
This is an electronic reprint of the original article.
This reprint may differ from the original in pagination and typographic detail.

Kurmanavičiūtė, Dovilė; Kataja, Hanna; Jas, Mainak; Vällilä, Anne; Parkkonen, Lauri
Target of selective auditory attention can be robustly followed with MEG

Published in:
Scientific Reports

DOI:
[10.1038/s41598-023-37959-4](https://doi.org/10.1038/s41598-023-37959-4)

Published: 06/07/2023

Document Version
Publisher's PDF, also known as Version of record

Published under the following license:
CC BY

Please cite the original version:
Kurmanavičiūtė, D., Kataja, H., Jas, M., Vällilä, A., & Parkkonen, L. (2023). Target of selective auditory attention can be robustly followed with MEG. *Scientific Reports*, 13(1), 1-10. Article 10959.
<https://doi.org/10.1038/s41598-023-37959-4>



OPEN Target of selective auditory attention can be robustly followed with MEG

Dovilė Kurmanavičiūtė^{1✉}, Hanna Kataja¹, Mainak Jas^{1,3}, Anne Väilä¹ & Lauri Parkkonen^{1,2}

Selective auditory attention enables filtering of relevant acoustic information from irrelevant. Specific auditory responses, measurable by magneto- and electroencephalography (MEG/EEG), are known to be modulated by attention to the evoking stimuli. However, such attention effects have typically been studied in unnatural conditions (e.g. during dichotic listening of pure tones) and have been demonstrated mostly in averaged auditory evoked responses. To test how reliably we can detect the attention target from unaveraged brain responses, we recorded MEG data from 15 healthy subjects that were presented with two human speakers uttering continuously the words “Yes” and “No” in an interleaved manner. The subjects were asked to attend to one speaker. To investigate which temporal and spatial aspects of the responses carry the most information about the target of auditory attention, we performed spatially and temporally resolved classification of the unaveraged MEG responses using a support vector machine. Sensor-level decoding of the responses to attended vs. unattended words resulted in a mean accuracy of $79\% \pm 2\%$ ($N = 14$) for both stimulus words. The discriminating information was mostly available 200–400 ms after the stimulus onset. Spatially-resolved source-level decoding indicated that the most informative sources were in the auditory cortices, in both the left and right hemisphere. Our result corroborates attention modulation of auditory evoked responses and shows that such modulations are detectable in unaveraged MEG responses at high accuracy, which could be exploited e.g. in an intuitive brain–computer interface.

Selective auditory attention enables filtering of relevant acoustic information from irrelevant and is often studied using dichotic listening^{1,2} where the listener is exposed to simultaneous but different auditory streams to each ear and is asked to follow one stream while suppressing the other, akin to the cocktail party problem³. Selectively attending to one stream manifests as changes in auditory evoked responses that can be measured non-invasively with electroencephalography (EEG) and magnetoencephalography (MEG)^{4–8}.

More recently, machine-learning methods have been applied to EEG/MEG data to study attention modulation of transient auditory evoked responses⁹, auditory steady-state responses^{10,11} or responses to continuous speech^{12–15}. Exploiting such attention modulation in a brain–computer interface has been probed in several studies^{16–25}, some of which have employed natural sounds as stimuli and yielded a useful-in-practice classification accuracy also when applied to patients that cannot communicate^{26,27}. However, these auditory speller-type BCI systems require extensive training that might be exhausting for a patient. Furthermore, in patients with disorders of consciousness, using this type of a BCI may exceed the capacity of their working memory^{25,28}, which could drastically drop the accuracy. In comparison to speller-BCIs, BCIs based on either speech tracking or on detecting infrequent and unexpected changes in auditory streams could be designed such that their working-memory load is limited^{21,29–31}. However, BCIs utilizing speech tracking often require long data spans (usually tens of seconds) to output one bit since the dynamics of continuous natural speech are complex and thus the responses less salient than those for isolated words or simple tone pips (see e.g., Ref.³²). Yet, near real-time performance has been demonstrated through advanced modelling³³.

Auditory streaming BCIs often employ oddball streams, comprising frequently-occurring stimuli (standard) and a rarely-occurring exception (deviant)³⁰. Selective attention then increases the amplitude of the response to a deviant compared to an unattended stimulus^{16,34,35}. However, this approach allows attention target to be determined only at the rate the deviants are presented, and this rate cannot be increased above 10 or 20% of all stimuli without diminishing the overall amplitude of the deviant responses. Therefore, the information transfer rate of such a BCI remains modest.

¹Department of Neuroscience and Biomedical Engineering, Aalto University, P.O. Box 12200, 00076 Aalto, Finland. ²Aalto Neuroimaging, Aalto University, 00076 Aalto, Finland. ³Athinoula A. Martinos Center for Biomedical Imaging, 149 Thirteenth Street, Charlestown, MA 02129, USA. ✉email: dovile.kurmanaviciute@aalto.fi

In this study, we propose a novel paradigm for eventual BCI applications that differs from the conventional cocktail party problem by employing simple, minimally overlapping word stimuli in two rapid sequences, thereby enabling fast tracking of the target of attention. By embedding sequence deviants, we can also include a task that allows behavioural quantification of the deployment of attention. Our aim was to provide a paradigm that could be efficiently used in a simple yet intuitive brain–computer interface.

To this end, we created an acoustically realistic scene with two concurrent auditory stimulus streams. Stimuli comprised of two human speakers uttering the words “Yes” and “No” in an interleaved manner at -40 and $+40$ degrees from the line forward from the subject, mimicking a real-life situation where two persons are speaking simultaneously on the sides of the subject. In each stream, the pitch of the word alternated (standard) but this implicit rule was occasionally broken by presenting two same-pitch versions of the stimulus word in succession (deviant). We measured MEG in 15 subjects while they were presented with these stimuli and were asked to covertly count these deviants in the attended stream and report the number at the end of each measurement block.

Results

Behavioral data. On average, the subjects reported $40 \pm 17.6\%$ (deviant probability 10%, $N = 5$) and $97 \pm 0\%$ (deviant probability 5%, $N = 6$) of the deviants in the stream they were instructed to attend to. Three subjects were not included in this analysis due to technical problems in collecting their deviant counts.

Sensor-level analysis. Time-resolved decoding was performed on the unaveraged epochs comprising all channels at each time point. At the group level, decoding “Attended No” vs. “Unattended No” and “Attended Yes” vs. “Unattended Yes” both showed peaks around 160 ms (Fig. 1a).

Spatially-resolved decoding indicated that the most informative signals arose from temporal regions; the patterns of decoding accuracy were qualitatively similar across the subjects; see Fig. 1b for a representative subject and for the group result.

