# Multiway canonical correlation analysis of brain data

Alain de Cheveigné [a,b,c,*], Giovanni M. Di Liberto [a,b], Dorothée Arzounian [a,b], Daniel D.E. Wong [a,b], Jens Hjortkjær [d,e], Søren Fuglsang [d], Lucas C. Parra [f]

[a] *Laboratoire des Systèmes Perceptifs, UMR 8248, CNRS, France*
[b] *Département d'Etudes Cognitives, Ecole Normale Supérieure, PSL University, Paris, France*
[c] *UCL Ear Institute, London, United Kingdom*
[d] *Hearing Systems Group, Department of Electrical Engineering, Technical University of Denmark, Denmark*
[e] *Danish Research Centre for Magnetic Resonance, Centre for Functional and Diagnostic Imaging and Research, Copenhagen University Hospital Hvidovre, Denmark*
[f] *City College New York, USA*

ABSTRACT

Brain data recorded with electroencephalography (EEG), magnetoencephalography (MEG) and related techniques often have poor signal-to-noise ratios due to the presence of multiple competing sources and artifacts. A common remedy is to average responses over repeats of the same stimulus, but this is not applicable for temporally extended stimuli that are presented only once (speech, music, movies, natural sound). An alternative is to average responses over multiple subjects that were presented with identical stimuli, but differences in geometry of brain sources and sensors reduce the effectiveness of this solution. Multiway canonical correlation analysis (MCCA) brings a solution to this problem by allowing data from multiple subjects to be fused in such a way as to extract components common to all. This paper reviews the method, offers application examples that illustrate its effectiveness, and outlines the caveats and risks entailed by the method.

## 1. Introduction

Stimulus-driven signals recorded with electroencephalography (EEG), magnetencephalography (MEG) and related techniques compete with much stronger sources within the brain, the body, and the environment. The signal of interest usually represents only a fraction of the power at the electrode or sensor, the rest consisting of noise and artifacts. To overcome those, a common practice is to present the same stimulus multiple times and average the responses over repeated presentations. Supposing that the response is the same for all presentations, and the noise is uncorrelated between presentations, the signal-to-noise power ratio (SNR) improves with the number of repeats. SNR can be further improved by combining data across sensors. This could, for example, be achieved by optimizing *spatial filters* based on assumptions about signal and noise, in particular the assumption that the signal should be the same on all repeats (de Cheveigné and Parra, 2014). Data SNR then benefits from both spatial filtering *and* averaging. Neither is applicable if the stimulus is presented only once, for example because it is long (e.g. a long sample of speech or music), or because one wishes to probe a phenomenon likely to fade with repetitions (e.g. surprise).

Instead of presenting the same stimulus multiple times to one subject, one can also present the same stimulus to multiple subjects just once. To the extent that different subjects' brains are functionally similar, we expect similar responses (Hasson et al., 2004; Dmochowski et al., 2012; Lankinen et al., 2014). Unfortunately, the position or orientation of neural sources relative to sensors or electrodes is likely to differ across subjects, so averaging over subjects in sensor space is suboptimal. In order to compare between subjects, or average over subjects, we first need some way to transform the data of each to a common representation that is comparable across subjects. This can be accomplished with spatial filters that are tuned to each individual subject (e.g. Haxby et al., 2011; Lankinen et al., 2014).

Canonical Correlation Analysis (CCA) is a powerful technique to find linear components that are correlated between two data matrices (Hotelling, 1936). Given two matrices $\mathbf{X}_1$ and $\mathbf{X}_2$ of size $T \times d_1$ and $T \times d_2$, CCA produces transform matrices $\mathbf{V}_1$ and $\mathbf{V}_2$ of sizes $d_1 \times d_0$ and $d_2 \times d_0$, where $d_0$ is at most equal to the smaller of $d_1$ and $d_2$. The columns of $\mathbf{Y}_1 = \mathbf{X}_1 \mathbf{V}_1$ are of norm 1 and mutually uncorrelated between each other, as are the columns of $\mathbf{Y}_2 = \mathbf{X}_2 \mathbf{V}_2$, while, more importantly, corresponding columns from each ("canonical correlate pairs") are maximally

---

correlated. The first pair of canonical correlates (CC) defines the linear combinations of each data matrix with the *highest possible correlation* between them. The next pair of CCs defines the most highly correlated combination orthogonal to the first pair, and so-on. Applied to data from two subjects, CCA can find spatial filters that maximize the brain activity common to both, transforming both subject's data so that they can more easily be compared or averaged. However, CCA does not address the issue of comparing or merging responses across more than two subjects.

Extensions to connect multiple data matrices have been proposed under names such as *multiple CCA* (Gross and Tibshirani, 2015; Witten and Tibshirani, 2009), *multiway CCA* (Sturm, 2016; Zhang et al., 2011), *multiset CCA* (Takane et al., 2008; Correa et al., 2010b, a; Hwang et al., 2012; Lankinen et al., 2014; Zhang et al., 2017; Via et al., 2005a,b; Li et al., 2009), *generalized CCA* (Kiers et al., 1994; Afshin-Pour et al., 2012; Melzer et al., 2001; Tenenhaus, 2011; Tenenhaus et al., 2015; Velden, 2011; Fu et al., 2017), or multiway orthonormalized partial least squares (Sun et al., 2009). This diversity in names covers a diversity of formulations (Kettenring, 1971) that all share the aim of finding components that are similar across data matrices. Recent progress addresses regularization (Tenenhaus, 2011), sparsity (Fu et al., 2017; Tenenhaus et al., 2015), missing data (van de Velden and Takane, 2012), nonlinearity (Melzer et al., 2001), or deep learning (Benton et al., 2017). Using similar techniques, independent Component Analysis (ICA) has been generalized under the name of group ICA (GICA) (Eichele et al., 2011; Calhoun and Adali, 2012; Huster et al., 2015; Huster and Raud, 2018).

CCA has been used extensively for brain data analysis and modality fusion (Sui et al., 2012; Dähne et al., 2015; Dmochowski et al., 2017), and several studies have applied multiway CCA (MCCA) and variants thereof to merge data across subjects (Correa et al., 2010b; Afshin-Pour et al., 2012, 2014; Lankinen et al., 2014, 2018; Li et al., 2009; Hwang et al., 2012; Karhunen et al., 2013; Haxby et al., 2011; Sturm, 2016; Zhang et al., 2017). This paper builds on those studies with the aim to better understand the range of applicability of the tool, what is achieved, and what are the caveats. We describe a simple formulation of MCCA that is easy to understand and explain.

We show that MCCA can be applied effectively to multi-subject datasets of EEG or fMRI, both to *denoise* the data prior to further analyses, and to *summarize* the data and reveal traits common across the population of subjects. MCCA-based denoising yields significantly better scores in an auditory stimulus-response classification task, and MCCA-based joint analysis of fMRI data reveals detailed subject-specific activation topographies. The aims of this paper are (a) to provide an intuitive understanding of MCCA, (b) investigate ways in which it can be put to use, and (c) demonstrate its effectiveness for a range of common tasks in the analysis of brain data.

## 2. Methods

In this section we describe a simple formulation of MCCA, show how it can be applied to a variety of tasks, and give details of the real and synthetic data sets used by the examples reported in the Results section.

### 2.1. Data analysis

**Data model.** Assume a data set consisting of $N$ data matrices, each comprised of a time series matrix $\mathbf{X}_n$ of dimensions $T$ (time) $\times d_n$ (channels). These could represent EEG, MEG or fMRI data recorded from $N$ different subjects in response to the same stimulus. They could also be data from multiple imaging modalities gathered from the same subject. Each matrix $\mathbf{X}_n$ consists of linear combinations of a set of sources $\mathbf{S}$ common to all data matrices, to which is added a "noise" matrix $\mathbf{N}_n$ of time series uncorrelated with $\mathbf{S}$, and uncorrelated with the noise matrices $\mathbf{N}_{n' \neq n}$ that were added to the other data matrices:

$$\mathbf{X}_n = \mathbf{A}_n \mathbf{S} + \mathbf{N}_n, \tag{1}$$

where $\mathbf{A}_n$ is a mixing matrix specific to subject $n$. The "shared sources" $\mathbf{S}$ might represent brain sources driven by the same stimulus similarly across different subjects. We are interested in finding these shared sources and suppressing the noise. Note that this model assumes that different subjects share the same source *time course*, but not necessarily the same mixing matrix.

It is useful to reason in terms of the *subspace* spanned by a set of source or electrode time series (consisting of all their linear combinations). In the following we use the term 'source' to designate a particular process in the brain or elsewhere (and by extension the time series it produces), 'signal' to designate sources of interest, and 'temporal pattern' or 'component' to designate more generally any one-dimensional time series within the subspace spanned by the data.

**A simple CCA formulation.** Consider two data matrices, $\mathbf{X}_1$ and $\mathbf{X}_2$ of size $T \times d$ where $T$ is time and $d$ the number of channels (the same for both matrices for simplicity). All data are assumed to have zero mean. Each matrix is spatially whitened by applying principal component analysis (PCA) and scaling each principal component (PC) to unit norm to obtain whitened matrices $\overline{\mathbf{X}}_1$ and $\overline{\mathbf{X}}_2$. Whitened data are then concatenated and submitted to a new PCA to obtain a matrix $\mathbf{Y} = [\mathbf{X}_1, \mathbf{X}_2]\mathbf{V}$ of size $T \times 2d$, where $\mathbf{V}$ combines the whitening and second PCA matrices (Fig. 1 left). The submatrices $\mathbf{V}_1$ and $\mathbf{V}_2$ formed of the first and last $d$ rows of $\mathbf{V}$ define transforms applicable to each data matrix:

$$\mathbf{Y}_1 = \mathbf{X}_1 \mathbf{V}_1, \mathbf{Y}_2 = \mathbf{X}_2 \mathbf{V}_2, \tag{2}$$

with $\mathbf{Y} = \mathbf{Y}_1 + \mathbf{Y}_2$ (Fig. 1 center).

The outcome of this analysis is equivalent to standard CCA, as explained in the Discussion, the first $d$ columns of $\mathbf{Y}_1$ and $\mathbf{Y}_2$ forming canonical pairs (within a scaling factor). Indeed, rotating $\overline{\mathbf{X}}_1$ and $\overline{\mathbf{X}}_2$ to maximize the correlation of the resulting $\mathbf{Y}_1$ and $\mathbf{Y}_2$, as required by the CCA objective, is equivalent to rotating with the goal of maximizing the norm of their sum, $\mathbf{Y}_1 + \mathbf{Y}_2$, as achieved by the second PCA (Fig. 1 right). The appeal of this formulation is that it is easily extendable to multiple data matrices.

