

# INVESTIGATION ON NEURAL RESPONSES RELATED TO THE LOCALIZATION OF NATURAL SOUNDS

Andrea Bellotti<sup>1,2</sup>

Marco Binelli<sup>2</sup>

Giovanni M. Di Liberto<sup>3,4</sup>

Luca Ascari<sup>3</sup>

Christopher Raymaekers<sup>1</sup>

Angelo Farina<sup>2</sup>

<sup>1</sup> Toyota Motor Europe, Zaventem, Belgium

<sup>2</sup> Department of Engineering and Architecture, University of Parma, Parma, Italy

<sup>3</sup> Camlin Italy, Parma, Italy

<sup>4</sup> École normale supérieure, PSL University, CNRS, Paris, France

andrea.bellotti@toyota-europe.com

## ABSTRACT

Spatial hearing allows the localization of sounds in complex acoustic environments. There is considerable evidence that this neural system rapidly adapts to changes in sensory inputs and behavioral goals. However, the mechanisms underlying this context-dependent coding are not well understood. In fact, previous studies on sound localization have mainly focused on the perception of simple artificial sounds, such as white-noise or pure tone bursts. In addition, previous research has generally investigated the localization of sounds in the frontal hemisphere while ignoring rear sources. However, their localization is evolutionary relevant and may show different neural coding, given the inherent lack of visual information. Here we present a pilot electroencephalography (EEG) study to identify robust indices of sound localization from participants listening to a short natural sound from eight source positions on the horizontal plane. We discuss a procedure to perform a within-subject classification of the perceived sound direction. Preliminary results suggest a pool of discriminative subject-specific temporal and topographical features correlated with the characteristics of the acoustic event. Our preliminary analysis has identified temporal and topographical features that are sensitive to spatial localization, leading to significant decoding of sounds direction for individual subjects. This pilot study adds to the literature a methodological approach that will lead to the objective classification of natural sounds location from EEG responses.

## 1. INTRODUCTION

Auditory processing in the human auditory cortex has been suggested to be underpinned by a dual neural system [1–6], with anterior areas largely engaged in decoding the content of a sound ('what'), whereas posterior temporal and parietal areas having a crucial role in the processing of spatial information ('where').

Previous research showed distinct cortical patterns when listeners were presented with sounds from various directions, demonstrating that non-invasive neural recordings such as electroencephalography (EEG) are sensitive to au-

ditory spatial processing [7, 8]. However, it is still unclear how accurately that spatial auditory signal can be decoded. The present work investigated the cortical processing of spatial auditory perception and assessed the possibility of decoding sound location from the EEG signals. Previous work on spatial hearing focused on artificial sounds, such as white-noise or pure tone bursts, even though it has been shown that, in some cases, natural sounds produce richer and stronger responses [9]. We adopted a wood cracking sound, which is a familiar natural sound with quasi-impulsive characteristics, and thus we hypothesized that it would produce more complex patterns, including some clear Evoked Potentials (EP; [10]), allowing stronger classification performances. The combination of a natural sound, real speakers, and a pure listening task (rather than a decision task; e.g. P300 paradigm) was used to evaluate the classification performance in a simulated realistic scenario.