To aim at the highest accuracy in determining the direction of attention, we also decoded using the entire epoch, that is, all time points and all channels at once. First, we tested with 1-s epochs ($-200 \dots 800$ ms) which yielded a mean accuracy of $79\% \pm 2\%$ (range 67–91%) for “Attended No” vs. “Unattended No” and $79\% \pm 2\%$ (range 68–91%) for “Attended Yes” vs. “Unattended Yes”; see Fig. 2. The group mean accuracy was not significantly different between the stimulus words (paired t-test; $p = 0.874$).

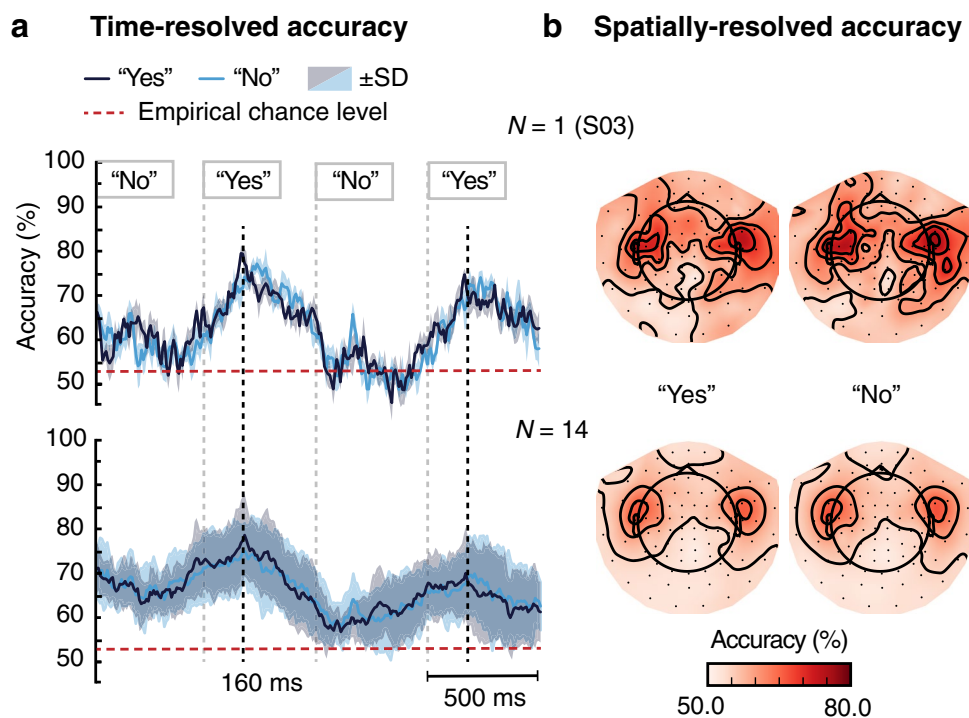


Figure 1. Temporally and spatially resolved decoding reveals highest decoding accuracy around 160 ms after each stimulus word and in the MEG channels above the auditory cortices. **(a)** Time-resolved decoding in a 2-s time window of attended vs. unattended word stimulus plotted for a representative subject (top) and for the group (bottom). The mean decoding accuracy is shown as a dark blue line for the “Yes” and as a light blue line for the “No” stimulus word. The standard deviation (SD), computed across the cross-validation folds of the classifier training and testing, is shown as dark/light blue shading. **(b)** Spatially-resolved decoding accuracy maps in a representative subject (top) and at the group level (bottom) for attended vs. unattended stimulus words. Prior to the decoding, epochs of the high- and low-pitch words were concatenated.

Prolonging the decoding epoch to 2 s (– 500 ... 1500 ms) resulted in an average decoding accuracy of $83\% \pm 2\%$ (range 71–94%) for “Attended Yes” vs. “Unattended Yes” and $83\% \pm 2\%$ (range 68–94%) for “Attended No” vs. “Unattended No”. Using the 2-s vs. 1-s epochs increased the classification accuracy in all 14 subjects, which is a statistically significant change (binomial test: $p = 0.000061$).

To further characterize the stimulus-related information present in the auditory evoked responses, we decoded also for the stimulus word (not for attention). The mean accuracy was $88\% \pm 2\%$ (range 75–96%) for the 1-s epochs and $93\% \pm 1\%$ (range 82–99%) for the 2-s epochs when decoding “Attended Yes” vs. “Attended No”. Similarly, when decoding for “Unattended Yes” vs. “Unattended No”, we obtained an average decoding accuracy of $88\% \pm 2\%$ (range 76–95%) for the 1-s epochs and $92\% \pm 2\%$ (range 77–98%) for the 2-s epochs. Thus, the accuracy of word-wise decoding did not depend significantly on whether the words were attended or not ($p = 0.80$ for the short epochs and $p = 0.53$ for the long epochs). Compared to attention decoding, this word-wise decoding gave statistically significantly higher accuracy for both the short ($p < 0.005$ for all four possible comparisons) and long ($p < 0.005$) epochs.

We also performed pitch-wise decoding (unaveraged evoked responses to high-pitch vs. low-pitch versions of the word stimuli), which yielded above chance-level decoding accuracy of $70\% \pm 3\%$ (“Attended No”), $71\% \pm 3\%$ (“Unattended No”), $71\% \pm 3\%$ (“Attended Yes”) and $72\% \pm 3\%$ (“Unattended Yes”) across all subjects ($N = 14$).

The average evoked responses to each attention condition (“Attended Yes”, “Unattended Yes”, “Attended No”, “Unattended No”) for a single subject and for the group can be found in Supplementary Fig. S1. In that figure, each condition represents pooled responses to the low- and high-pitch stimuli. These average evoked responses were computed only for the sensor- and source-level visualizations, and all decoding was performed on unaveraged (single-trial) responses.

The group-averaged evoked responses peaked at 250 ms (at channel ‘MEG 1322’) after the stimulus onset for the “Attended Yes” and at 136 ms (‘MEG 1322’) for the “Unattended Yes” condition. For the condition “Attended No”, the responses peaked at 340 ms (‘MEG 0242’) and for “Unattended No” at 350 ms (‘MEG 0212’). The planar gradient strength maps (Supplementary Fig. S1) are compatible with sources in auditory cortices.

Source-level analysis. The group-level source estimates depicted in the Fig. 3 show the responses to attended and unattended word stimuli at three different latencies. In the right hemisphere, the activation peaked at 270 ms after the onset of the attended “Yes”. Activation to the attended “No” peaked at 330 ms in the left hemisphere after the stimulus onset. The interindividual variation in the response latencies and amplitudes was considerably higher in the left vs. right hemisphere, which led to smearing of the group-average source dynamics in the left hemisphere (Fig. 3).

Spatial-searchlight decoding (Fig. 4) revealed that the source signals giving the highest decoding accuracy arose from the auditory cortices but significantly also from sensorimotor cortex that has been associated with auditory stimuli processing as well³⁶. However, the spatial peaks of accuracy, as show in Fig. 4, did not align well in time and space across subjects, which led to low group-average accuracy for any single location on the cortex.