**A simple MCCA formulation.** Consider $N$ data matrices $\mathbf{X}_n$ each of size $T \times d$ with zero mean. Each data matrix is spatially whitened by applying PCA and scaling all PCs to unit norm to obtain whitened matrices $\overline{\mathbf{X}}_n$. Whitened data are then concatenated along the component dimension and submitted to a second PCA to obtain a matrix $\mathbf{Y} = [\mathbf{X}_1 \dots \mathbf{X}_N]\mathbf{V}$ of size $T \times D$, $D = Nd$, where $\mathbf{V}$ combines the whitening and second PCA matrices (Fig. 2 left). The submatrices $\mathbf{V}_n$ of size $d \times D$ formed by extracting successive $d$-row blocks of $\mathbf{V}$ define transforms applicable to each data matrix:

$$\mathbf{Y}_n = \mathbf{X}_n \mathbf{V}_n, \tag{3}$$

with $\mathbf{Y} = \sum_n \mathbf{Y}_n$ (Fig. 2, right). If data matrices have different numbers of channels $d_n$, then $\mathbf{V}_n$ has size $d_n \times D$ where $D = \sum_n d_n$. We call the columns of $\mathbf{Y}_n$ *canonical correlates* (CCs) by analogy with CCA, and those of $\mathbf{Y}$ *summary components* (SC). Each SC is the sum over data sets of the CCs of same rank. Columns of $\mathbf{Y}$ are mutually orthogonal by virtue of the final PCA, but the same is usually not true for the $D > d$ columns of $\mathbf{Y}_n$ which form an *overcomplete basis* of the temporal patterns spanned by $\mathbf{X}_n$. This formulation of MCCA is equivalent to the MAXVAR formulation of Kettenring (1971) as explained in the Discussion (Parra, 2018). The appeal of this formulation is that it is conceptually and computationally straightforward. PCs can be discarded from the initial PCAs, so as to control dimensionality and limit overfitting effects (next section).

The variances of the summary components (the columns of $\mathbf{Y}$) reflect the degree to which temporal patterns are shared between data matrices (Fig. 3) – the variance of each SC corresponding to the degree of correlation of each shared dimension found in the data. If the data matrices $\mathbf{X}_n$ share no components, the variances of all SCs are one (Fig. 3 a). If a component is shared by all $N$ data matrices, the norm of the first SC is $N$ (Fig. 3 d). In real data, spurious correlations typically cause the variance
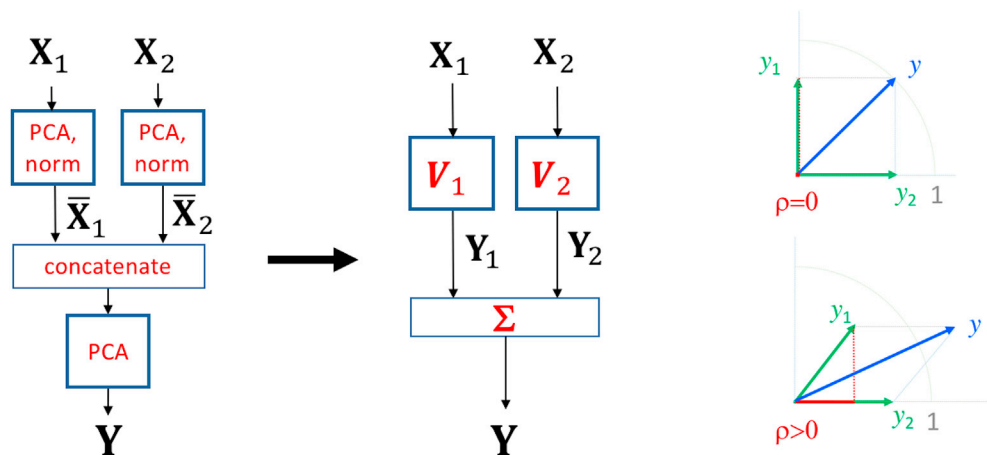
**Fig. 1.** Block diagram of the simple CCA formulation. Left: each data matrix is whitened by PCA followed by normalization. Normalized PCs from both data matrices are concatenated side by side and submitted to a final PCA. Center: transforms $\mathbf{V}_1$ and $\mathbf{V}_2$ (combining whitening and PCA matrices) transform each data set into canonical correlates (CC) $\mathbf{Y}_1 = \mathbf{X}_1\mathbf{V}_1$ and $\mathbf{Y}_2 = \mathbf{X}_2\mathbf{V}_2$. Their sum $\mathbf{Y} = \mathbf{Y}_1 + \mathbf{Y}_2$ is the matrix of summary components (SC). Right: rotating vectors $y_1$ and $y_2$ to maximize the norm of their sum is equivalent to maximizing their correlation coefficient $\rho$, symbolized by the projection of $y_1$ on $y_2$ (red line).
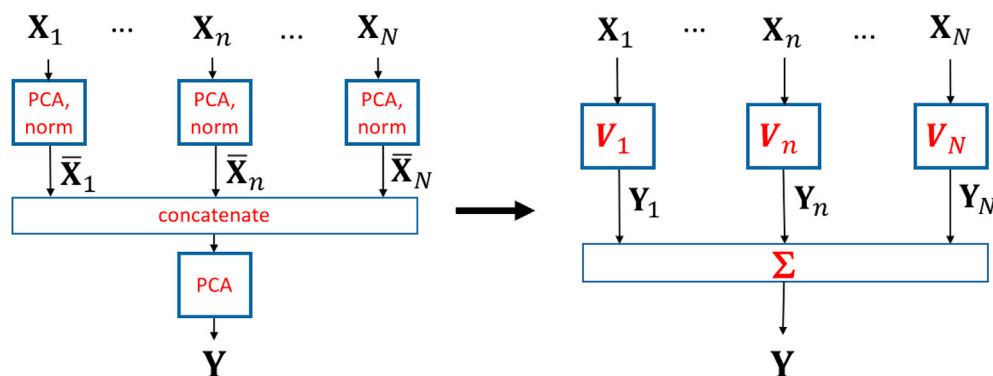


**Fig. 2.** Block diagram of the simple MCCA formulation. Left: each data matrix $\mathbf{X}_n$ is whitened by PCA followed by normalization. Normalized PCs from all data matrices are concatenated side by side and submitted to a final PCA. Right: the matrix $\mathbf{Y}$ of summary components (SC) can be expressed as the sum of individual transforms $\mathbf{Y}_n = \mathbf{X}_n\mathbf{V}_n$ (canonical correlates, CC).
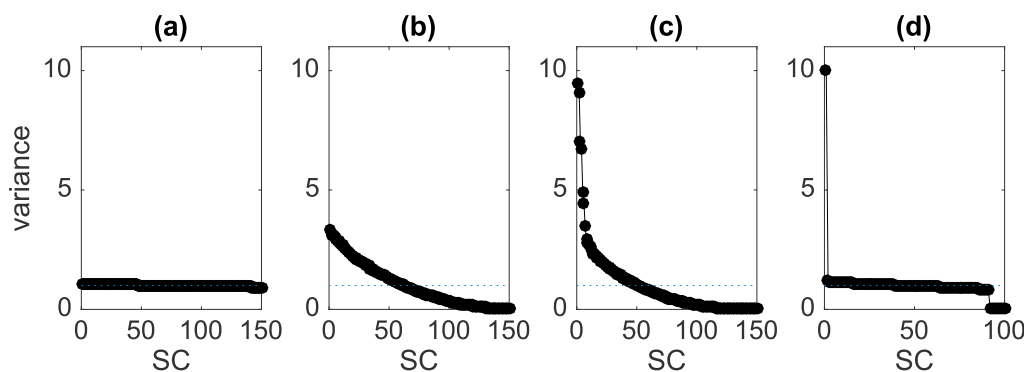


**Fig. 3.** SC variance as a function of order for MCCA analyses applied to 4 different types of dataset, each involving 10 data matrices. (a) Each data matrix consisted of an independent $10000 \times 15$ matrix of Gaussian white noise. In this case the SC variance profile is flat since there is no (or little) correlation between data matrices. (b) Each data matrix consisted of a $165 \times 15$ matrix of independent and uncorrelated Gaussian noise. In this case the SC variance profile is skewed, reflecting spurious numerical correlations between the statistically independent columns. (c) Each data matrix consisted of a $165 \times 15$ matrix of values derived from fMRI responses of 10 subjects in response to 165 sounds. Prior to MCCA the 6309 voxels were reduced to 15 channels using SVD (see description of Example 6 in the Methods). (d) Each data matrix consisted of a $10000 \times 10$ matrix of Gaussian white noise with an embedded sinusoid (Example 1, Fig. 4) that was the same in all data matrices. In the last two examples, only a small subset of the MCCA components reflect shared activity as evident by the low SC variance at higher MCCA orders.
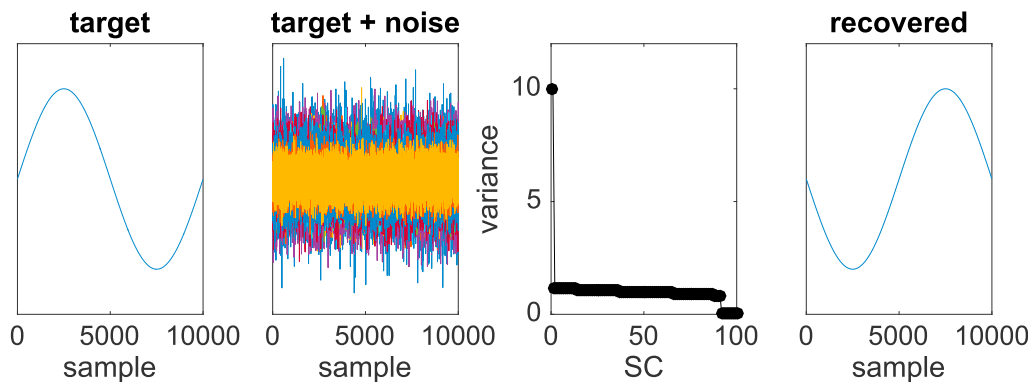
**Fig. 4.** Simulation with separable noise. Left: target signal. Next to left: target in noise at $SNR = 10^{-20}$. Next to right: variance of SCs as a function of order. The variance of the first SC is equal to 10 as target is perfectly shared across subjects and mixed in separable noise. Right: target recovered by MCCA (with arbitrary sign).

profile to follow a smooth curve (Fig. 3 b). Genuine shared activity often shows up as components with variance elevated relative to this smooth trend (Fig. 3 c). Elevated variance for the first components is accompanied by reduced variance for the last components (as in Fig. 3 d) corresponding to directions *least shared* between data matrices.

**Dimensionality reduction.** If one or more data matrices $\mathbf{X}_n$ are rank-deficient (for example artifacts were projected out during preprocessing) or include numerous low-variance dimensions, likely to be dominated by noise, it may be useful to reduce their dimensionality to $\dot{d} < d$. This can be done easily in the initial PCA by discarding the $d - \dot{d}$ PCs with smallest variance. The MCCA transform matrices $\mathbf{V}_n$ are then of size $d \times \dot{D}$, $\dot{D} = N\dot{d}$, and the CC and SC matrices of size $T \times \dot{D}$. Dimensionality reduction avoids computational issues with rank-deficient data, reduces computation and memory requirements, and serves as regularization to help avoid overfitting in later stages.

**Dealing with data matrices with more channels than samples.** CCA fails if the data matrices have fewer samples than channels ($T \leq d$), as is typically the case for fMRI or calcium imaging data for which there are many more voxels or pixels than observation samples (Asendorf, 2015). A simple solution is to replace each data matrix $\mathbf{X}_n$ (size $T \times d$) by a matrix $\dot{\mathbf{X}}_n$ of size $T \times \dot{T}$ with $\dot{T} < T$ columns that capture the principal temporal patterns spanned by $\mathbf{X}_n$. This can be done by applying singular value decomposition (SVD) to express the data as

$$\mathbf{X}_n = \mathbf{U}\mathbf{S}^\mathsf{t}\mathbf{V} \tag{4}$$

and setting $\dot{\mathbf{X}}_n = \dot{\mathbf{U}}$ where $\dot{\mathbf{U}}$ consists of the first $\dot{T}$ columns of $\mathbf{U}$. Since the $\dot{\mathbf{X}}_n$ have more samples than channels there is no obstacle to applying MCCA to them. This sequence of operations can be represented by a set of transform matrices $\mathbf{V}_n$ of size $d \times N\dot{T}$. Applying them to the data yields canonical correlate and summary matrices of size $T \times N\dot{T}$. Using this approach, it is straightforward to apply MCCA to datasets with a large number of "channels" such as data from calcium imaging or fMRI.