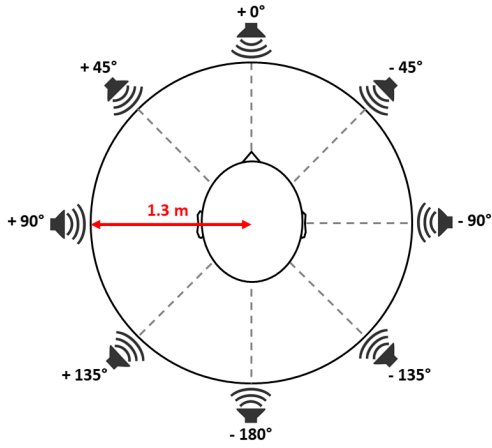
## 2. MATERIALS AND METHODS

### 2.1 Participants

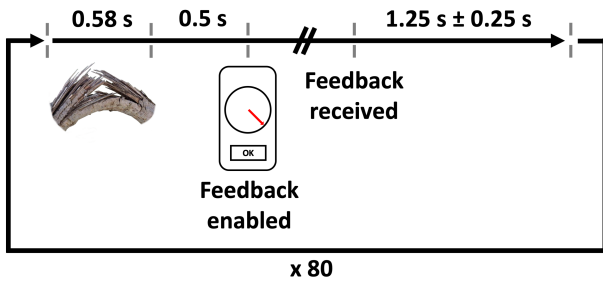
Five healthy volunteers (1 female, one left-handed, mean age 32.2 years ranging between 29 and 39 years) participated in the study. All participants provided voluntary information consent and reported to have normal hearing abilities and no known neurological or psychiatric diseases. One subject was excluded from the analysis as the experiment could not be completed.

### 2.2 Data acquisition

EEG measurements were recorded at a sampling rate of 500 Hz using a g.Nautilus PRO system (gtec, Austria) equipped with 32 active dry electrodes (g.SAHARA) positioned according to the 10-20 system. Reference and ground were placed at the two mastoids. The acquisitions were performed in four different days in the dimly lit, acoustically treated, listening room of the electroacoustic laboratory of *Casa del Suono* (Parma, Italy).



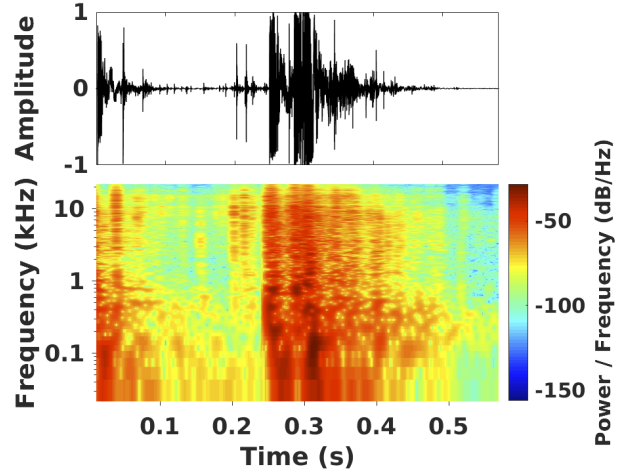
**Figure 1:** Experimental setup. The natural sound was played from 8 evenly separated loudspeakers positioned on the horizontal plane at a distance of 1.3 meters from the sweet spot. The listening level at the sweet spot was adjusted to peak values of 80 dBA.



**Figure 2:** Experimental paradigm. Each recording session was composed of two runs consisting of 80 trials. Each trial consisted of the presentation of the sound from a given direction. The randomized sequence of 80 directions within a run was balanced (same number of presentations from each direction). A subjective feedback was provided after each trial through a mobile application and then a random inter-trial time interval was implemented to avoid inter-trial phase-locking effects.

### 2.3 Experimental setting

Subjects sat in a comfortable chair surrounded by eight loudspeakers evenly spaced every 45° (+0°, -45°, -90°, -180°, +135°, +90° and +45°) at ear height and 1.3 meters from the sweet spot as shown in Fig.3. The loudspeakers were controlled by Max 8 (Cycling '74) hosted on a Windows 7 computer and driven through a dedicated soundcard (details about room and setup can be found in [11]). The speakers were equalized using inverse filters of the impulse response computed applying the Kirkeby regularization [12]. Moreover, as suggested in [13], the listening level at the sweet spot was calibrated to peak values of 80 dBA, which corresponds to the subjective preferred level. The experiment consisted in a pure sound localization task where the subjects were asked to confirm the perceived direction of a natural sound randomly played from one of



**Figure 3:** Characteristics of the adopted natural sound. Top plot shows the raw signal (mono track sampled at 44100 Hz), and the bottom one represents its spectrogram (time vs frequency).

