machines really learn. *Nature communications*, 10(1):1–8, 2019.
- [7] Christoph Molnar. *Interpretable Machine Learning*. 2019. <https://christophm.github.io/interpretable-ml-book/>.
- [8] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should i trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 1135–1144, 2016.
- [9] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualizing image classification models and saliency maps. In *Proc. of ICLR*, 2014.
- [10] Mirco Ravanelli and Yoshua Bengio. Speaker recognition from raw waveform with sincnet. In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pages 1021–1028. IEEE, 2018.
- [11] Mirco Ravanelli and Yoshua Bengio. Interpretable convolutional filters with sincnet. *arXiv preprint arXiv:1811.09725*, 2018.
- [12] Hong Zeng, Zhenhua Wu, Jiaming Zhang, Chen Yang, Hua Zhang, Guojun Dai, and Wanzeng Kong. Eeg emotion classification using an improved sincnet-based deep learning model. *Brain sciences*, 9(11):326, 2019.
- [13] Davide Borra, Silvia Fantozzi, and Elisa Magosso. Interpretable and lightweight convolutional neural network for eeg decoding: application to movement execution and imagination. *Neural Networks*, 129:55–74, 2020.
- [14] Stephen Nowicki. Manual for the receptive tests of the diagnostic analysis of nonverbal accuracy 2. *Atlanta, GA: Department of Psychology, Emory University*, 2000.
- [15] Arnaud Delorme, Makoto Miyakoshi, Tzyy-Ping Jung, and Scott Makeig. Grand average erp-image plotting and statistics: A method for comparing variability in event-related single-trial eeg activities across subjects and conditions. *Journal of neuroscience methods*, 250:3–6, 2015.
- [16] Juan Manuel Mayor Torres, Tessa Clarkson, Kathryn M Hauschild, Christian C Luhmann, Matthew D Lerner, and Giuseppe Riccardi. Facial emotions are accurately encoded in the neural signal of those with autism spectrum disorder: A deep learning approach. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 2021.
- [17] Andrea Mognon, Jorge Jovicich, Lorenzo Bruzzone, and Marco Buiatti. Adjust: An automatic eeg artifact detector based on the joint use of spatial and temporal features. *Psychophysiology*, 48(2):229–240, 2011.
- [18] Nima Bigdely-Shamlo, Tim Mullen, Christian Kothe, Kyung-Min Su, and Kay A Robbins. The prep pipeline: standardized preprocessing for large-scale eeg analysis. *Frontiers in neuroinformatics*, 9:16, 2015.
- [19] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.
- [20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [21] Jaime Pineda, Ashley Juavinett, and Michael Datko. Self-regulation of brain oscillations as a treatment for aberrant brain connections in children with autism. *Medical hypotheses*, 79(6):790–798, 2012.
- [22] Elisabeth VC Friedrich, Aparajithan Sivanathan, Theodore Lim, Neil Suttie, Sandy Louchart, Steven Pillen, and Jaime A Pineda. An effective neurofeedback intervention to improve social interactions in children with autism spectrum disorder. *Journal of autism and developmental disorders*, 45(12):4084–4100, 2015.
- [23] Geraldine Dawson, Sara Jane Webb, and James McPartland. Understanding the nature of face processing impairment in autism: insights from behavioral and electrophysiological studies. *Developmental neuropsychology*, 27(3):403–424, 2005.
- [24] Melissa H Black, Nigel TM Chen, Kartik K Iyer, Ottmar V Lipp, Sven Bölte, Marita Falkmer, Tele Tan, and Sonya Girdler. Mechanisms of facial emotion recognition in autism spectrum disorders: Insights from eye tracking and electroencephalography. *Neuroscience & Biobehavioral Reviews*, 80:488–515, 2017.