### 2.2. Applications of MCCA

**Quantifying correlation between *N* data matrices.** The variance of each column of $\mathbf{Y}$ indicates the degree to which a component is shared across data matrices. The value is 1 if the data matrices are perfectly uncorrelated, and $N$ if all data matrices include that component (Fig. 3). Components with variance close to zero are minimally correlated between data matrices. The profile of variances over SCs thus offers a measure of "sharedness" between data matrices (but see Caveats).

**Summarizing a set of data matrices.** The first few columns of $\mathbf{Y} = \sum_n \mathbf{Y}_n$ represent temporal patterns that capture most of the correlation across data matrices $\mathbf{X}_n$. They form a basis of the signal subspace that contains those shared temporal patterns.

**Denoising.** Each data matrix $\mathbf{X}_n$ may be denoised by projecting it to the overcomplete basis of CCs, selecting the first $\dot{D} < D$ components, and projecting back. We refer to this procedure as "denoising", as it can be used to attenuate components that are least shared across subjects, likely to be noise. This can be summarized by a denoising matrix $\mathbf{D}_n$ product of the first $\dot{D}$ columns of $\mathbf{V}_n$ by the first $\dot{D}$ rows of its pseudoinverse. The denoised data are then obtained as $\tilde{\mathbf{X}}_n = \mathbf{X}_n\mathbf{D}_n$.

**Dimension weighting.** Dimensionality reduction is often performed by applying PCA to a data matrix and truncating the PC series (Cunningham and Yu, 2014). However, this equates relevance to variance, which may not be appropriate because noise components can have high variance and useful targets small variance. MCCA can be used to weight dimensions according to their *consistency across data matrices*, which may be a better criterion than variance.

**Outlier detection.** Temporally-local glitches and artifacts may interfere with data interpretation and analysis. Analysis algorithms based on least-squares are particularly sensitive to high-amplitude artifacts. MCCA can be used to derive a cross-subject 'consensus' response, so that individual subject's data points that deviate greatly from the consensus can be flagged as outliers and excluded from analysis.

### 2.3. Details of the evaluation examples

The methods are evaluated using six datasets, including synthetic data, EEG, and fMRI.

**Example 1 - sinusoidal target in separable noise.** Synthetic data for this example consisted of 10 data matrices, each of dimensions 10000 samples × 10 channels. Each was obtained by multiplying 9 Gaussian noise time series (independent and uncorrelated) by a 9 × 10 mixing matrix with random Gaussian coefficients. To this background of noise was added a "target" consisting of a sinusoidal time series (Fig. 4, left) multiplied by a 1 × 10 mixing matrix with random coefficients. The target was the same for all data matrices, but the mixing matrices differed, as did the noise matrices. The SNR was set to $10^{-20}$, i.e. a very unfavorable SNR. The noise is of rank 9 and the signal of rank 1, so signal and noise are in principle linearly separable.

**Example 2 - sinusoidal target in non-separable noise.** Synthetic data for this example consisted of 10 matrices of dimensions 10000 samples × 10 channels, each obtained by multiplying 10 Gaussian noise time series (independent and uncorrelated) by a 10 × 10 mixing matrix with random coefficients. To this background was added a sinusoidal target as in the previous example, with SNR varied as a parameter. The noise here is of rank 10 so signal and noise are *not* linearly separable.

**Example 3 - sinusoidal target in EEG noise.** Data for this example used EEG to simulate realistic background noise. EEG data were recorded during approximately 20 min from one subject in the absence of any task, from 40 electrodes (32 standard positions plus additional electrodes on forehead and temple) at 2048 Hz sampling rate with a BioSemi system. A

robust polynomial detrending routine (de Cheveigné and Arzounian, 2018) was used to remove slow drifts. Ten "data matrices" were produced by selecting 3-s intervals of EEG data with random offsets, removing their means, and adding a target consisting of 4 cycles of a 4 Hz sinusoid multiplied by a $1 \times 40$ mixing matrix with random Gaussian coefficients, renewed for each data matrix. The analysis was repeated for each value of the SNR parameter.

**Example 4 - EEG response to tones.** Data for this example were borrowed from a study on auditory attention (Southwell et al., 2017). EEG data were recorded at a sampling rate of 128 Hz using a 64-channel EEG system in response to 120 repetitions of a 1 kHz tone pip with interstimulus interval (ISI) randomized between 750 and 1550 ms. They were recorded for the purpose of locating electrodes responsive to sound, prior to the main experiment in that study. Data from the first 10 subjects (to be consistent with the simulations) were selected and detrended using a robust detrending routine, bad channels were interpolated using spherical interpolation (EEGLAB), and the data were filtered between 2 and 45 Hz. A peristimulus epoch of duration 1.2 s (starting 0.2 s prestimulus) was defined for each trial, and the corresponding data were extracted as a 3D matrix of dimensions time $\times$ channel $\times$ trial. For each channel, the 0.2 s prestimulus waveform was averaged over trials and subtracted from that channel's waveform (a.k.a. "baseline correction"). PCA was applied to the data of each subject, the first 30 PCs (arbitrary number) were retained and the remainder discarded.

Two analyses were performed on these data to try to extract the cortical response to the 1 kHz tone from the background EEG noise. In the first, repetition over trials was exploited to design a spatial filter for each subject using the joint diagonalization algorithm (JD) that maximizes the ratio of trial-averaged variance to total variance (de Cheveigné and Simon, 2008; de Cheveigné and Parra, 2014). Combining the initial PCA matrices (size $64 \times 30$) and subsequent JD matrices (size $30 \times 30$) led to 10 analysis matrices of size $64 \times 30$, one for each subject. In the second analysis, MCCA was applied, using 30 PCs from each subject in the first PCA. Multiplying the initial PCA matrices (size $64 \times 30$) by the subsequent subject-specific MCCA matrices (size $30 \times 300$) led to 10 subject-specific analysis matrices of size $64 \times 300$.

For each subject, the first column of the JD analysis matrix defines the best linear combination of channels to maximize repeat-reliability across trials, while the first column of the MCCA analysis matrix defines the best linear combination of channels of that subject to maximize correlation with the other subjects.

**Example 5 - EEG response to speech.** Data for this example were taken from a study on auditory cortical responses to natural speech (Di Liberto et al., 2015). The same data were also used in a recent study on the application of CCA to speech/EEG decoding (de Cheveigné et al., 2018). We borrowed the decoding methods and evaluation metrics from the second study, with the purpose of evaluating the benefit of introducing a denoising stage based on MCCA before the speech/EEG decoding stage.

In brief, EEG data were recorded from 8 subjects using a 128-channel BioSemi system with standard electrode layout, at 512 Hz sampling rate. Each subject listened to 32 speech excerpts, each of duration 155 s, from an audio book, presented diotically via headphones, for a total of approximately 1.4 h. The database included both the audio stimuli and the EEG responses. Further details about the stimulus and recording are available in Di Liberto et al. (2015). The EEG were preprocessed (downsampling to 64 Hz, detrending, artifact removal), and the stimulus temporal envelope calculated as described in de Cheveigné et al. (2018).

A stimulus-response model (de Cheveigné et al., 2018; Dmochowski et al., 2017) was evaluated according to several metrics: correlation, d-prime, and percent-correct classification scores for a match vs mismatch classification task. The classification task consisted in deciding whether a segment of EEG matched the stimulus that produced it (match) or some unrelated stimulus segment (mismatch). The duration of the segment was varied as a parameter from 1 to 64 s. The decoding model used CCA to relate the stimulus to the EEG response, producing multiple stimulus-response CC pairs that were used for discrimination. Further details of the decoding model, classification task, and metrics can be found in de Cheveigné et al. (2018). Here, we are only interested in knowing if scores for decoding are improved by introducing a stage of EEG denoising based on MCCA. Statistical significance is evaluated using a Wilcoxon rank-sum test.

For this denoising, the EEG data of each subject were submitted to MCCA, keeping 40 PCs in the first PCA, resulting in a $128 \times 320$ analysis matrix for each subject. The first 110 columns of this matrix were multiplied by the first 110 rows of its pseudoinverse to yield a $128 \times 128$ subject-specific denoising matrix. This has the effect of attenuating activity that is *least* correlated with the other subjects. The numbers of components to retain (40 PCs, 110 CCs) were chosen by trial and error.

**Example 6 - fMRI response to natural sounds.** Data for this example were taken from a study that measured fMRI responses to natural sounds (Norman-Haignere et al., 2015). Responses were gathered from 10 subjects to each of 165 sounds belonging to 11 categories including speech, music, animal vocalizations, and others. For each subject, the recording session was repeated either twice or 3 times. Data were projected to a common anatomical space, see Norman-Haignere et al. (2015) for further details. For the present analysis, data for each subject were averaged over repeats and organized as a matrix $\mathbf{X}_n$ of 165 sounds $\times$ 6309 voxels (voxels from both hemispheres were used, and voxels outside a subject-specific region of interest that included primary and secondary auditory cortex were set to zero). In this analysis we are interested in finding particular profiles of response over sounds (for example speech vs non-speech, or music vs non-music) and also the brain areas within the region of interest associated with such profiles in each subject.

As there are more "channels" (voxels) than samples ($T < d$), an SVD was used as described in the Methods and the first 10 dimensions were used for MCCA. The columns of $\dot{\mathbf{X}}_n$ are white so the first PCA is unnecessary. Matrices $\dot{\mathbf{X}}_n$ were concatenated and subjected to the second-step PCA of the MCCA algorithm, and the 15 first columns (arbitrary number) of the SC matrix were selected as a basis spanning the profiles over sounds that were most similar across subjects.

To find profiles specific to particular sound categories (e.g. speech, music, etc.), JD (de Cheveigné and Parra, 2014) was used to find a linear transform applicable to the 15-column basis to maximize the variance over the selected category, relative to the other categories. This can be seen as a rotation of the basis so as to isolate activity specific to processing of that sound category. This $165 \times 1$ activation profile was then cross-correlated with the $165 \times 6309$ matrix of fMRI response data of each subject to find the topography specific to that subject (Haufe et al., 2014).

## 3. Results

The MCCA method is evaluated first with synthetic data to get an understanding of its basic properties and capabilities, and then with real EEG and fMRI data to see whether these extend to situations of practical use. Additional examples and scripts are available at http://audition.ens.fr/adc/NoiseTools/src/NoiseTools/EXAMPLES/MCCA_EXAMPLE_SCRIPTS/.

### 3.1. Synthetic data

**Example 1 - sinusoidal target in separable noise.** The data consist of 10 matrices including a sinusoidal target (Fig. 4, left) common to all data matrices, with added noise distinct across matrices (see Methods). At the unfavorable SNR of $10^{-20}$ the target is not visible in the raw data matrices (Fig. 4 center), and it cannot be extracted by averaging because of the extremely low SNR and the fact that the mixing coefficients are of random sign. Since the data are separable (the rank of the noise is only 9), the target *can* be recovered by applying the appropriate demixing matrix

(inverse of the mixing matrix), however that matrix is unknown.