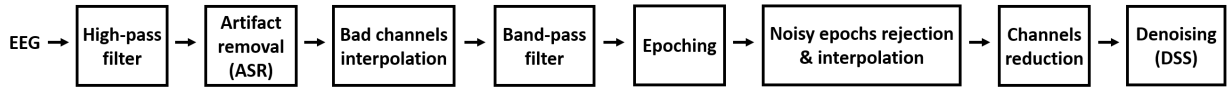
the eight directions. The sound used in the study (shown in gray in Fig.5) was a 0.58 seconds long wood cracking sound exhibiting two main impulsive components. The participants provided their feedback after each sound through a custom mobile application running on a smartphone that they held on their hands throughout the whole experiment. Each subject participated in four sessions that were performed in four different days. A session was composed of two runs consisting of 80 trials each (10 repetitions per direction), interleaved by a 5-minute break.

The sound directions of each run were defined by a different random sequence and the timings of each trial followed the protocol shown in Fig.2. After the end of the sound reproduction, the interaction with the mobile application was disabled for 0.5 seconds to avoid undesired movements close to the time window of interest. Once enabled, the participants could confirm their perceived direction without timing constraints to allow them to make the best choice. The inter-trial time interval between the feedback confirmation and the beginning of the following sound reproduction randomly varied between 1 second and 1.5 seconds. The randomization of the inter-trial interval was implemented to avoid artificially induced phase-locking in the epochs.

### 2.4 Pre-processing

Pre-processing was performed offline by using EEGLAB (version 14.1.1b) and custom Python code. At first, EEG recordings were visually inspected to reject noisy channels and time intervals. As a result, channel Oz was excluded from the analysis due to the abnormal fluctuations present in multiple subjects' sessions.

After this preliminary step, the pre-processing pipeline shown in Fig.4 was implemented. Raw EEG signals were pre-processed with the Artifact Subspace Reconstruction (ASR) algorithm to increase the SNR [14, 15]. Specifically, signals were high-pass filtered using the default filter



**Figure 4:** Pre-processing pipeline.

(non-causal FIR filter, Kaiser window, 0.25 Hz - 0.75 Hz transition band, 80 dB attenuation) and the clean calibration data was automatically extracted from the recordings. The calibration data is used in the algorithm to estimate the channels covariance used both to identify and to interpolate noisy intervals. All parameters of the ASR method were set to the default values, with the exception of a more stringent channel correlation criterion (minimum channel correlation allowed equal to 0.75 compared to the default of 0.85), a less constraining subject-specific standard deviation threshold ranging from 7 to 12 (compared to the default of 5) to remove only particularly large artifacts, and a more aggressive criterion for bad channels identification (maximum tolerated fraction with respect to the total recording duration set to 0.2 compared to the default of 0.5). The channels identified as bad by ASR were replaced by a spherical interpolation of all the remaining channels using the FieldTrip standard 10-5 channel locations.

After the artifact removal phase, cleaned signals were band-pass filtered (one-pass, zero-phase, non-causal FIR filter, Hamming window with 0.0194 pass-band ripple and 53 dB stop-band attenuation, 3 Hz - 35 Hz) and then epoched retaining the 0.58 after the acoustic stimulus onset. The single channel signals exceeding the threshold range  $\pm 45 \mu\text{V}$  were spherically interpolated. If, after this interpolation step, any of the single channel signals of a trial were still exceeding the threshold, the trial was rejected. The percentage of epochs exceeding the threshold for the four subjects was  $4.7\% \pm 4.3\%$ , and the single channel signals interpolation successfully corrected  $91.6\% \pm 13.9\%$  of them.

Finally, the Denoising Source Separation (DSS) approach [16] was used to denoise the EEG data. This method rotates the data into a component space that maximizes the separability among the classes. As a result, it facilitates the subsequent data analysis by reducing the within-class and increasing the between-class variability. Being this method strongly dependent on the average of all the channels, prior to its usage we discarded the high-amplitude frontal channels Fp1 and Fp2. The optimal number of DSS components to retain was selected through cross-validation in the data analysis phase for each subject and every spatial configuration. In most of the cases the best value found was in the range between 20 and 25 components out of 29, thus rejecting from 4 to 9 components.