Discussion

In this study, we recorded brain signals while presenting simple, minimally overlapping spoken-word stimuli and demonstrated that the target of selective auditory attention to these concurrent streams can robustly (accuracy on average 79% and up to 91% in the best-performing subject) be decoded from just 1 s of MEG data. The decoding

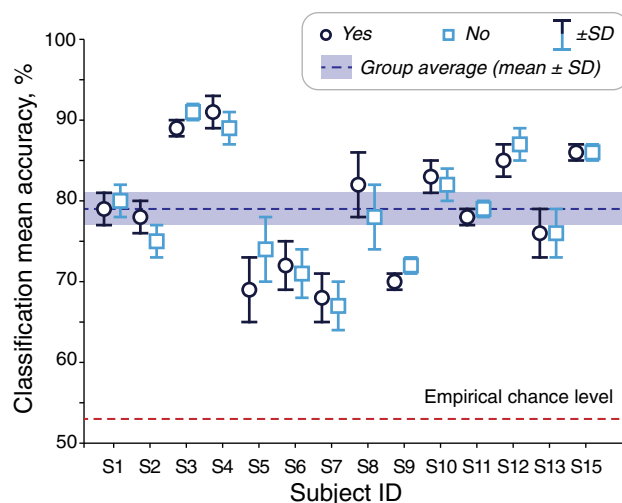


Figure 2. Target of selective auditory attention could be reliably detected in all subjects. The dark blue circles and light blue squares indicate the accuracy of the entire-epoch (all data points in a 1-s window and all channels given to the decoder) classification of responses to attended vs. unattended “Yes” and “No” word-stimuli for all subjects. The standard deviation (SD) was computed over the five cross-validation folds of the decoder and is shown as plot whiskers for each subject.

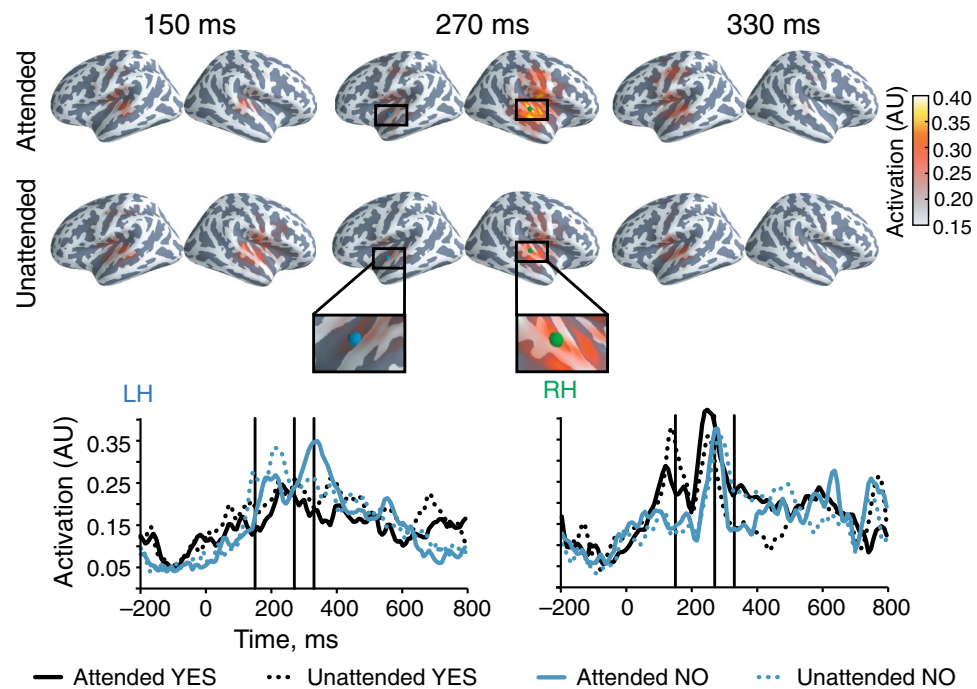


Figure 3. Source estimation of the MEG evoked responses corroborates attention modulation in auditory cortical regions. *Top:* Source estimates of the evoked responses to attended and unattended word stimuli; estimates averaged across the group ($N = 11$). The colour represents the source amplitude first normalized to the absolute peak value of each individual source estimate and then averaged across subjects. *Bottom:* The temporal dynamics of left (LH) and right (RH) auditory-cortex activation to attended and unattended stimulus words (“Yes”/“No”), extracted from the source estimate at the coloured dots (green/blue in the top panel).

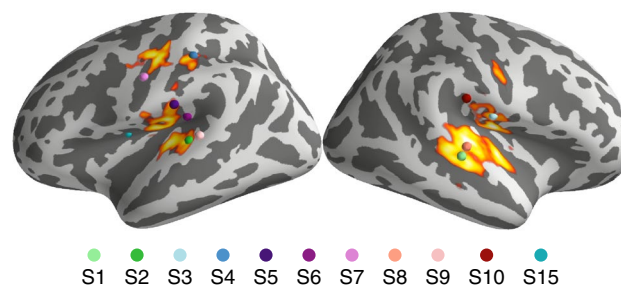


Figure 4. Auditory cortical regions generate the signals most informative of the attended stimulus stream. The colour gradient (yellow highest) represents the source-space spatial searchlight decoding result averaged across the subject group ($N = 11$). Each color dot represents the accuracy peak in one subject.

accuracy peaked around 160 ms after the onset of both stimulus words and remained above chance level for several hundred milliseconds. The highest accuracy was obtained with signals arising from the auditory cortices.

Previous studies have shown that non-semantic acoustic properties, such as sound-source-specific pitch³⁷, are crucial for solving the cocktail party problem at the perceptual level^{38,39}. Although our paradigm did not strictly adhere to the conventional cocktail-party definition as our stimulus words overlapped only minimally, the cortical representation of these properties and the MEG signals evoked by them may still contribute to our ability to decode the target of attention. The presence of pitch-related information in our data was demonstrated by the above-chance-level accuracy when decoding for the pitch (low vs. high pitch) instead of attention. Similarly, we obtained a high accuracy when decoding for the stimulus word (group mean of 88% for the 1-s epochs) or for the stimulus pitch (range 70–72%) instead of attention, which also speaks for the presence acoustic information in the MEG responses; differing stimulus durations (the two words and their pitch variants are of slightly different lengths) and the corresponding variation in the evoked responses is probably the most important feature that the decoder uses to achieve this higher accuracy.

Our behavioural results showed that by doubling the number of deviants in the stimulus sequence (10% instead of 5%), their detection rate dropped drastically, suggesting that the higher deviant rate was too demanding a task. The load theory of attention⁴⁰—mostly studied in the visual domain—suggests that with a higher

attention demand, the task performance drops, which is in line with our data. The current view of load theory⁴¹ is that the load on working memory and cognitive control processes would hamper target detection, whereas load to visual short-term memory would do the opposite, that is, would reduce detecting distractors. However, the debate is still ongoing (for a review, see⁴²).