MCCA applied to the dataset produced projection matrices $\mathbf{V}_n$ that recover the target from $\mathbf{X}_n$ (Fig. 4 right). This benefit is similar to that of methods that leverage multiple repetitions to blindly discover spatial filters to improve SNR (de Cheveigné and Simon, 2008; de Cheveigné and Parra, 2014), but instead of repetitions, MCCA leverages the fact that the same target is mixed into multiple data matrices. To summarize, MCCA can reveal a target common across data matrices despite an extremely unfavorable SNR.

**Example 2 - sinusoidal target in non-separable noise.** Data are the same as in the previous example, except that the noise is full rank (10 independent time series mixed in 10 channels) so the target is no longer linearly separable, and one cannot expect to recover the target perfectly, especially at extremely low SNRs. Nonetheless, at a moderately unfavorable SNR ($10^{-20}$ in power) MCCA can recover an estimate of the target that is noisy (Fig. 5 center) but much cleaner than the raw data (not shown). Fig. 5 (right) shows the proportion of residual noise in the signal recovered by MCCA as a function of SNR, together with the same proportion for the best raw channel. MCCA provides a clear benefit over a range of SNRs. Two factors can contribute to failure: non-separability per se, and the fact that the algorithm fails to find the ideal demixing matrix. To get an idea of their contributions, Fig. 5 (right) also shows the proportion of residual noise for the ideal demixing matrix (yellow). The MCCA-derived matrix performs only slightly less well than the ideal matrix, and both much better than the best channel of the raw data. To summarize, MCCA is also of use when the data are not separable.

**Example 3 - sinusoidal target in real EEG noise.** EEG background noise differs from the white Gaussian noise that was used in the previous simulations in several ways: it usually has full rank (in particular because of electrode-specific noise), but the variance is unequally distributed across dimensions. It is also temporally structured, with strong temporal correlation and an overall low-pass spectrum. The first component recovered by MCCA is plotted in Fig. 6 (right) for several values of SNR. For SNRs of 0.1 or better the target is almost perfectly recovered. At SNR $= 0.03$ the recovered waveform is somewhat noisy, and at SNR $= 0.01$ or below the target is lost. For comparison Fig. 6 (left) shows the time course of a raw data channel (the channel that showed the largest correlation with the target). For SNR $= 10$ the target waveform is obvious in the raw data, and for SNR $= 1$ it can be guessed, but for smaller values of SNR it is lost in the EEG noise. Comparing Fig. 6 left and right, there is a range of SNRs (roughly 0.03 to 1) for which MCCA provides a clear benefit. Below SNR $= 0.03$ the algorithm switched to some other component within the data (Fig. 6 right, lowest trace) that happened to be similar across data matrices because of random correlations.

### 3.2. Real data

**Example 4 - EEG response to tones.** In this example, contrary to the previous one, the target is not known. However, since the data were collected in response to multiple repeats *and* for multiple subjects, we can apply two different methods (JD and MCCA) and see if they corroborate each other. JD finds a linear transform that optimizes SNR assuming that the signal repeats over trials. Fig. 7 (top) shows the result of applying the JD analysis to the data of one subject. In the plot on the top left, the blue line shows the mean over repeats of the first component, and the gray band shows $\pm$ 2 SD of a bootstrap resampling of this mean. On the top right is the topography associated with this component computed as the map of cross-correlation coefficients between the component and each channel (Parra et al., 2005; Haufe et al., 2014). MCCA can similarly be used to design a subject-specific spatial filter that improves SNR. The plots on the bottom of Fig. 7 show the result of applying the subject-specific matrix derived from the MCCA analysis for the same subject. Despite the different criteria used by the two analyses (consistency over trials for JD, consistency between subjects for MCCA) the temporal patterns are remarkably similar. The correlation between the two ranges over subjects from 0.6 to 0.996, with a mean of 0.91. To summarize, it appears that MCCA can exploit between-subject consistency to find a spatial filter that is as effective as that found by JD using between-trial consistency. This is useful for experimental designs that do not involve repeated trials.

The subject-specific MCCA analysis matrices ($\mathbf{V}_n$) transform each subject's data ($\mathbf{X}_n$) into CCs ($\mathbf{Y}_n$) that are well correlated across subjects so that it makes sense to average them across subjects and interpret the SCs ($\mathbf{Y}$) as reflecting shared activity. Fig. 8 top left shows the trial- and subject-averaged time course of the first SC, which can be interpreted as our best estimate of stimulus-evoked activity common to all subjects. It benefits from several stages of enhancement: (a) spatial filtering within each subject, (b) averaging over trials, (c) averaging across subjects. Also shown in Fig. 8 are the ten subject-specific topographies associated with this component. Despite some differences, topographies are quite similar across most subjects except S1. The bottom left plot shows the maximum over electrodes of the correlation coefficient between the first SC and each electrode (trial-averaged). Correlation coefficients are relatively high except for Subject 1.

**Example 5 - EEG response to speech.** For stimuli presented once only, such as speech or music, one cannot use repetition to distinguish the brain response from the noise. Instead, systems identification techniques (Lalor et al., 2009; Holdgraf et al., 2017; Crosse et al., 2016) have been used to fit an encoding model to estimate the part of brain response that is driven by the stimulus. The part of the response that fits the model can be taken as the "true" response, and the rest discarded as noise.

However, this requires a prior choice of a stimulus representation (e.g. envelope or spectrogram) that can be linearly related to the brain signals. An imperfect fit could signify either unpredictible noise in the brain signal, or inadequacy of the stimulus representation. MCCA can help in this context by providing a less noisy "consensus" brain response that can be used to evaluate the adequacy of the stimulus features.

Here, EEG were recorded in response to continuous speech (see Methods), and a model was fit to stimulus and response to capture their correlation (de Cheveigné et al., 2018; Dmochowski et al., 2017). The
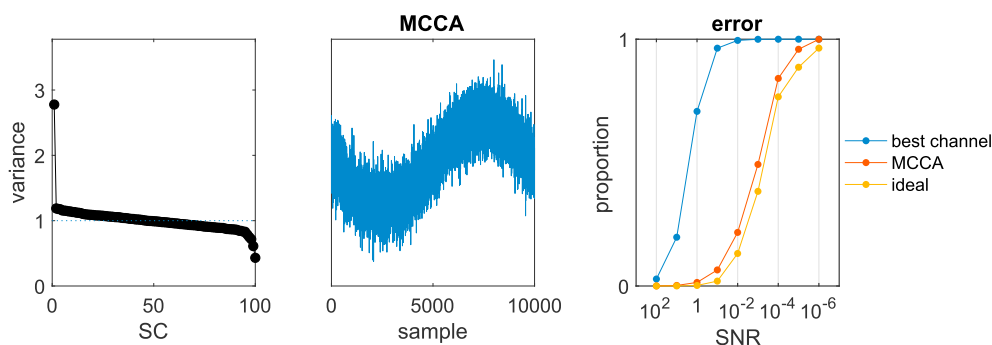


**Fig. 5.** Simulation with inseparable noise. Left: variance of SCs as a function of their order at SNR $= 10^{-2}$. Center: target signal recovered from mixture at SNR $= 10^{-2}$. Right: proportion of residual noise power as a function of SNR for the raw data (blue), first SC (red) or ideal demixing matrix (yellow).
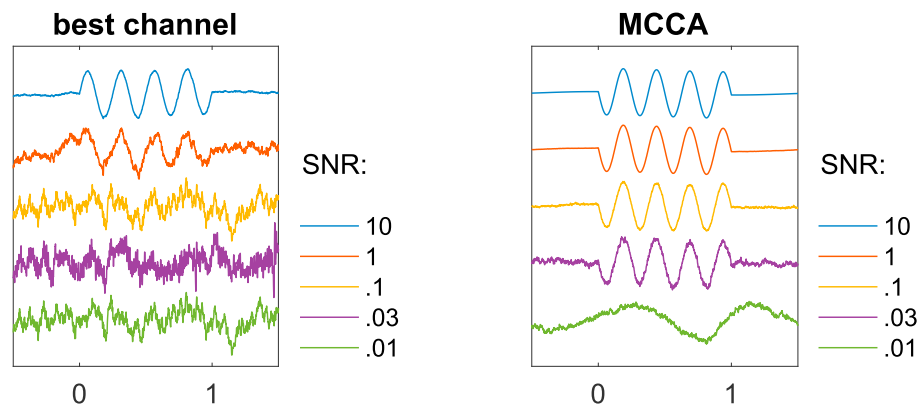
**Fig. 6.** Simulation with EEG noise. Left: time course of the best raw data channel for several values of SNR. Right: time course of the first MCCA component for several values of SNR.
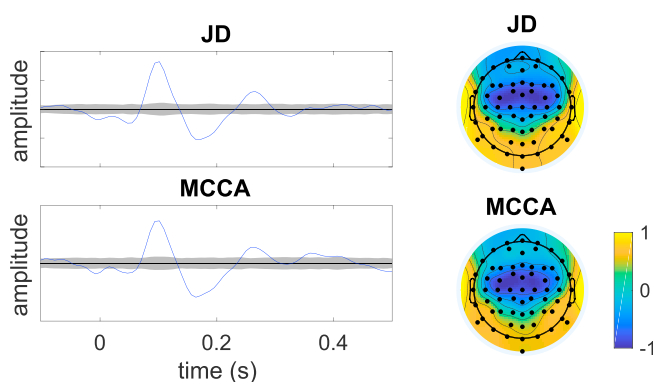


**Fig. 7.** Comparison between JD solution (within-subject repeat-reliability) and MCCA solution (between-subject similarity) for one subject among ten. Data were in response to repeated tones. Left: average over trials (blue) and ± 2 SD of a bootstrap resampling (gray) of the first JD component, which maximizes reliability across trials (top), or first subject-specific CC (bottom). Right: associated topographies (correlation between trial-averaged component and trial-averaged electrode waveforms).

model used CCA to form pairs of maximally-correlated linear transforms of the audio stimulus features and of the EEG respectively (audio-EEG CCs). Note that this usage of CCA is unrelated to our usage of MCCA to merge data across subjects. The quality of that model was evaluated using a match vs mismatch classification task (see Methods). We compute *correlation*, *d-prime* and *percent correct* classification scores to evaluate the benefit of inserting a stage of MCCA-based denoising within the EEG preprocessing pipeline, before classification.

Fig. 9 (a) shows the correlation between the first audio-EEG CC pair (thick blue line) and subsequent pairs (thin lines), with and without MCCA-based denoising, for one subject. To the extent that correlation is limited in part by EEG noise, the higher scores on the right suggest that denoising was effective. The d-prime metric measures the degree of separation between distributions of correlation scores for matched and mismatched segments. Fig. 9 (b) shows the d-prime metric for the first pair (thick blue) and subsequent pairs (thin lines), with and without MCCA-based denoising for segments of duration 64 s. The dotted line shows the d-prime metric for the multivariate distributions of audio-EEG CC pairs. The larger d-prime scores with MCCA-based denoising suggest that it can effectively contribute to improved discrimination. Fig. 9 (c) shows classification scores as a function of segment duration with (red) and without (black) MCCA-based denoising. The higher scores with MCCA-based denoising show its benefit for this task. Fig. 9 (d) shows that a similar benefit is found in all subjects (p = 0.00016, rank sum test). The thick lines are scores for a duration of 16 s, whereas the thin lines are for

segments of 2 s (lowest lines) or 64 s (highest lines). To summarize, MCCA is of benefit as a denoising tool for EEG responses to speech.