## 2.5 Data analysis

After pre-processing, the extracted trials were used to train subject-specific classification models for different localization configurations (shown in Table 1, 2, and 3). In particular, given the limited amount of trials per class at disposal, we used the Random Forest ensemble method leveraging

bagging to reduce model variance. Moreover, we experimented with localization configurations involving aggregations of directions (see first three columns of Table 3) in order to increase the number of trials per class hypothesizing common patterns in the brain responses associated to spatially related directions. We tuned the hyperparameters of the classifiers through a grid search in a 5-fold cross-validation configuration with a 70% train, 20% validation and 10% test splitting.

EEG epochs associated to the different directions and time-locked to the stimulus onset were aggregated using the median operator to obtain the corresponding event-related potentials (ERPs) [10]. All further analyses were conducted on 500 ms epochs that started before the arising of the N1 component in the ERP. This window of interest started at 80 ms, as shown by the black dashed vertical lines in Fig. 5. The reason for this choice is that the N1 component showed a delay (latency of around 132 ms after stimulus onset) compared to the typical N1 latency of a sound onset response (80-120 ms; [17, 18]). Future experiments will tackle this issue by recording the played sounds with a microphone to extract more precise timing information.

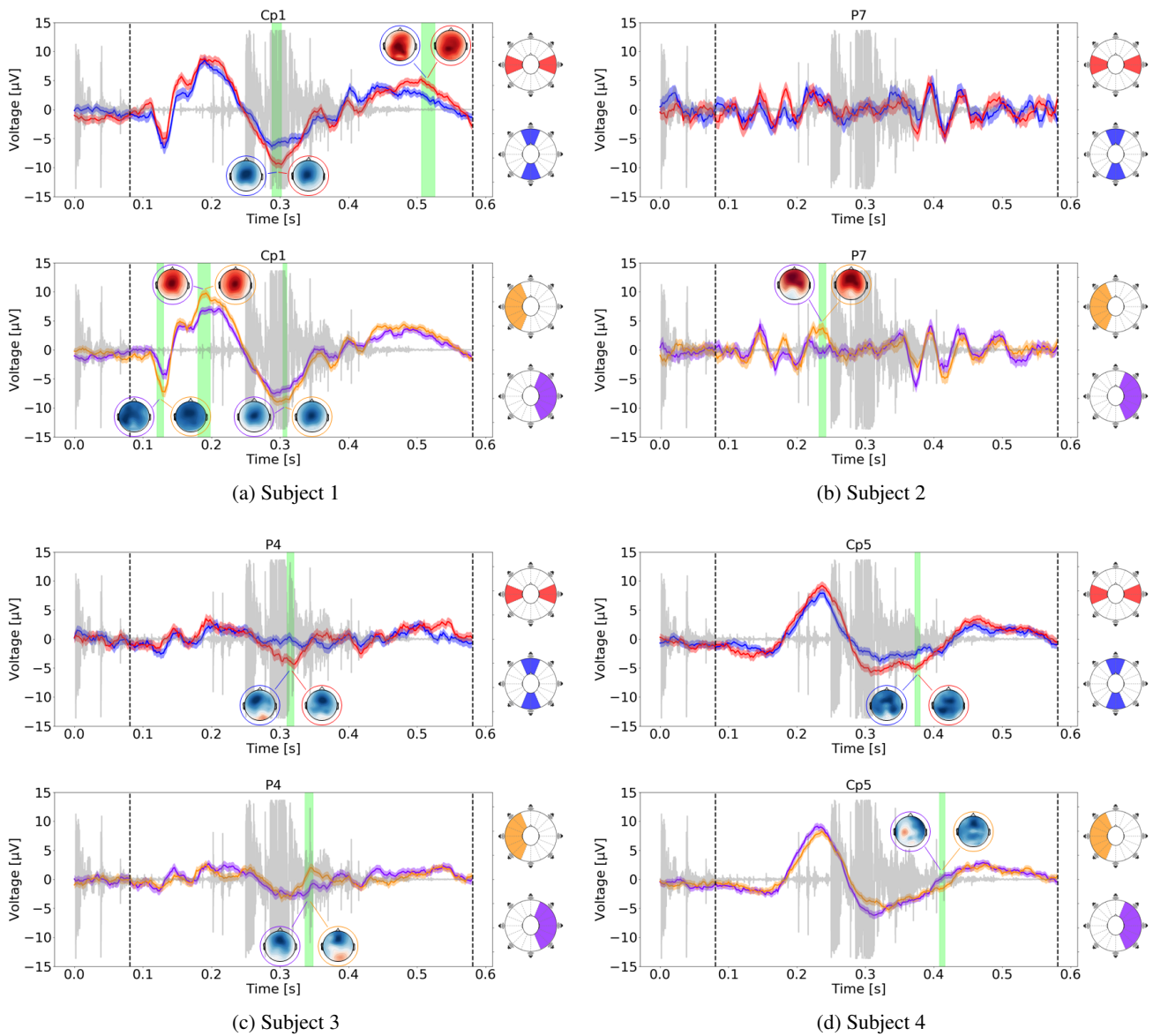
To build the input feature vector, for each channel we partitioned this interval and we used the average value within the sub-windows as features. The size of the sub-windows was selected through a cross-validation procedure and in most cases was found to be 3 samples, resulting in a feature vector length of 2407 (29 channels  $\times$  83 sub-windows). With more data, we would automatically select only the statistically significant intervals like those highlighted in green in Fig. 5. We didn't apply this procedure because it requires a greater dataset in order to obtain reliable outcomes about the significance, especially considering the need for cross-validation. In this pilot, in order to be as fair as possible, we rather decided to blindly use all possible features, leaving to the classifier the burden of identifying the best ones in a supervised fashion.