Earlier studies have demonstrated that rich naturalistic stimuli, compared to monotonous tones, not only improve the users' ergonomic evaluation of the situation but also yield higher decoding accuracy^{21,43}. Further, it has been shown that subjects perform better on selective attention tasks when presented with naturalistic speech in comparison to other kind of naturalistic stimuli⁴⁴. Thus, the naturalistic, spoken-word stimuli that we used have likely contributed to the high classification accuracy.

In our study, the pitch difference between the spoken-word streams due to speaker gender (male and female voices) likely helped focusing attention to one speaker. Yet, the pitch itself did not seem to play a role in either stream alone as the attention decoding accuracy for each word-stream was very similar (Fig. 2).

Stimulus timing likely has an effect on decoding accuracy as it influences the amplitude and latency of the attention-modulated evoked responses. Stimulus onset asynchrony (SOA) has been shown to affect accuracy in decoding attention to simple tones by Höhne et al.⁴⁵; they found that SOA of 1000 ms gave the best decoding accuracy but the highest information transfer rate was achieved with short SOAs (87–175 ms). Other studies that used virtual sound stimuli observed that SOA of 400–600 ms provided the best decoding accuracy^{28,46}. Given those previous studies, our SOA of 1000 ms was likely optimal in terms of decoding accuracy but probably would not have yielded the highest information transfer rate, if our paradigm was applied in a brain–computer interface (BCI).

Typically, using a longer span of data for decoding improves accuracy if all data are informative; for example, Maße and colleagues have demonstrated this in the BCI context⁴⁷. Also our results showed that using the long (2-s) instead of the short (1-s) epoch increased the accuracy of decoding the target of attention in all 14 subjects. Again, for a BCI, the long epochs may not be the optimal choice to maximize the information transfer rate.

Spatial searchlight decoding across the cortex yielded accuracy peaks at locations similar to those of the largest differences in the source estimates of the evoked responses to the two attention conditions. This agreement of the two analysis methods further supports the notion that the selective auditory attention-modulated cortical activity is mostly in the primary auditory cortex⁴⁸. Regardless of the roughly similar cortical location of the most attention-informative source in each subject, these locations did not fully overlap, which led to a dispersed group average even though interindividual variation of cortical anatomy was reduced by surface-based morphing of the individual brains to an average brain. This variation—although minor—in the location and orientation of the source providing the highest decoding accuracy likely means that classifiers do not generalize well across subjects but that the classifier should be trained separately for each subject if one is aiming to the highest accuracy.

Left and right hemispheres are differently specialised to process auditory stimuli. Language-specific areas are typically lateralised to the left hemisphere⁴⁹. For instance, left hemisphere has been found to respond more than the right to the temporal aspects of auditory stimuli⁵⁰. For comparison, right hemisphere has been found to be more involved in spectral processing of e.g. tones and music⁵¹. Previous studies found that right hemisphere responds to the manipulation of pitch in human speech^{50,52,53}.

Based these previous findings and our data, we suggest that selective attention is engaged to follow the regular pitch alternation and thus to support the detection of its infrequent deviants in the indicated stimulus sequence. In our data, such an engagement was manifested in the responses in the right hemisphere (attended vs. unattended stimuli) at around 270 ms after the stimulus onset. The left hemisphere was activated later (at around 330 ms) and had higher activity for the attended vs. unattended stimuli while such a difference was not as clear in the right hemisphere. The peak of left-hemisphere source could potentially be related to the processing of lexical/semantic information as, e.g., for visual word stimuli⁵⁴.

It is conceivable that in our experiment participants applied different strategies of keeping their attention to one auditory stream or they might have even changed their strategy in the course of the experiment. This possibility could be studied by training the decoder by samples from specific parts of the recording (e.g. only from the beginning) and comparing the obtained classification accuracy. In addition, the influence of stimulation rate to selective attention and its decoding from brain signals could be tested. Moreover, future studies could assess individual differences in response latency and spatial patterns on the MEG sensor array that may limit across-subject generalization.

Using the current experimental paradigm, one could test how robustly the observed attention modulation of brain responses could be detected by EEG instead of MEG. Spatial separability of cortical sources is typically poorer in EEG compared to MEG^{55,56} and thus the accuracy of decoding the attended word stream from EEG would likely be lower; yet, the accuracy could remain at a level which enables a portable and intuitive brain–computer interface.

Conclusions

We showed that the attended spoken-word stream can reliably be decoded from just one-second epochs of unaveraged MEG data. The achieved high decoding accuracy shall enable future investigations on the neural mechanisms of attentional selection and it may also be exploited in a MEG- or EEG-based streaming brain–computer interface.

Materials and methods

Participants. Fifteen healthy adult volunteers (4 females, 11 males; mean age 28.8 ± 3.8 years, range 23–38 years) participated in our study. Two subjects were left-handed and the rest right-handed. Participants did not report hearing problems or history of psychiatric disorders. The study was approved by the Aalto University

Research Ethics Committee. The research was carried out in accordance with the guidelines of the Declaration of Helsinki, and the subjects gave written informed consent prior the measurements.

Stimuli and experimental protocol. The subjects were presented with two auditory streams, one comprising the spoken word “Yes” and the other the word “No”. The words alternated such that the words did not overlap. In each stream (“Yes” and “No”), the stimulus onset asynchrony (SOA) was 1000 ms, and the duration of the stimulus words were 450–550 ms depending on the word and its pitch variant. Fig. 5 illustrates the stimulus timing (Fig. 5a), positioning in the acoustic scene (Fig. 5b) and the structure of the stimulus sequence Fig. 5c.

To create a realistic acoustic scene, the stimuli were recorded with a dummy head (Mk II, Cortex Instruments GmbH, Germany) at the center of a room with dimensions comparable to those of the magnetically shielded room where the MEG recordings were performed later. The speakers were standing at about at -40° and $+40^\circ$ degrees from the front-line of the dummy head at a distance of 1.13 m. The word “Yes” was uttered by a female and the word “No” by a male speaker. Thus, in the experiment, the sound of each speaker was presented to both ears of the subject as recorded with the dummy head; see Fig. 5a.

Subjects were asked to attend to one of the two asynchronously alternating spoken word streams at a time by a visual cue (“LEFT-YES” or “RIGHT-NO”) shown next to the fixation cross and accordingly on left or right side of the screen. Each stream had two alternating pitches of the spoken word (denoted as [..., yes, YES, yes, YES, ...] and as [..., no, NO, no, NO, ...]). The original voice recordings were used as the low-pitch stimuli, and the pitch was increased by 13% and 15% for the high-pitch versions of “Yes” and “No”, respectively. In each spoken-word stream, occasional violations (deviants) of otherwise regular pitch alternation (standards) occurred. The subjects were instructed to count the deviants in order to keep their attention to the indicated spoken word stream stimuli. The stimulus sequence always started with the low-pitch word; see Fig. 5c.