The match-vs-mismatch task is used here as a metric to evaluate the quality of the stimulus-response decoding model, that is involved also in the more complex Auditory Attention Detection (AAD) task that aims to determine which of two concurrent voices is the focus of a listener's attention (Ding and Simon, 2012; Fuglsang et al., 2017; Lalor et al., 2009; Khalighinejad et al., 2017; Koskinen and Seppä, 2014; Martin et al., 2014; Mesgarani and Chang, 2012; Mirkovic et al., 2015; O'Sullivan et al., 2014; Tiitinen et al., 2012; Zion Golumbic et al., 2013). That task is actively studied for its potential use for the "cognitive control" of an external device such as a hearing aid.

**Example 6 - fMRI responses to natural sounds** Data were taken from a study that investigated fMRI responses to natural sounds (Norman-Haignere et al., 2015), in which 10 subjects listened to a set of 165 sounds belonging to 11 different classes. MCCA was applied to find evidence of selectivity to sound common across subjects as explained in the Methods. In brief, the $165 \times 6309$ matrix of voxel activations for each subject was reduced to a $165 \times 12$ matrix using SVD, and the reduced matrices concatenated and submitted to PCA to obtain a $165 \times 120$ matrix of SCs. Their variances are plotted in Fig. 10 (top left). The first 10 SCs were then subjected to a JD analysis to enhance the contrast between musical sounds (classes 'Music' + 'VocalMusic') and other sounds as explained in the Methods.

The profile of activation over sounds of the first JD component is plotted in Fig. 10 (top right), with sounds ordered by class and coded as different colors. Activations of the first two classes ('Music' + 'VocalMusic') are clearly distinct from those of the other classes. The corresponding topography of activation over voxels for each subject can be calculated by cross-correlating this component with the profile of activation over sounds of each voxel. Topographies for the left hemisphere for all subjects are plotted in Fig. 10 (bottom). To a first approximation, topographies are consistent in that a dorso-frontal concentration of activity is found in most subjects. To a second approximation, each topography includes additional regions, suggesting a wider area of activation that is more subject-specific. Such subject-specific details would be smoothed out by averaging over subjects. A similar analysis to enhance speech-specific activation revealed spatial patterns with more ventral topographies (not shown). Norman-Haignere et al. (2015) also found distinct music-related and speech-related components using an ICA-related technique.

The benefit of MCCA here can be interpreted in terms of dimensionality weighting, based here on *consistency across subjects* rather than variance as with PCA. Weighting based on MCCA, followed by PCA and truncation, allowed the final JD analysis to be performed on a matrix of size $165 \times 12 \times 10$ rather than $165 \times 6309 \times 10$, making it more effective by reducing overfitting. Using MCCA rather than just PCA ensures that these 12 dimensions are dominated by activity similar across
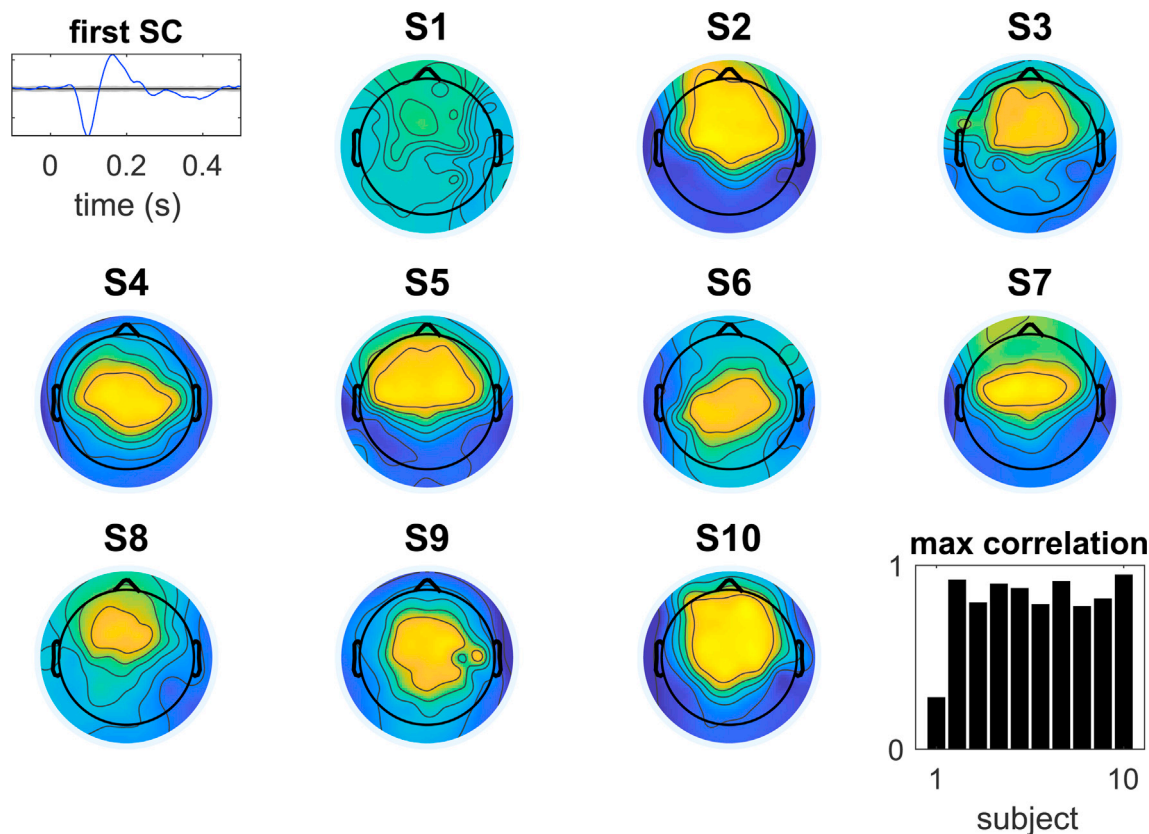
**Fig. 8.** MCCA analysis of tone response, summary over 10 subjects. Top left: trial-averaged time course of the first SC. Bottom right: maximum absolute value of correlation between that component and each electrode, for each subject. Other panels: topography of correlation values (of the SC with each electrode) for each subject (the color code is the same as in Fig. 7, bottom).
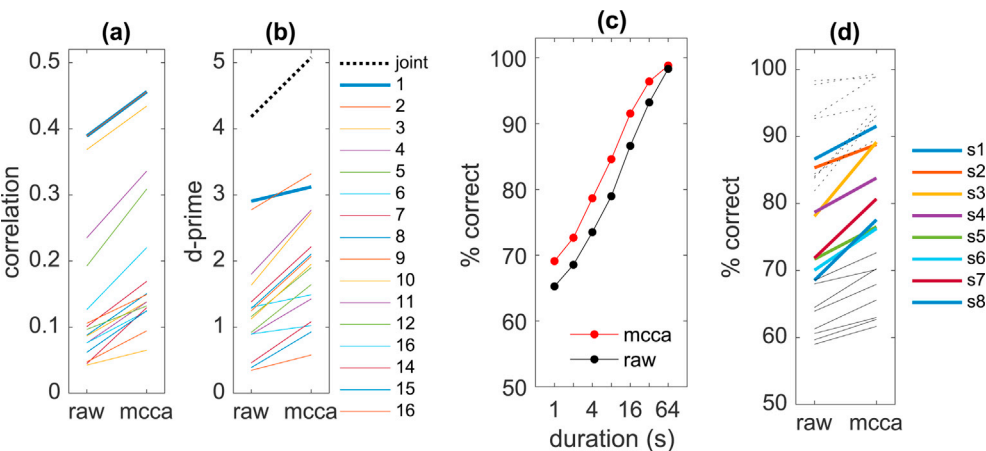


**Fig. 9.** Speech-EEG decoding. (a) Correlation coefficient for the audio-EEG first CC pair (thick blue line) and subsequent pairs (thin lines) for a CCA model, with and without MCCA-based denoising. (b) d-prime metric for a classification task for the first audio-EEG CC pair (thick blue line) and subsequent pairs (thin lines), with and without MCCA-based denoising. The dotted line is for multivariate classification based on all CC pairs. (c) Percentage correct classification as a function of segment duration, with and without MCCA-based denoising. (d) Percentage correct for segments of duration of 16 s (thick lines) for 8 subjects, with and without MCCA-based denoising. Thin lines are scores for 64 s (uppermost) or 2 s (lowermost).

subjects, likely to be relevant because all subjects heard the same stimuli. This example demonstrates that MCCA can be applied also to data with more channels (pixels or voxels) than data points.

### 4. Discussion

MCCA finds a linear transform applicable to each data matrix within a data set to align them to common coordinates and reveal shared temporal patterns. It can be used in several ways: as a *denoising* tool applicable to an individual data matrix, as a tool for *dimensionality weighting*, as a tool to *align* data matrices within a common space to allow comparisons, or as a tool to *summarize* data and reveal temporal patterns that are general

across data matrices. As formulated here, MCCA is easy to understand, straightforward to apply, and computationally cheap. Care is nonetheless required when applying it, in particular to avoid phenomena such as overfitting.

**What is new?** As reviewed in the Introduction, several versions of MCCA have been proposed in the literature and applied to the analysis of brain data. The contributions of this paper are the following. First, the formulation as a cascade of PCA, normalization, concatenation, and PCA offers an intuitive explanation that may help practitioners gain insight into this method. Past formulations are hard to follow for the non-mathematically inclined, and their sheer number is bewildering. We used a similar 2-step formulation in a recent tutorial on joint
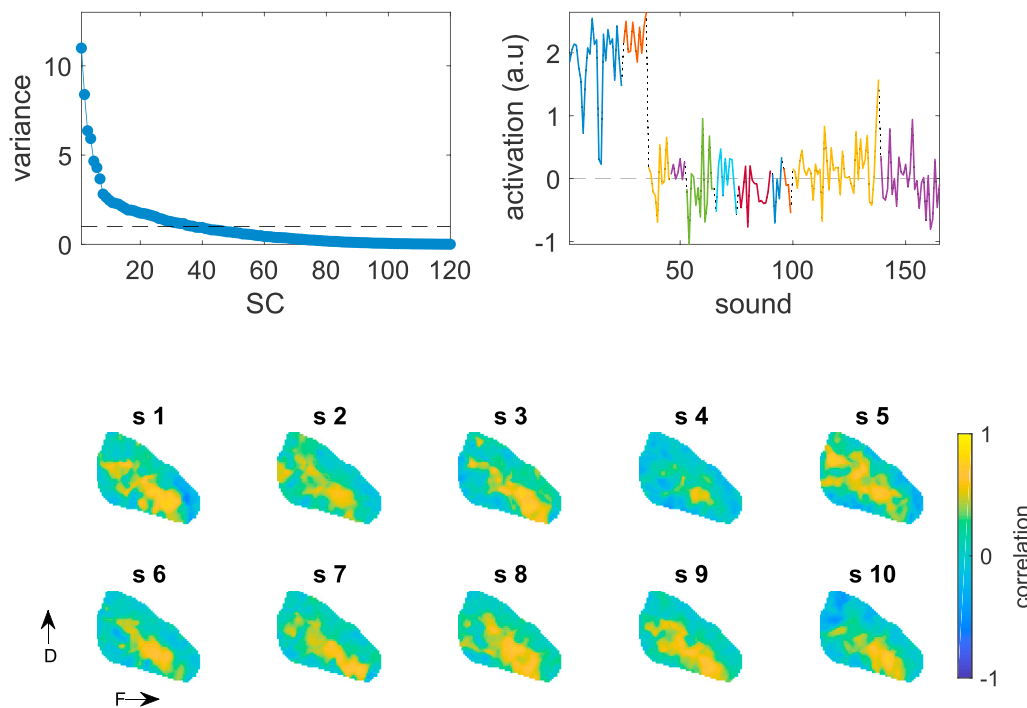
**Fig. 10.** MCCA of fMRI responses to natural sounds. Top left: SC variance as a function of order. Top right: activation as a function of sound of a component selective for music obtained by applying JD to the first 15 SCs (see text). Each color represents a different sound category; the first two categories are 'music' and 'vocal music'. Bottom: topographies of correlation between the music-selective JD component and the profile of response over sound of each voxel of the left hemisphere, for each subject. D: dorsal, F: frontal.

decorrelation (de Cheveigné and Parra, 2014), and we hope that the present paper too will have tutorial value. Second, our usage of MCCA as a denoising tool, to attenuate noise within individual subjects based on across-subject consistency by projection on the overcomplete basis of its SCs, seems to be new. Third, we provide tutorial examples that may encourage researchers to put MCCA to work for a wider range of tasks, including denoising, outlier detection, summarization, and cross-subject statistics.