The other hyperparameters tuned through cross-validation were number and depth of the trees, set to 1000 and 12 respectively and, as discussed in Section 2.4, the number of components to retain in the DSS denoising phase. This parameter was tuned separately for each of the classification analyses (i.e. for each subject and localization configuration). The optimal value was generally between 20 and 25 components out of 29, with a modal value of 21 (value used in the plots of Fig. 5).

## 3. RESULTS

### 3.1 Physiological responses

Fig. 5 compares for each subject the ERP response of a representative channel for two localization configura-



**Figure 5:** ERP response to the sound onset of the four subjects in the two binary localization configurations providing the greatest discrimination accuracy (configuration S, longitudinal vs latitudinal, top figure, and configuration R, left area vs right area, bottom figure). For each subject, only one representative channel is presented. The gray line in the background represents the reproduced natural sound and the green areas show the statistically significant time windows ( $p < 0.05$ , Welch's t-test with Bonferroni correction) within the area of interest delimited by the black dashed vertical lines. Scalp topographies for the two classes are presented to show the spatial distribution in the statistically significant time windows of each localization configuration.

	A)	B)	C)	D)	E)	F)	G)	H)
Subject 1	61.2 ± 11.8	60.9 ± 9.6	56.9 ± 12.0	62.5 ± 12.0	58.4 ± 15.1	56.9 ± 9.9	59.4 ± 16.9	63.1 ± 12.2
Subject 2	57.8 ± 9.4	61.9 ± 13.2	53.1 ± 13.3	57.5 ± 11.8	41.2 ± 13.5	47.2 ± 12.4	54.1 ± 12.5	63.1 ± 11.7
Subject 3	53.1 ± 10.2	63.7 ± 9.2	61.6 ± 8.2	58.7 ± 10.5	49.4 ± 11.2	53.1 ± 8.3	64.1 ± 11.8	62.5 ± 10.4
Subject 4	60.9 ± 10.4	60.3 ± 10.7	55.6 ± 13.7	58.4 ± 11.6	53.7 ± 9.1	62.5 ± 9.0	64.7 ± 12.4	60.3 ± 11.6
Average	58.3 ± 10.5	<b>61.7 ± 10.7</b>	56.8 ± 11.8	59.3 ± 11.5	50.7 ± 12.2	54.9 ± 9.9	60.5 ± 13.4	<b>62.3 ± 11.5</b>

**Table 1:** Single-trial classification performance achieved in binary localization configuration (single individual directions). Bold values indicate the configurations where the average performance among the subjects minus the standard deviation is greater than the chance level (50%).