Deviants were presented with the probability of 10% in both streams for the first seven subjects and with probability of 5% for the rest of the subjects. The deviant frequency was decreased based on subject feedback to reduce the mental load of memorizing the deviant count.

The experiment comprised 8 blocks, each lasting about 135 s. Two seconds before a block started, the subject was instructed to direct his/her attention to one of the streams by the cues “LEFT-YES” or “RIGHT-NO” on the screen. The task of the subject was to focus on the indicated word stream, covertly count the deviants, maintain gaze at the fixation cross displayed on the screen and verbally report the count at the end of the block.

During the cue “LEFT-YES”, the evoked responses to the word “Yes” were assigned to the condition “Attended Yes” and the evoked responses to “No” were assigned to the condition “Unattended No”. Similarly, during the cue “RIGHT-NO”, evoked responses to “No” were assigned to the condition “Attended No” and, accordingly, evoked responses to “Yes” were assigned to the condition “Unattended Yes”.

The experiment always started with a block with the cue “LEFT-YES” and was followed by a block with the cue “RIGHT-NO”. The order of the remaining six blocks was randomized across subjects. The first blocks were not randomised due to our main goal to use the first two blocks for training the classifier. The total length of the experiment was 50–60 min including the breaks between the blocks.

PsychoPy version 1.79.01^{57,58} Python package was used for controlling and presenting the auditory stimuli and visual instructions. The stimulation was controlled by a computer running Windows 2003 for the first nine subjects and Linux Ubuntu 14.04 for the rest. Auditory stimuli were delivered by a professional audio card (E-MU 1616m PCIe, E-MU Systems, Scotts Valley, CA, USA), an audio power amplifier (LTO MACRO 830, Sekaku

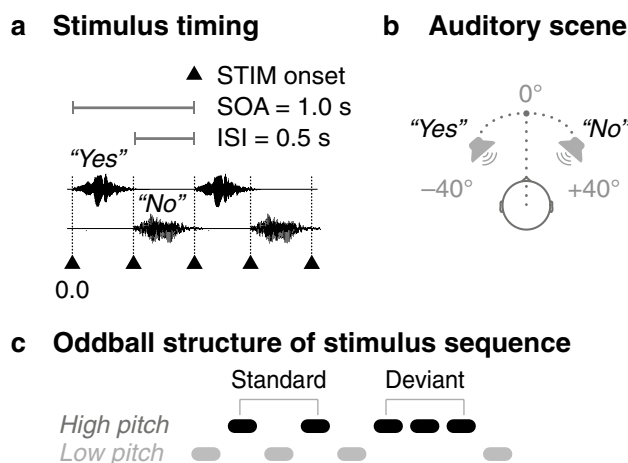


Figure 5. Experimental design: By a text cue on the screen (“LEFT-YES” or “RIGHT-NO”), subjects were instructed to attend either the “Yes” or the “No” stimulus stream. **(a)** Stimulus timing in 2-s time window with the onsets of the stimulus words marked as black triangles. **(b)** The virtual arrangement of the speakers uttering the words “Yes” and “No” with respect to the subject. **(c)** The structure of each stimulus stream. The low- and high-pitch versions of the word stimuli alternate (standard) but occasionally three high-pitch versions are presented (deviant).

Electron Industry Co., Ltd, Taichung, Taiwan), and custom-built loudspeaker units outside of the shielded room and plastic tubes conveying the stimuli separately to the ears. Sound pressure was adjusted to a comfortable level for each subject individually. Due to timing inaccuracies in the stimulus presentation system, the delay from the trigger to sound onset for the “Yes” stimuli varied with a standard deviation of 7 ms while that for “No” varied with a standard deviation of 11 ms.

MEG data acquisition. MEG measurements were performed with a whole-scalp 306-channel Elekta—Neuromag VectorView MEG system (MEGIN Oy, Helsinki, Finland) at the MEG Core of Aalto Neuroimaging, Aalto University. During acquisition, the data were filtered to 0.1–330 Hz and sampled at 1 kHz. Prior to the MEG recording, anatomical landmarks (nasion, left and right preauricular points), head-position indicator coils, and additional scalp-surface points (around 100) were digitized using an Isotrak 3D digitizer (Polhemus Navigational Sciences, Colchester, VT, USA). Bipolar electrooculogram (EOG) with electrodes positioned around the right eye (laterally and below) was recorded. Fourteen of the 15 subjects were recorded with continuous head movement tracking. All subjects were measured in the seated position. The back-projection screen for delivering the visual instructions was 1 m from the eyes of the subject. If needed, vision was corrected by nonmagnetic goggles.

The MEG recording of one subject had technical problems and this dataset had to be dropped from the analysis.

Data pre-processing. The MaxFilter software (version 2.2.10; MEGIN Oy, Helsinki, Finland) was applied to all MEG data (magnetometers and planar gradiometers) to suppress external interference using temporal signal space separation and to compensate for head movements⁵⁹. Further analysis was performed using the MNE-Python^{60,61}, version 0.21; and ScikitLearn⁶², version 0.23.2; software packages.

Infinite-impulse-response filters (4th-order Butterworth, applied both forward and backward in time) were employed to filter the unaveraged MEG data to 0.1–30 Hz for visualization of the evoked responses and for sensor- and source-level decoding. Ocular artifacts were suppressed by removing those independent components (1–4 per subject, on average 3) that correlated most with the EOG signal.

For the subsequent data analysis, only planar gradiometers were used due to the straightforward interpretation of their spatial pattern; they show the maximum signal right above the active source.

Epochs with two different pre- (200 ms and 500 ms) and post-stimulus (800 ms and 1500 ms) periods were extracted from the MEG data at every word stimulus. Epochs were rejected if any of the planar gradiometer signals exceeded 4000 fT/cm. Deviant epochs were excluded from data analysis. The trial counts were equalized, and the responses averaged across each condition (“Attended Yes”, “Attended No”, “Unattended Yes” and “Unattended No”) for visualization and source estimation.

Source estimation. Head models were constructed based on individual magnetic resonance images (MRIs) by applying the watershed algorithm implemented in the FreeSurfer software^{63–65}, version 5.3. Using the MNE software, single-compartment boundary element models (BEM) comprising 5120 triangles were then created based on the inner skull surface. In addition to the one subject with technical problems in MEG recording, the MRIs of three subjects were not available, leaving 11 subjects for the source estimation.

For the source space, the cortical mantle was segmented from MRIs using FreeSurfer and the resulting triangle mesh was subdivided to 4098 sources per hemisphere. The dynamic statistical parametric mapping⁶⁶, *dSPM*; variant of minimum-norm estimation was applied to model the activity at these sources. The noise covariance was estimated from the 2-min resting-state measurement of each subject. These data were pre-processed similarly as the task-related data.

The source amplitudes for the attention conditions “Attended Yes”, “Unattended Yes”, “Attended No” and “Unattended No” were estimated for all subjects individually. For the group-level source estimate, the obtained source amplitudes were first normalized such that the absolute peak value of the attended condition became one, the estimates were morphed to the FreeSurfer average brain and then averaged across subjects. The morphing procedure from individual brains to the average brain is described by Greve et al.⁶⁷.