**How does it work?** The effect of the processing steps is schematized in Fig. 11. Multiple data matrices contain the same source component **S**, illustrated here as a color gradient, mixed into two 2-dimensional data matrices (Fig. 11 a). Each point represents a sample in time (row of the data matrix) and the two axes represent two channels (columns of the data matrix). The color could represent a hidden sensory response that is similar across two subjects. The initial PCAs sphere each data matrix (b), so that the cloud of points is free to rotate in any direction. However, concatenating the sphered data matrices creates a cloud (in a 4-dimensional space) that is not spherical because of the shared component correlation along some direction in 4-D space (projected to 2D in panel (c)). The second PCA finds this direction of correlation between the data matrices and aligns it with the vertical axis (d), in the process transforming each data matrix so that it is optimally aligned with the other (e).

**Relation with other formulations of CCA and MCCA** As explained by Parra (2018), the aim of MCCA is to find projection vectors $\mathbf{v}_n$

applicable to $\mathbf{X}_n$ that maximize the ratio of between-set to within-set covariance:

$$\rho = \frac{1}{N-1} \frac{r_B}{r_W} \tag{5}$$

with:

$$r_B = \sum_n \sum_{n' \neq n} {}^{t}\mathbf{v}_n \mathbf{R}_{nn'} \mathbf{v}_{n'}$$

$$r_W = \sum_n {}^{t}\mathbf{v}_n \mathbf{R}_{nn} \mathbf{v}_n.$$

where $\mathbf{R}_{nn} = {}^{t}\mathbf{X}_n \mathbf{X}_n$ and $\mathbf{R}_{nn'} = {}^{t}\mathbf{X}_n \mathbf{X}_{n'}$ are covariance and cross-covariance matrices of the data. The divisor $1 - N$ ensures that $\rho$ scales between 0 and 1. Setting to zero the derivative of $\rho$ with respect to $\mathbf{v}_n$, the solution is obtained by solving the equation

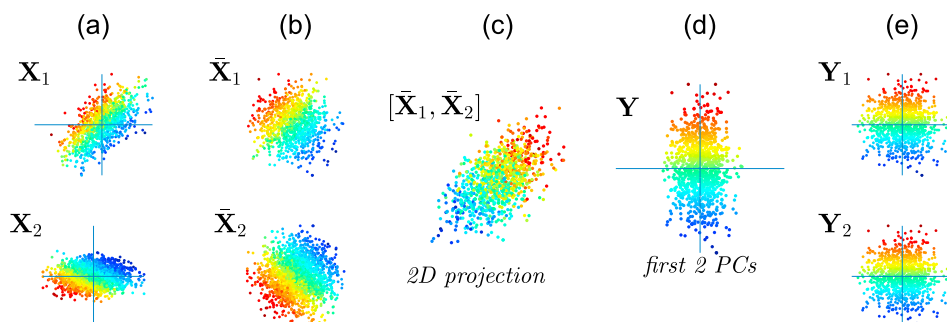$$\mathbf{R}\mathbf{v} = \mathbf{D}\mathbf{v}\lambda, \tag{6}$$

with



**Fig. 11.** Principle of MCCA. (a) Several data matrices share a common component (coded as color) but its orientation and nature are unknown. (b) Whitening makes the data matrices free to rotate. (c) Concatenation creates a cloud in 4D space (projected here to 2D) with a direction of greater correlation/variance due to the shared component. (d) The second PCA aligns this direction with the axes. (e) In the process, the whitened data matrices are rotated such that shared dimensions are maximally aligned.

$$\mathbf{R} = \begin{bmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} & \cdots & \mathbf{R}_{1N} \\ \mathbf{R}_{21} & \mathbf{R}_{22} & \cdots & \mathbf{R}_{1N} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{R}_{N1} & \mathbf{R}_{N2} & \cdots & \mathbf{R}_{NN} \end{bmatrix}, \mathbf{D} = \begin{bmatrix} \mathbf{R}_{11} & 0 & \cdots & 0 \\ 0 & \mathbf{R}_{22} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \mathbf{R}_{NN} \end{bmatrix}, \tag{7}$$

where $\lambda = \rho/(N-1) + 1$. Now, first decompose $\mathbf{D} = \mathbf{U}\mathbf{\Lambda}^t\mathbf{U}$. Because $\mathbf{D}$ is the block-diagonal matrix of the covariances in each data set, this decomposition implies doing PCA on each data set separately, i.e whitening each data set. With this decomposition Eq. (6) can be rewritten as:

$$\begin{aligned} \mathbf{R}\mathbf{v} &= & \mathbf{U}\mathbf{\Lambda}^t\mathbf{U}\mathbf{v}\lambda \\ \mathbf{\Lambda}^{-1/2t}\mathbf{U}\mathbf{R}\mathbf{v} &= & \mathbf{\Lambda}^{1/2t}\mathbf{U}\mathbf{v}\lambda \\ [\mathbf{\Lambda}^{-1/2t}\mathbf{U}\mathbf{R}\mathbf{U}\mathbf{\Lambda}^{-1/2}][\mathbf{\Lambda}^{1/2t}\mathbf{U}\mathbf{v}] &= & [\mathbf{\Lambda}^{1/2t}\mathbf{U}\mathbf{v}]\lambda \\ \tilde{\mathbf{R}}\tilde{\mathbf{v}} &= & \tilde{\mathbf{v}}\lambda \end{aligned} \tag{8}$$

where $\tilde{\mathbf{R}} = \mathbf{\Lambda}^{-1/2t}\mathbf{U}\mathbf{R}\mathbf{U}\mathbf{\Lambda}^{-1/2}$ is the covariance of the whitened concatenated data. Equation (8) thus corresponds to performing PCA on the concatenated whitened data. In summary, the two-step PCA described in the Methods ('simple MCCA formulation') maximizes summed correlation between data sets. The solution for MCCA in terms of an eigenvalue problem (6) is equivalent to the MAXVAR solution described by Kettenring (1971) (see, e.g. Vía et al., 2005a,b). Note, however, that here we have derived it as an optimization of the summed correlations, a criterion known as SUMCOR. The ambiguity in terminology results from the fact that the solution to the MCCA problem does not only depend on the criterion being optimized, but also on the constraints placed on the solutions (Nielsen, 2002). Kettenring introduced these terms only to specify different optimality criteria, yet the same solution can be obtained under different criteria depending on the choice of constraint (Asendorf, 2015).

**MCCA vs CCA** MCCA is understood as a generalization of CCA but some differences are worth noting. For CCA the focus is usually on the CCs $\mathbf{Y}_n$ ($n = 1, 2$), whereas for MCCA it may also be on the SCs $\mathbf{Y}$. For standard CCA the projection matrices are restricted to $d$ (or $\min_n d_n$) columns for each data set, whereas for MCCA it may be useful to consider more than $d$ columns (as in Example 5). If the objective were to capture sources common to *all* data matrices, $d$ components might suffice, but to capture also sources shared by *several data matrices but not all*, more than $d$ columns may be required. For CCA the $d$ columns of $\mathbf{Y}_1$ are mutually uncorrelated as are those of $\mathbf{Y}_2$, whereas for MCCA the $D$ columns of $\mathbf{Y}_n$ are correlated between each other in general. Columns of their sum $\mathbf{Y}$ are uncorrelated, however.

The larger number ($D > d$) and non-orthogonality of the columns of $\mathbf{Y}_n$ might be disconcerting for the researcher familiar with CCA. The method may be modified such that $\mathbf{Y}_n$ is instead constituted of $d$ orthogonal columns. For this, MCCA is applied as above, for each $n$ the first column of $\mathbf{Y}_n$ is projected out of $\mathbf{X}_n$, and MCCA applied again. This deflationary procedure terminates after $d$ steps because the dimensionality of each data matrix is then exhausted. Smaller matrices with orthogonal columns might be convenient in certain situations, but as pointed out they might not capture all shared sources. The procedure described in the Methods is better in this respect.

**Individual vs group** Studies that gather data from multiple subjects are confronted with the question of how to analyze and summarize the complex multisubject data that result. Why observe multiple subjects? Why summarize them? The answers depend on whether we are interested by *general patterns* robust with respect to intersubject differences and noise, or by the detail of *individual patterns* within each subject. Adding more subjects gives more statistical power to establish the reality of the former, and more examples of the latter.

Individual patterns are fascinating, but costly to report and in some cases hard to distinguish from noise. For those reasons, many studies focus on finding effects for which inferences can be made at the population level, i.e. generalizable from the sample of subjects observed to a wider population of unseen subjects. A tension remains nonetheless between this goal, and that of doing justice to the individual variations that are known to exist and are potentially important for understanding

(Molenaar, 2004; Gates and Molenaar, 2012; Karch et al., 2015).

Given the goal of finding a subject-general pattern, gathering data from multiple subjects in response to the same stimulus serves several purposes. First, to counteract variability by increasing the number of observations, analogous to averaging over repeated trials. Second, to allow data-dependent analysis to improve SNR based on similarity between subjects, analogous to methods that improve SNR based on similarity between trials (de Cheveigné and Parra, 2014). Third, to make inferences at the population level via group-level statistical analysis.

The conventional strategy of calculating a "grand average", with corresponding channels or voxels of each subject being averaged together (Choi et al., 2013; Luck, 2005), is hampered by inter-subject differences in source-to-sensor mapping. The problem is mild for sources with broad topographies (as in Fig. 8), but for sources with more local spatial characteristics a mismatch between subjects might result in destructive summation. A similar problem affects measures of inter-subject correlation (ISC) applied directly to channels or voxels (Hasson et al., 2004), or to linear combinations that assume the same mixing vectors for all subjects (Dmochowski et al., 2012; Parra et al., 2018).

A simple expedient is to select, for each subject, a group of channels based on responses to a "localizer" stimulus or task, calculate a root mean square average waveform over those channels, and then average over subjects (e.g. Chait et al. (2010)). However, this packs the multidimensional cortical activity into a single time course from which it may be hard to infer the richer dynamics of cortical activity. Another approach is to apply inverse modeling to map the activity to a source space common across subjects (Litvak and Friston, 2008). However, this requires accurate anatomical information for each subject and is subject to the validity of the reconstruction models (Mahjoory et al., 2017), as well as between-subject variability in source positions and orientations (Lio and Boulinguez, 2016).