	I)	J)	K)	L)	M)	N)	O)	P)
Subject 1	49.2 ± 8.1	49.2 ± 8.6	45.8 ± 8.9	45.6 ± 8.9	42.9 ± 8.5	40.8 ± 10.9	37.7 ± 11.3	41.0 ± 9.4
Subject 2	37.3 ± 9.4	39.4 ± 10.0	42.5 ± 11.1	43.1 ± 9.3	38.9 ± 9.2	32.3 ± 8.2	36.2 ± 8.9	37.7 ± 8.2
Subject 3	41.4 ± 9.4	45.2 ± 8.0	48.1 ± 9.4	50.6 ± 10.6	42.1 ± 8.0	38.7 ± 11.5	43.7 ± 9.1	39.2 ± 10.2
Subject 4	46.9 ± 10.3	51.7 ± 10.5	52.3 ± 11.2	43.3 ± 12.4	47.9 ± 9.9	35.4 ± 10.0	36.9 ± 6.3	47.7 ± 10.5
Average	<b>43.7 ± 9.3</b>	<b>46.3 ± 9.3</b>	<b>47.2 ± 10.1</b>	<b>45.7 ± 10.3</b>	<b>43.0 ± 8.9</b>	36.8 ± 10.1	38.6 ± 8.9	41.4 ± 9.6

**Table 2:** Single-trial classification performance achieved in ternary localization configurations (single individual directions). Chance level is 33%.

	Q)	R)	S)	T)	U)	V)	W)	X)
Subject 1	63.8 ± 6.2	64.8 ± 7.6	65.6 ± 9.0	37.8 ± 6.0	33.0 ± 6.4	28.7 ± 4.6	27.7 ± 6.1	18.0 ± 4.2
Subject 2	55.4 ± 8.0	62.4 ± 5.3	55.1 ± 6.6	32.0 ± 7.0	28.1 ± 8.4	23.4 ± 6.8	25.6 ± 6.0	14.8 ± 4.7
Subject 3	56.0 ± 4.7	67.3 ± 6.9	62.5 ± 7.6	37.8 ± 8.2	33.6 ± 6.2	23.9 ± 5.4	25.7 ± 8.2	18.3 ± 4.8
Subject 4	57.9 ± 5.3	62.1 ± 5.1	68.1 ± 7.9	40.3 ± 9.8	34.7 ± 5.7	27.9 ± 5.3	32.2 ± 8.2	19.1 ± 4.9
Average	<b>58.3 ± 6.1</b>	<b>64.1 ± 6.2</b>	<b>62.8 ± 7.8</b>	<b>37.0 ± 7.7</b>	<b>32.3 ± 6.7</b>	<b>26.0 ± 5.5</b>	<b>27.8 ± 7.1</b>	<b>17.6 ± 4.7</b>

**Table 3:** Single-trial classification performance achieved in binary localization configurations (aggregated directions), and quaternary, quinary and octonary localization configurations (single individual directions). Chance levels are 50%, 25%, 20%, and 12.5% respectively.

tions. The epoch window actually used for the classification is delimited by the two black dashed vertical lines and the green areas indicate time-windows where the ERPs to the two directions were statistically significant ( $p < 0.05$ , Welch's t-test with Bonferroni correction). Some of these windows were statistically significant for multiple spatially related channels, resulting in regions of interest (see the scalp topographies shown in the Figure 5). This topographical information was not explicitly investigated in the present study but will be considered in future analyses. The ERP results also show a qualitative difference between time-locked responses to the first and second sound click. This phenomenon may be a consequence of the processing of a natural sound, whose processing involves phenomena such as temporal prediction and adaptation (speech perception; [19]).