Decoding. *Sensor-level decoding.* A linear support vector machine⁶⁸, *SVM*; classifier implemented in the Scikit-learn package⁶² was applied to unaveraged epochs to decode the conditions “Attended Yes” vs. “Unattended Yes” and “Attended No” vs. “Unattended No”. For comparison, decoding was also performed stimulus-word-wise, i.e. “Attended Yes” + “Unattended Yes” vs. “Attended No” + “Unattended No”. In addition, pitch-wise and single pitch variant attention-wise decoding were performed.

The pre-processed MEG data (filtered to 0.1–30 Hz) were down-sampled by a factor of 8 to a sampling rate of 125 Hz to reduce the number of features while preserving sufficient temporal information. Amplitudes of the planar gradiometer channels were concatenated to form the feature vector. Shuffled five-fold cross-validation (CV) was applied with an 80/20 split; 80% of data were used for training and the rest for testing. The empirical chance level was around 55% for our sample size of 500 epochs in this two-class decoding task⁶⁹. We also verified the empirical chance level by the method by Ojala and Garriga⁷⁰ and in our case it was 53%.

Decoding was separately performed on data of (1) the entire 2-s epoch (250 time points × 204 channels; *long-epoch decoding*), (2) the entire 1-s epoch (125 time points × 204 channels; *short-epoch decoding*), (3) one time point (1 time point × 204 channels; *time-resolved decoding*), and (4) one channel (250 time points × 1 channel; *spatially-resolved decoding*).

Source-level decoding. A linear SVM decoder with five-fold cross-validation (80%/20% split for training/testing) was applied to the individual source estimates for the attention conditions “Attended Yes” vs. “Unattended Yes” and “Attended No” vs. “Unattended No” calculated from all MEG planar gradiometer channels. A spatial searchlight decoding across the source space was used on the 1-s (– 200 to 800 ms after stimulus onset) epochs, and the resulting accuracy maps were morphed to the FreeSurfer average brain and averaged across the subjects ($N = 11$). The accuracy maps for attention conditions “Attended Yes” vs. “Unattended Yes” and “Attended No” vs. “Unattended No” were then averaged to obtain a general accuracy map.

Data availability

The datasets generated and analysed during the current study are not publicly available due to the local legislation on research on humans but are available from the corresponding author on reasonable request.

Received: 20 March 2023; Accepted: 30 June 2023

Published online: 06 July 2023

References

1. Soveri, A. *et al.* Modulation of auditory attention by training: Evidence from dichotic listening. *Exp. Psychol.* **60**, 44–52. <https://doi.org/10.1027/1618-3169/a000172> (2013).
2. Tallus, J., Soveri, A., Hämäläinen, H., Tuomainen, J. & Laine, M. Effects of auditory attention training with the dichotic listening task: Behavioural and neurophysiological evidence. *PLoS ONE* **10**, e0139318. <https://doi.org/10.1371/journal.pone.0139318> (2015).
3. Cherry, E. C. Some experiments on the recognition of speech, with one and with two ears. *J. Acoust. Soc. Am.* **25**, 975–979 (1953).
4. Hillyard, S. A., Hink, R. F., Schwent, V. L. & Picton, T. W. Electrical signs of selective attention in the human brain. *Science* **182**, 177–180 (1973).
5. Picton, T. W. & Hillyard, S. A. Effects of attention. Human auditory evoked potentials. II. *Electroencephalogr. Clin. Neurophysiol.* **36**, 191–200. [https://doi.org/10.1016/0013-4694\(74\)90156-4](https://doi.org/10.1016/0013-4694(74)90156-4) (1974).
6. Woods, D. L., Hillyard, S. A. & Hansen, J. C. Event-related brain potentials reveal similar attentional mechanisms during selective listening and shadowing. *J. Exp. Psychol. Hum. Percept. Perform.* **10**, 761–777 (1984).
7. Woldorff, M. G. & Hillyard, S. A. Modulation of early auditory processing during selective listening to rapidly presented tones. *Electroencephalogr. Clin. Neurophysiol.* **79**, 170–191 (1991).
8. Woldorff, M. G. *et al.* Modulation of early sensory processing in human auditory cortex during auditory selective attention. *Proc. Natl. Acad. Sci. U.S.A.* **90**, 8722–8726. <https://doi.org/10.1073/pnas.90.18.8722> (1993).
9. King, J. R. *et al.* Single-trial decoding of auditory novelty responses facilitates the detection of residual consciousness. *Neuroimage* **83**, 726–738. <https://doi.org/10.1016/j.neuroimage.2013.07.013> (2013).
10. Kim, D.-W. *et al.* Classification of selective attention to auditory stimuli: Toward vision-free brain–computer interfacing. *J. Neurosci. Methods* **197**, 180–185. <https://doi.org/10.1016/j.jneumeth.2011.02.007> (2011).
11. Kaongoen, N. & Jo, S. A novel hybrid auditory BCI paradigm combining ASSR and P300. *J. Neurosci. Methods* **279**, 44–51. <https://doi.org/10.1016/j.jneumeth.2017.01.011> (2017).
12. Ding, N. & Simon, J. Z. Neural coding of continuous speech in auditory cortex during monaural and dichotic listening. *J. Neurophysiol.* **107**, 78–89. <https://doi.org/10.1152/jn.00297.2011> (2012).
13. Mirkovic, B., Debener, S., Jaeger, M. & De Vos, M. Decoding the attended speech stream with multi-channel EEG: Implications for online, daily-life applications. *J. Neural Eng.* **12**, 046007. <https://doi.org/10.1088/1741-2560/12/4/046007> (2015).
14. Biesmans, W., Das, N., Francart, T. & Bertrand, A. Auditory-inspired speech envelope extraction methods for improved EEG-based auditory attention detection in a cocktail party scenario. *IEEE Trans. Neural Syst. Rehabil. Eng.* **25**, 402–412. <https://doi.org/10.1109/TNSRE.2016.2571900> (2017).
15. Fuglsang, S. A., Dau, T. & Hjortkjaer, J. Noise-robust cortical tracking of attended speech in real-world acoustic scenes. *Neuroimage* **156**, 435–444. <https://doi.org/10.1016/j.neuroimage.2017.04.026> (2017).
16. Furdea, A. *et al.* An auditory oddball (P300) spelling system for brain–computer interfaces. *Psychophysiology* **46**, 617–625. <https://doi.org/10.1111/j.1469-8986.2008.00783.x> (2009).
17. Halder, S. *et al.* An auditory oddball brain–computer interface for binary choices. *Clin. Neurophysiol.* **121**, 516–523. <https://doi.org/10.1016/j.clinph.2009.11.087> (2010).
18. Schreuder, M., Blankertz, B. & Tangermann, M. A new auditory multi-class brain–computer interface paradigm: Spatial hearing as an informative cue. *PLoS ONE* **5**, e9813. <https://doi.org/10.1371/journal.pone.0009813> (2010).
19. Höhne, J., Schreuder, M., Blankertz, B. & Tangermann, M. A novel 9-class auditory ERP paradigm driving a predictive text entry system. *Front. Neurosci.* **5**, 99. <https://doi.org/10.3389/fnins.2011.00099> (2011).
20. Nambu, Isao *et al.* Estimating the intended sound direction of the user: Toward an auditory brain–computer interface using out-of-head sound localization. *PLoS ONE* **8**, e57174. <https://doi.org/10.1371/journal.pone.0057174> (2013).
21. Hill, N. J. *et al.* A practical, intuitive brain–computer interface for communicating “yes” or “no” by listening. *J. Neural Eng.* **11**, 035003. <https://doi.org/10.1088/1741-2560/11/3/035003> (2014).
22. Simon, N. *et al.* An auditory multiclass brain–computer interface with natural stimuli: Usability evaluation with healthy participants and a motor impaired end user. *Front. Hum. Neurosci.* **8**, 1039. <https://doi.org/10.3389/fnhum.2014.01039> (2015).
23. Halder, S. *et al.* An evaluation of training with an auditory P300 brain–computer interface for the Japanese Hiragana syllabary. *Front. Neurosci.* **10**, 446. <https://doi.org/10.3389/fnins.2016.00446> (2016).
24. Hübner, D., Schall, A., Prange, N. & Tangermann, M. Eyes-closed increases the usability of brain–computer interfaces based on auditory event-related potentials. *Front. Hum. Neurosci.* **12**, 391. <https://doi.org/10.3389/fnhum.2018.00391> (2018).
25. Halder, S., Leinfelder, T., Schulz, S. M. & Kübler, A. Neural mechanisms of training an auditory event-related potential task in a brain–computer interface context. *Hum. Brain Mapp.* **40**, 2399–2412. <https://doi.org/10.1002/hbm.24531> (2019).
26. Kübler, A. *et al.* A brain–computer interface controlled auditory event-related potential (P300) spelling system for locked-in patients. *Ann. N. Y. Acad. Sci.* **1157**, 90–100. <https://doi.org/10.1111/j.1749-6632.2008.04122.x> (2009).
27. Lulé, D. *et al.* Probing command following in patients with disorders of consciousness using a brain–computer interface. *Clin. Neurophysiol.* **124**, 101–106. <https://doi.org/10.1016/j.clinph.2012.04.030> (2013).
28. Käthner, I. *et al.* A portable auditory P300 brain–computer interface with directional cues. *Clin. Neurophysiol.* **124**, 327–338. <https://doi.org/10.1016/j.clinph.2012.08.006> (2013).
29. Hill, N. J. *et al.* Classifying EEG and ECoG signals without subject training for fast BCI implementation: Comparison of nonparalyzed and completely paralyzed subjects. *IEEE Trans. Neural Syst. Rehabil. Eng.* **14**, 183–186. <https://doi.org/10.1109/TNSRE.2006.875548> (2006).
30. Hill, N. J., Moinuddin, A., Häuser, A.-K., Kienle, S. & Schalk, G. Communication and control by listening: Toward optimal design of a two-class auditory streaming brain–computer interface. *Front. Neurosci.* **6**, 181. <https://doi.org/10.3389/fnins.2012.00181> (2012).