Data-driven methods such as MCCA are attractive in that they find a mapping between subjects based only on shared temporal aspects of the data. MCCA and related methods have been widely used for fMRI data (Li et al., 2009; Correa et al., 2010b; Hwang et al., 2012; Afshin-Pour et al., 2012, 2014; Karhunen et al., 2013; Haxby et al., 2011) and EEG/MEG (Lankinen et al., 2014; Sturm, 2016; Zhang et al., 2017). In contrast to MCCA, which finds variance dimensions that are similar across subjects with no attempt to ensure that they correspond to sources within the brain, ICA-based approaches attempt to isolate sources common across subjects based on criteria of statistical independence (Calhoun and Adali, 2012; Eichele et al., 2011; Huster et al., 2015; Chen et al., 2016; Madsen et al., 2016; Huster and Raud, 2018). Group ICA (GICA) as formulated by Eichele et al. (2011) can be seen as a concatenation of MCCA (as described here) with ICA. Isolating the MCCA step, as we do here, is useful conceptually and avoids the computational cost and assumptions associated with ICA. Hyperalignment, as used by Haxby et al. (2011), is conceptually the same as MCCA but restricting the transformations to rotations, i.e. Procrustes analysis (Xu et al., 2012). Hyperalignment has the advantage to maintain metric distance of patterns in the original and transformed space, but the disadvantage that it cannot favor channels with higher inter-subject correlation. Methods for group analysis of data from multiple subjects are reviewed by Correa et al. (2010a, b); Calhoun and Adali (2012); Sui et al. (2012); Afshin-Pour et al. (2014); Dähne et al. (2015); Chen et al. (2016); Huster and Raud (2018).

MCCA allows spatial patterns to differ across subjects, and thus accomodates individual differences in source or sensor geometry. On the other hand, it requires temporal patterns to be shared between subjects so as to distinguish them from noise, and thus it might miss a pattern present in one individual but not others. The problem is analogous to that encountered in studies that average responses over trials: the analysis is blind to eventual variations in response from trial to trial. MCCA is however sensitive to temporal patterns shared by subsets of individuals of size $> 1$. In this sense, it does not completely suppress subject individuality.

**Denoising and dimensionality weighting.** As described in the Methods and illustrated in the Results, data from single subjects can be denoised by projecting on the overcomplete basis of $D$ CCs, truncating, and projecting back. Data dimensions that are not shared with other subjects are *downweighted* but not removed, so in general the rank of the data remains the same. Setting the cutoff $\dot{D} < D$ to a relatively high order suppresses only those components that are very different from those found in other subjects, most likely to be noise. In Example 5, the set of 40 PCs that represented each subject was transformed into 320 CCs, of which 110 were selected before being projected back to obtain "denoised" data, yielding the benefit shown in Fig. 9. The CCs that were rejected absorbed some of the subject-specific temporal patterns of noise, improving the outcome. The number of CCs to retain may be tuned by means of a cross-validated parameter search procedure.

It is often useful to reduce the dimensionality of the data for computational reasons (to reduce memory or computation time), or to avoid overfitting. One standard procedure is to apply PCA, and truncate the series of PCs. However this implicitly equates variance to relevance, which may not be justified, as artifact components may have high variance, and useful sources may be weak. MCCA allows an alternative: whiten each data matrix, transform to CCs, truncate (to $\dot{D} < D$), and transform back. These 'denoised' data are no longer spatially white: applying PCA will rotate them such that the first components are most consistent with the other data matrices. MCCA thus allows to replace the variance criterion by a criterion of consistency with other data. The cutoff $\dot{D}$ is an arbitrary parameter that can be chosen by trial and error or grid search.

As a tool to analyze or denoise the data of a single subject, MCCA is comparable to data-driven linear analysis techniques such as PCA, ICA, JD, CCA and others. The fact that it uses a different criterion makes it *complementary* to those methods as a denoising or dimensionality weighting tool (e.g one can apply MCCA before or after ICA, JD, etc.).

**Caveats and cautions.** A risk, common to other data-driven methods such as ICA or JD, is circularity of the analysis (Kriegeskorte et al., 2009). The method is designed to optimize correlation between data matrices, and therefore the observation that the components that it finds *are* correlated between data matrices carries little weight unless corroborated by careful cross-validation. Related to this issue is overfitting: each SC depends on $D = \sum_n d_n$ parameters, a number that can be large if there are many data matrices involved. Overfitting can be detected using resampling and cross-validation methods, and the risk of overfitting can be reduced by dimensionality reduction or other regularization techniques (Wong et al., 2018).

MCCA can easily latch on to artefacts and noise components shared across data matrices. Uninteresting linear or polynomial trends (for example EEG drift potentials) may thus appear among the first MCCA components. More generally, MCCA can be biased towards narrowband or low-frequency components common across data matrices, *even if their*

*phase is not aligned*, particularly if the noise is spectrally-shaped or contains narrow-band components. This is illustrated in Fig. 12 that shows the result of applying MCCA to ten "data matrices", each of 12 s duration, extracted at random from the same 40-channel EEG data that was used as background noise in Example 3. No known signal is common across these data matrices, nonetheless the lowest-order SCs have narrow spectra (Fig. 12 left) and quasi-sinusoidal waveforms (right). It is easy to understand why MCCA might take such components to be shared: a sinusoid of arbitrary phase can be expressed as the weighted sum of a sine and a cosine, and thus narrowband activity can be approximated as resulting from two sinusoidal components in quadrature phase. As this is the case for all datasets, MCCA will select the two-component sinusoidal basis as common. Such spurious components compete with genuine shared activity, complicating the analysis.

A difference in latency of a brain response between different subjects might interfere with the ability of MCCA to extract this activity (Lankinen et al., 2018). A common outcome in that case is two components, one with a shape similar to the source time series averaged over subjects, and the other similar to their difference (or derivative). MCCA can readily be extended to include time-lags to account for differences in response latency between subjects, although this comes at the expense of a greater number of parameters and a greater risk of overfitting. MCCA is obviously of no benefit in the absence of synchronous temporal patterns, for example it is not well suited for analyzing resting-state data of a group of subjects. MCCA yields both CCs and SCs, either of which can be exploited. When reporting, it is important to specify which, to avoid confusion. As an example, the phrase 'MCCA was applied as a preprocessing step' is not sufficient to specify what was done.

**Applicability to real-time processing.** This work was motivated in part by the need to steer an auditory assistive device using brain signals. A common paradigm for that purpose is to correlate each of several competing sound streams with the user's brain activity, to determine which is the focus of attention, and then amplify that sound stream at the user's ears (Wronkiewicz et al., 2016). An obstacle is the high-level of noise and artifacts in the EEG data, and analysis and denoising methods are essential for the success of this application. To be useful, a method must be applicable to *real-time* processing, whereas MCCA as described here works in batch mode. It may nonetheless be of use in the following fashion. EEG data is recorded from a pool of subjects to a calibration sample of speech, and MCCA is used to derive a "canonical" EEG response to that sample. To adapt the system to a new user, EEG data are recorded in response to the calibration sample, and a spatial filter is designed (for example using CCA) to maximize similarity between the subject's and the canonical response. This spatial filter is then used in the real-time processing pipeline. This suggests that MCCA can also be put to use in a practical application such as cognitive control of a hearing aid.
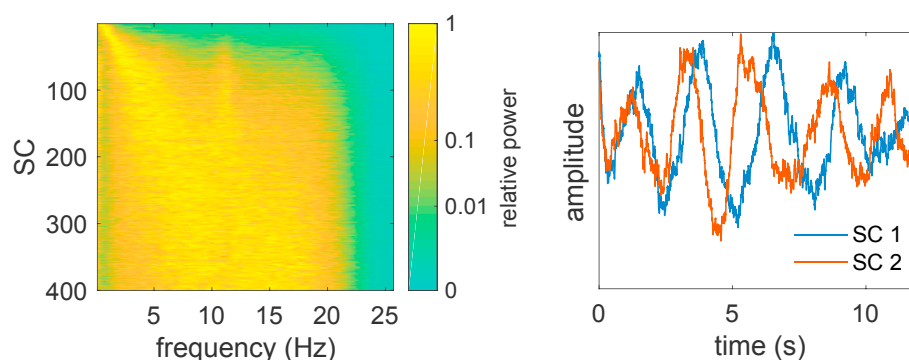


**Fig. 12.** MCCA is biased towards narrowband and low-frequency activity. Left: power spectra of SCs derived from an MCCA analysis of 10 EEG "data matrices" of duration 12 s randomly sampled from 40-channel EEG data. Power is coded as color. Right: time course of the first two SCs.

## 5. Conclusion

Multiway CCA is a powerful tool for analysis of multi-subject multivariate datasets. It can be used both to design spatial filters to denoise data of each individual subject, and to summarize data across subjects. Many related methods have been proposed in the literature, but the processing principles behind them, and the range of tasks that they can be used for, are not widely appreciated. The use of MCCA (or similar techniques) should be more prevalent given the ubiquitous need for merging data across subjects. In this paper we presented a formulation of MCCA that is relatively easy to understand, illustrated in detail how it works, and showed how it can be put to use for a wide range of common tasks in multi-subject multivariate data analysis.

## Acknowledgements

## References

Afshin-Pour, B., Grady, C., Strother, S., 2014. Evaluation of spatio-temporal decomposition techniques for group analysis of fMRI resting state data sets. Neuroimage 87, 363–382.

Afshin-Pour, B., Hossein-Zadeh, G.A., Strother, S.C., Soltanian-Zadeh, H., 2012. Enhancing reproducibility of fMRI statistical maps using generalized canonical correlation analysis in NPAIRS framework. Neuroimage 60, 1970–1981.

Asendorf, N.A., 2015. Informative Data Fusion : beyond Canonical Correlation Analysis Ph.D. Diss. University of Michigan.

Benton, A., Khayrallah, H., Gujral, B., Reisinger, D.A., Zhang, S., Arora, R., 2017. Deep Generalized Canonical Correlation Analysis. ArXiv arXiv:1702.

Calhoun, V.D., Adali, T., 2012. Multisubject independent component analysis of fMRI: a decade of intrinsic networks, default mode, and neurodiagnostic discovery. IEEE Rev. Biomed. Eng. 5, 60–73.

Chait, M., de Cheveigné, A., Poeppel, D., Simon, J.Z., 2010. Neural dynamics of attending and ignoring in human auditory cortex. Neuropsychologia 48, 3262–3271.

Chen, X., Wang, Z.J., McKeown, M., 2016. Joint blind source separation for neurophysiological data analysis: multiset and multimodal methods. IEEE Signal Process. Mag. 33, 86–107.

Choi, I., Rajaram, S., Varghese, L.A., Shinn-Cunningham, B.G., 2013. Quantifying attentional modulation of auditory-evoked cortical responses from single-trial electroencephalography. Front. Hum. Neurosci. 7.

Correa, N.M., Adali, T., Li, Y.O., Calhoun, V.D., 2010a. Canonical correlation analysis for data fusion and group inferences. IEEE Signal Process. Mag. 39–50. July.