### 3.2 Classification results

Behavioral assessment of the subjects' localization ability indicated, on average, a 98.9% correct identification of the sounds direction. Given the strong localization ability common to all subjects, in future experiments we might avoid the collection of subjective feedback for the perceived sound direction. This would result in a shorter experimental protocol allowing the collection of a higher number of trials in the same amount of time.

The classification performance of the subject-specific models on the test set were the result of a 20 randomized executions of a train-validation-test procedure. Table 1, 2, and 3 show the single-trial results in terms of accuracy, computed as the number of correctly classified trials divided by the total number of trials. Each column of a table represents a different localization configuration and we divided them in three groups based on the arity (i.e. the number of elements in a set) of the classification problem and the size of the classes in terms of directions. Table 1

reports the results of the binary classifications involving couples of individual directions, Table 2 the results of the ternary classifications, and Table 3 the results of the quaternary, quinary and octonary classifications (T, U, V, W, and X) and the results of binary classifications involving groups of directions (Q, R, and S). Some configurations were modeled better than others with the selected approach (bold values in Table 1, 2, and 3). The criterion used for their identification consists in the fact that the average of the subject-specific performance minus the standard deviation is greater than the chance level. The mathematical chance level (defined by the number of classes involved in the classification) was validated empirically by means of a permutation test based on random shuffling. As expected, given that the number of trials per class was balanced, the empirical baseline matched the mathematical one.

The results of the localization configurations involving aggregations of directions based on their spatial distribution are slightly higher and present a lower variance with respect to counterparts involving single directions (probably due to the greater number of trials available for training and validation). It was not validated statistically, nevertheless it seems to give credit to the hypothesis of common patterns in the brain responses associated to spatially related directions. This aspect will be taken into account in the design of the experiment and future data analyses for a proper investigation.

#### 4. CONCLUSIONS

We investigated the localization of a natural sound from EEG signals. The single-trial classification results, albeit drawn from a small dataset, demonstrate the effectiveness of the approach both from a methodological and practical point of view. In particular, all four subjects achieved significant classification accuracies, especially in the localization configurations involving aggregations of directions. We hypothesize this to be due to the greater number of trials per class available in those configurations, therefore encouraging the optimization of the experimental protocol allowing the acquisition of more trials. Considering the extremely high localization scores seen in this pilot, one option would be to avoid the collection the subjective feedback. Likewise, these great performances in the configurations involving aggregations of directions seem to indicate the presence of common patterns in the EEG responses associated to spatially related directions.

In the experiment we employed a particular natural sound featuring two prominent clicks separated by 250 ms, thus evoking two overlapping ERPs that, interestingly, elicited responses with different temporal patterns. This preliminary result suggests that the evoked responses to the two clicks may reflect distinct cortical contributions that were overlooked by previous studies with artificial auditory stimuli, thus suggests further extensive investigation with more complex natural sounds.

#### 5. ACKNOWLEDGMENTS

The authors gratefully acknowledge the European Commission for its support of the Marie Skłodowska Curie program through the H2020 ETN PBNv2 project (GA 721615) as well as Toyota Motor Europe and the University of Parma for providing the infrastructure needed for this experiment.