31. Hill, N. J. & Schölkopf, B. An online brain–computer interface based on shifting attention to concurrent streams of auditory stimuli. *J. Neural Eng.* **9**, 026011. <https://doi.org/10.1088/1741-2560/9/2/026011> (2012).
32. O’Sullivan, J. A. *et al.* Attentional selection in a cocktail party environment can be decoded from single-trial EEG. *Cereb. Cortex* **25**, 1697–1706. <https://doi.org/10.1093/cercor/bht355> (2015).
33. Miran, S. *et al.* Real-time tracking of selective auditory attention from M/EEG: A Bayesian filtering approach. *Front. Neurosci.* **12**, 262. <https://doi.org/10.3389/fnins.2018.00262> (2018).
34. Sellers, E. W. & Donchin, E. A P300-based brain–computer interface: Initial tests by ALS patients. *Clin. Neurophysiol.* **117**, 538–548. <https://doi.org/10.1016/j.clinph.2005.06.027> (2006).
35. Klobassa, D. S. *et al.* Toward a high-throughput auditory P300-based brain–computer interface. *Clin. Neurophysiol.* **120**, 1252–1261. <https://doi.org/10.1016/j.clinph.2009.04.019> (2009).
36. Preisig, B. C., Riecke, L. & Hervais-Adelman, A. Speech sound categorization: The contribution of non-auditory and auditory cortical regions. *Neuroimage* **258**, 119375. <https://doi.org/10.1016/j.neuroimage.2022.119375> (2022).
37. Woods, D. L. & Alain, C. Conjoining three auditory features: An event-related brain potential study. *J. Cogn. Neurosci.* **13**, 492–509. <https://doi.org/10.1162/08999290152001916> (2001).
38. Bee, M. A. & Micheyl, C. The cocktail party problem: What is it? How can it be solved? And why should animal behaviorists study it? *J. Comp. Psychol.* **122**, 235–251. <https://doi.org/10.1037/0735-7036.122.3.235> (2008).
39. Bronkhorst, A. W. The cocktail-party problem revisited: Early processing and selection of multi-talker speech. *Attent. Percept. Psychophys.* **77**, 1465–1487. <https://doi.org/10.3758/s13414-015-0882-9> (2015).
40. Lavie, N., Hirst, A., De Fockert, J. W. & Viding, E. Load theory of selective attention and cognitive control. *J. Exp. Psychol. Gen.* **133**, 339–354. <https://doi.org/10.1037/0096-3445.133.3.339> (2004).
41. Konstantinou, N. & Lavie, N. Dissociable roles of different types of working memory load in visual detection. *J. Exp. Psychol. Hum. Percept. Perform.* **39**, 919–924. <https://doi.org/10.1037/a0033037> (2013).
42. Brockhoff, L., Schindler, S., Bruchmann, M. & Straube, T. Effects of perceptual and working memory load on brain responses to task-irrelevant stimuli: Review and implications for future research. *Neurosci. Biobehav. Rev.* **135**, 104580. <https://doi.org/10.1016/j.neubiorev.2022.104580> (2022).
43. Höhne, J., Krenzlin, K., Dähne, S. & Tangermann, M. Natural stimuli improve auditory BCIs with respect to ergonomics and performance. *J. Neural Eng.* **9**, 045003. <https://doi.org/10.1088/1741-2560/9/4/045003> (2012).
44. Renvall, H. *et al.* Selective auditory attention within naturalistic scenes modulates reactivity to speech sounds. *Eur. J. Neurosci.* **54**, 7626–7641. <https://doi.org/10.1111/ejn.15504> (2021).
45. Höhne, J. & Tangermann, M. How stimulation speed affects event-related potentials and BCI performance. In *2012 Annual International Conference of the IEEE Engineering in Medicine and Biology Society* 1802–1805. <https://doi.org/10.1109/EMBC.2012.6346300> (2012).
46. Sugi, M. *et al.* Improving the performance of an auditory brain–computer interface using virtual sound sources by shortening stimulus onset asynchrony. *Front. Neurosci.* **12**, 108. <https://doi.org/10.3389/fnins.2018.00108> (2018).
47. Maÿe, A., Rauterberg, R. & Engel, A. K. Instant classification for the spatially-coded BCI. *PLoS ONE* **17**, e0267548. <https://doi.org/10.1371/journal.pone.0267548> (2022).
48. Bidet-Caulet, A. *et al.* Effects of selective attention on the electrophysiological representation of concurrent sounds in the human auditory cortex. *J. Neurosci.* **27**, 9252–9261. <https://doi.org/10.1523/JNEUROSCI.1402-07.2007> (2007).
49. Wernicke, C. Der Aphasische symptomkomplex: Eine psychologische studie auf anatomischer basis (M. Cohn und Weigart, 1874).
50. Zatorre, R. J. Spectral and temporal processing in human auditory cortex. *Cereb. Cortex* **11**, 946–953. <https://doi.org/10.1093/cercor/11.10.946> (2001).
51. Zatorre, R. J. & Baum, S. R. Musical melody and speech intonation: Singing a different tune. *PLoS Biol.* **10**, e1001372. <https://doi.org/10.1371/journal.pbio.1001372> (2012).