#### 6. REFERENCES

- [1] K. van der Heijden, J. P. Rauschecker, B. de Gelder, and E. Formisano, “Cortical mechanisms of spatial hearing,” *Nature Reviews Neuroscience*, vol. 20, no. 10, pp. 609 – 623, 2019.
- [2] C. Alain, S. R. Arnott, S. Hevenor, S. Graham, and C. L. Grady, ““what” and “where” in the human auditory system,” *Proceedings of the National Academy of Sciences*, vol. 98, no. 21, pp. 12301–12306, 2001.
- [3] P. P. Maeder, R. A. Meuli, M. Adriani, A. Bellmann, E. Fornari, J.-P. Thiran, A. Pittet, and S. Clarke, “Distinct pathways involved in sound recognition and localization: A human fmri study,” *NeuroImage*, vol. 14, no. 4, pp. 802 – 816, 2001.
- [4] J. Ahveninen, I. P. Jääskeläinen, T. Raij, G. Bonmassar, S. Devore, M. Hämäläinen, S. Levänen, F.-H. Lin, M. Sams, B. G. Shinn-Cunningham, T. Witzel, and J. W. Belliveau, “Task-modulated “what” and “where” pathways in human auditory cortex,” *Proceedings of the National Academy of Sciences*, vol. 103, no. 39, pp. 14608–14613, 2006.
- [5] D. J. Barrett and D. A. Hall, “Response preferences for “what” and “where” in human non-primary auditory cortex,” *NeuroImage*, vol. 32, no. 2, pp. 968 – 977, 2006.
- [6] S. Arnott, M. Binns, C. Grady, and C. Alain, “Assessing the auditory dual-pathway model in humans,” *NeuroImage*, vol. 22, pp. 401–8, 06 2004.
- [7] J. Lewald and S. Getzmann, “When and where of auditory spatial processing in cortex: A novel approach using electrotomography,” *PLOS ONE*, vol. 6, pp. 1–17, 09 2011.
- [8] L. Y. Deouell, A. S. Heller, R. Malach, M. D’Esposito, and R. T. Knight, “Cerebral responses to change in spatial location of unattended sounds,” *Neuron*, vol. 55, no. 6, pp. 985 – 996, 2007.
- [9] S. Getzmann and J. Lewald, “Effects of natural versus artificial spatial cues on electrophysiological correlates of auditory motion,” *Hearing Research*, vol. 259, no. 1, pp. 44 – 54, 2010.
- [10] S. Sur, V. K. Sinha, and J. Horváth, “Event-related potential: An overview,” *Industrial psychiatry journal*, vol. 18, no. 1, pp. 70–73, 2009.

- [11] A. Farina, A. Capra, P. Martignon, S. Fontana, F. Adriaensen, P. Galaverna, and D. Malham, “Three-dimensional acoustic displays in a museum employing WFS (Wave Field Synthesis) and HOA (High Order Ambisonics),” in *ICSV14*, 2007.
- [12] M. Binelli and A. Farina, “Digital equalization of automotive sound systems employing spectral smoothed fir filters,” 10 2008.
- [13] Y. Ando, *Concert Hall Acoustics*. Springer, 1985.
- [14] T. R. Mullen, C. A. E. Kothe, Y. M. Chi, A. Ojeda, T. Kerth, S. Makeig, T. Jung, and G. Cauwenberghs, “Real-time neuroimaging and cognitive monitoring using wearable dry eeg,” *IEEE Transactions on Biomedical Engineering*, vol. 62, pp. 2553–2567, Nov 2015.
- [15] C. Chang, S. Hsu, L. Pion-Tonachini, and T. Jung, “Evaluation of artifact subspace reconstruction for automatic artifact components removal in multi-channel eeg recordings,” *IEEE Transactions on Biomedical Engineering*, pp. 1–1, 2019.
- [16] A. de Cheveigné and L. C. Parra, “Joint decorrelation, a versatile tool for multichannel data analysis,” *NeuroImage*, vol. 98, pp. 487 – 505, 2014.
- [17] T. Onitsuka, N. Oribe, and S. Kanba, “Chapter 13 - neurophysiological findings in patients with bipolar disorder,” in *Application of Brain Oscillations in Neuropsychiatric Diseases* (E. Başar, C. Başar-Eroğlu, A. Özerdem, P. Rossini, and G. Yener, eds.), vol. 62 of *Supplements to Clinical Neurophysiology*, pp. 197 – 206, Elsevier, 2013.
- [18] J. D. Rollnik, “Chapter 11 - clinical neurophysiology of neurologic rehabilitation,” in *Clinical Neurophysiology: Diseases and Disorders* (K. H. Levin and P. Chauvel, eds.), vol. 161 of *Handbook of Clinical Neurology*, pp. 187 – 194, Elsevier, 2019.
- [19] J. Y. Choi and T. K. Perrachione, “Time and information in perceptual adaptation to speech,” *Cognition*, vol. 192, p. 103982, 2019.