52. Jamison, H. L., Watkins, K. E., Bishop, D. V. M. & Matthews, P. M. Hemispheric specialization for processing auditory nonspeech stimuli. *Cereb. Cortex* **16**, 1266–1275. <https://doi.org/10.1093/cercor/bhj068> (2006).
53. Hyde, K. L., Peretz, I. & Zatorre, R. J. Evidence for the role of the right auditory cortex in fine pitch resolution. *Neuropsychologia* **46**, 632–639. <https://doi.org/10.1016/j.neuropsychologia.2007.09.004> (2008).
54. Pyllkänen, L. & Marantz, A. Tracking the time course of word recognition with MEG. *Trends Cogn. Sci.* **7**, 187–189. [https://doi.org/10.1016/S1364-6613\(03\)00092-5](https://doi.org/10.1016/S1364-6613(03)00092-5) (2003).
55. Barth, D. S., Sutherling, W., Broffman, J. & Beatty, J. Magnetic localization of a dipolar current source implanted in a sphere and a human cranium. *Electroencephalogr. Clin. Neurophysiol.* **63**, 260–273. [https://doi.org/10.1016/0013-4694\(86\)90094-5](https://doi.org/10.1016/0013-4694(86)90094-5) (1986).
56. Stok, C. J. The influence of model parameters on EEG/MEG single dipole source estimation. *IEEE Trans. Biomed. Eng.* **34**, 289–296. <https://doi.org/10.1109/TBME.1987.326090> (1987).
57. Peirce, J. W. PsychoPy–psychophysics software in Python. *J. Neurosci. Methods* **162**, 8–13. <https://doi.org/10.1016/j.jneumeth.2006.11.017> (2007).
58. Peirce, J. W. Generating stimuli for neuroscience using PsychoPy. *Front. Neuroinform.* **2**, 10. <https://doi.org/10.3389/neuro.11.010.2008> (2008).
59. Taulu, S. & Kajola, M. Presentation of electromagnetic multichannel data: The signal space separation method. *J. Appl. Phys.* **97**, 124905. <https://doi.org/10.1063/1.1935742> (2005).
60. Gramfort, A. *et al.* MNE software for processing MEG and EEG data. *Neuroimage* **86**, 446–460. <https://doi.org/10.1016/j.neuroimage.2013.10.027> (2014).
61. Gramfort, A. *et al.* MEG and EEG data analysis with MNE–Python. *Front. Neurosci.* **7**, 267. <https://doi.org/10.3389/fnins.2013.00267> (2013).
62. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011).
63. Dale, A. M., Fischl, B. & Sereno, M. I. Cortical surface-based analysis. I. Segmentation and surface reconstruction. *Neuroimage* **9**, 179–194. <https://doi.org/10.1006/nimg.1998.0395> (1999).
64. Fischl, B., Sereno, M. I. & Dale, A. M. Cortical surface-based analysis. II: Inflation, flattening, and a surface-based coordinate system. *Neuroimage* **9**, 195–207. <https://doi.org/10.1006/nimg.1998.0396> (1999).
65. Fischl, B. FreeSurfer. *NeuroImage* **62**, 774–781. <https://doi.org/10.1016/j.neuroimage.2012.01.021> (2012).
66. Dale, A. M. *et al.* Dynamic statistical parametric mapping: Combining fMRI and MEG for high-resolution imaging of cortical activity. *Neuron* **26**, 55–67. [https://doi.org/10.1016/s0896-6273\(00\)81138-1](https://doi.org/10.1016/s0896-6273(00)81138-1) (2000).
67. Greve, D. N. *et al.* A surface-based analysis of language lateralization and cortical asymmetry. *J. Cogn. Neurosci.* **25**, 1477–1492. https://doi.org/10.1162/jocn_a_00405 (2013).
68. Cortes, C. & Vapnik, V. Support-vector networks. *Mach. Learn.* **20**, 273–297. <https://doi.org/10.1023/A:1022627411411> (1995).
69. Combrisson, E. & Jerbi, K. Exceeding chance level by chance: The caveat of theoretical chance levels in brain signal classification and statistical assessment of decoding accuracy. *J. Neurosci. Methods* **250**, 126–136. <https://doi.org/10.1016/j.jneumeth.2015.01.010> (2015).
70. Ojala, M. & Garriga, G. C. Permutation tests for studying classifier performance. *J. Mach. Learn. Res.* **11**, 1833–1863 (2010).

Acknowledgements

This research was supported by Academy of Finland, Grant No. 295075 “NeuroFeed”, and European Research Council, Grant No. 678578 “HRMEG”. The content is solely the responsibility of the authors and does not necessarily represent the official views of the funding organizations. The measurements were conducted at the MEG Core of Aalto Neuroimaging Infrastructure, Aalto University, Finland, and financially supported by Aalto Brain Centre. The authors thank prof. Ville Pulkki and Aalto Acoustics Lab, Aalto University, for the loan and guidance on the use of the dummy head.

Author contributions

D.K., H.K., M.J., A.V. and L.P. designed the experiment. D.K., H.K. and M.J. implemented the experiment. D.K. and H.K. conducted the measurements. D.K. analysed the results and wrote the manuscript. All authors reviewed the manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-37959-4>.

Correspondence and requests for materials should be addressed to D.K.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023