Correa, N.M., Eichele, T., Adalı, T., Li, Y.O., Calhoun, V.D., 2010b. Multi-set canonical correlation analysis for the fusion of concurrent single trial ERP and functional MRI. Neuroimage 50, 1438–1445.

Crosse, M.J., Di Liberto, G.M., Bednar, A., Lalor, E.C., 2016. The multivariate temporal response function (mTRF) toolbox: a MATLAB toolbox for relating neural signals to continuous stimuli. Front. Hum. Neurosci. 10, 1–14.

Cunningham, J.P., Yu, B.M., 2014. Dimensionality reduction for large-scale neural recordings. Nat. Neurosci. 17, 1500Ð1509.

Dähne, S., Bießmann, F., Samek, W., Haufe, S., Goltz, D., Gundlach, C., Villringer, A., Fazli, S., Müller, K.R., 2015. Multivariate machine learning methods for fusing functional multimodal neuroimaging data. Proc. IEEE 103, 1–22.

de Cheveigné, A., Arzounian, D., 2018. Robust detrending, rereferencing, outilier detection, and inpainting of multichannel data. Neuroimage 172, 903–912.

de Cheveigné, A., Parra, L.C., 2014. Joint decorrelation, a versatile tool for multichannel data analysis. Neuroimage 98, 487–505.

de Cheveigné, A., Simon, J.Z., 2008. Denoising based on spatial filtering. J. Neurosci. Methods 171, 331–339.

de Cheveigné, A., Wong, D., Liberto, G.M.D., Hjortkjaer, J., Slaney, M., Lalor, E., 2018. Decoding the auditory brain with canonical correlation analysis. Neuroimage 172, 206–216.

Di Liberto, G.M., O'Sullivan, J.A., Lalor, E.C., 2015. Low-frequency cortical entrainment to speech reflects phoneme-level processing. Curr. Biol. : CB 25, 2457–2465.

Ding, N., Simon, J.Z., 2012. Neural coding of continuous speech in auditory cortex during monaural and dichotic listening. J. Neurophysiol. 107, 78–89.

Dmochowski, J.P., Ki, J.J., DeGuzman, P., Sajda, P., Parra, L.C., 2017. Extracting multidimensional stimulus-response correlations using hybrid encoding-decoding of neural activity. Neuroimage 1–13.

Dmochowski, J.P., Sajda, P., Dias, J., Parra, L.C., 2012. Correlated components of ongoing EEG point to emotionally Laden attention a possible marker of engagement? Front. Hum. Neurosci. 6, 112.

Eichele, T., Rachakonda, S., Brakedal, B., Eikeland, R., Calhoun, V.D., 2011. EEGIFT: group independent component analysis for event-related EEG data. Comput. Intell. Neurosci. 2011 https://doi.org/10.1155/2011/129365.

Fu, X., Huang, K., Hong, M., Sidiropoulos, N.D., So, A.M.C., 2017. Scalable and flexible multiview MAX-VAR canonical correlation analysis. IEEE Trans. Signal Process. 65, 4150–4165.

Fuglsang, S.A., Dau, T., Hjortkjær, J., 2017. Noise-robust cortical tracking of attended speech in real-world acoustic scenes. Neuroimage 156, 435–444.

Gates, K.M., Molenaar, P.C.M., 2012. Group search algorithm recovers effective connectivity maps for individuals in homogeneous and heterogeneous samples. Neuroimage 63, 310–319.

Gross, S.M., Tibshirani, R., 2015. Collaborative regression. Biostatistics 16, 326–338.

Hasson, U., Nir, Y., Levy, I., Fuhrmann, G., Malach, R., 2004. Natural vision. Science 303, 1634–1640.

Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J.D., Blankertz, B., Bießmann, F., 2014. On the interpretation of weight vectors of linear models in multivariate neuroimaging. Neuroimage 87, 96–110.

Haxby, J.V., Guntupalli, J.S., Connolly, A.C., Halchenko, Y.O., Conroy, B.R., Gobbini, M.I., Hanke, M., Ramadge, P.J., 2011. A common, high-dimensional model of the representational space in human ventral temporal cortex. Neuron 72, 404–416.

Holdgraf, C.R., Rieger, J.W., Micheli, C., Martin, S., Knight, R.T., Theunissen, F.E., 2017. Encoding and decoding models in cognitive electrophysiology. Front. Syst. Neurosci. 11.

Hotelling, H., 1936. Relations between two sets of variates. Biometrika 28, 321–377.

Huster, R.J., Plis, S.M., Calhoun, V.D., 2015. Group-level component analyses of EEG: validation and evaluation. Front. Neurosci. 9, 1–14.

Huster, R.J., Raud, L., 2018. A tutorial review on multi-subject decomposition of EEG. Brain Topogr. 31, 1–14.

Hwang, H., Jung, K., Takane, Y., Woodward, T.S., 2012. Functional multiple-set canonical correlation analysis. Psychometrika 77, 48–64.

Karch, J.D., Sander, M.C., von Oertzen, T., Brandmaier, A.M., Werkle-Bergner, M., 2015. Using within-subject pattern classification to understand lifespan age differences in oscillatory mechanisms of working memory selection and maintenance. Neuroimage 118, 538–552.

Karhunen, J., Hao, T., Ylipaavalniemi, J., 2013. Finding dependent and independent components from related data sets: a generalized canonical correlation analysis based method. Neurocomputing 113, 153–167.

Karch, J., Hao, T., Ylipaavalniemi, J., 1971. Canonical analysis of several sets of variables. Biometrika 58, 433–451.

Khalighinejad, B., Cruzatto da Silva, G., Mesgarani, N., 2017. Dynamic encoding of acoustic features in neural responses to continuous speech. J. Neurosci. 37, 2176–2185.

Kiers, H.A.L., Cléroux, R., Ten Berge, J.M.F., 1994. Generalized canonical analysis based on optimizing matrix correlations and a relation with IDIOSCAL. Comput. Stat. Data Anal. 18, 331–340.

Koskinen, M., Seppä, M., 2014. Uncovering cortical MEG responses to listened audiobook stories. Neuroimage 100, 263–270.

Kriegeskorte, N., Simmons, W.K., Bellgowan, P.S.F., Baker, C.I., 2009. Circular analysis in systems neuroscience: the dangers of double dipping. Nat. Neurosci. 12, 535–540.

Lalor, E.C., Power, A.J., Reilly, R.B., Foxe, J.J., 2009. Resolving precise temporal processing properties of the auditory system using continuous stimuli. J. Neurophysiol. 102, 349–359.

Lankinen, K., Saari, J., Hari, R., Koskinen, M., 2014. Intersubject consistency of cortical MEG signals during movie viewing. Neuroimage 92, 217–224.

Lankinen, K., Saari, J., Hlushchuk, Y., Tikka, P., Parkkonen, L., Hari, R., Koskinen, M., 2018. Consistency and similarity of meg-and fmri-signal time courses during movie viewing. Neuroimage 173, 361–369.

Li, Y.O., Adali, T., Wang, W., Calhoun, V.D., 2009. Joint blind source separation by multiset canonical correlation analysis. IEEE Trans. Signal Process. 57, 3918–3929.

Lio, G., Boulinguez, P., 2016. How does sensor-space group blind source separation face inter-individual neuroanatomical variability? Insights from a simulation study based on the PALS-B12 Atlas. Brain Topogr. 31, 1–14.

Litvak, V., Friston, K., 2008. Electromagnetic source reconstruction for group studies. Neuroimage 42, 1490–1498.

Luck, S.J., 2005. *An Introduction to the Event-related Potential Technique* the. MIT Press.

Madsen, K.H., Churchill, N.W., M£rup, M., 2016. Quantifying functional connectivity in multi?subject fmri data using component models. Hum. Brain Mapp. 38, 882–899.

Mahjoory, K., Nikulin, V.V., Botrel, L., Linkenkaer-Hansen, K., Fato, M.M., Haufe, S., 2017. Consistency of EEG source localization and connectivity estimates. Neuroimage 152, 590–601.

Martin, S., Brunner, P., Holdgraf, C., Heinze, H.J., Crone, N.E., Rieger, J., Schalk, G., Knight, R.T., Pasley, B.N., 2014. Decoding spectrotemporal features of overt and covert speech from the human cortex. Front. Neuroeng. 7, 14.

Melzer, T., Reiter, M., Bischof, H., 2001. Nonlinear feature extraction using generalized canonical correlation analysis. ICANN 2001. LNCS 2130, 353–360.

Mesgarani, N., Chang, E.F., 2012. Selective cortical representation of attended speaker in multi-talker speech perception. Nature 485, 233–236.

Mirkovic, B., Debener, S., Jaeger, M., Vos, M.D., 2015. Decoding the attended speech stream with multi-channel EEG: implications for online, daily-life applications. J. Neural. Eng. 12, 046007.

Molenaar, P.C.M., 2004. A manifesto on psychology as idiographic science: bringing the person back into scientific psychology, this time forever. Measurement: Interdiscipl. Res. Perspect. 2, 201–218.

Nielsen, A.A., 2002. Multiset canonical correlations analysis and multispectral, truly multitemporal remote sensing data. IEEE Trans. Image Process. 11, 293–305.

Norman-Haignere, S., Kanwisher, N.G., McDermott, J.H., 2015. Distinct cortical pathways for music and speech revealed by hypothesis-free voxel decomposition. Neuron 88, 1281–1296.

O'Sullivan, J.A., Power, A.J., Mesgarani, N., Rajaram, S., Foxe, J.J., Shinn-Cunningham, B.G., Slaney, M., Shamma, S.A., Lalor, E.C., 2014. Attentional selection in a cocktail party environment can Be decoded from single-trial EEG. Cerebr. Cortex 25, 1697–1706.

Parra, L., 2018. Multi-set Canonical Correlation Analysis Simply Explained arXiv arXiv: 1802.03759.

Parra, L., Haufe, S., Dmochowski, J., 2018. Correlated Components Analysis - Extracting Reliable Dimensions in Multivariate Data arXiv arXiv:1801.08881v2.

Parra, L.C., Spence, C.D., Gerson, A.D., Sajda, P., 2005. Recipes for the linear analysis of EEG. Neuroimage 28, 326–341.

Southwell, R., Baumann, A., Gal, C., Barascud, N., Friston, K., Chait, M., 2017. Is predictability salient? A study of attentional capture by auditory patterns. Phil. Trans. Biol. Sci. 372, 20160105.

Sturm, I., 2016. Analyzing the Perception of Natural Music with EEG and ECoG. Ph.D. diss. Technischen Universität Berlin.

Sui, J., Adali, T., Yu, Q., Chen, J., Calhoun, V.D., 2012. A review of multivariate methods for multimodal fusion of brain imaging data. J. Neurosci. Methods 204, 68–81.

Sun, L., Ji, S., Yu, S., Ye, J., 2009. On the equivalence between canonical correlation analysis and orthonormalized partial least squares. In: International Joint Conference on Artificial Intelligence. IJCAI 2009, p. 1230.

Takane, Y., Hwang, H., Abdi, H., 2008. Regularized multiple-set canonical correlation analysis. Psychometrika 73, 753–775.

Tenenhaus, A., 2011. Regularized generalized canonical correlation analysis and PLS path modeling. Psychometrika 76, 257–284.

Tenenhaus, A., Philippe, C., Frouin, V., 2015. Kernel generalized canonical correlation analysis. Comput. Stat. Data Anal. 90, 114–131.

Tiitinen, H., Miettinen, I., Alku, P., May, P.J.C., 2012. Transient and sustained cortical activity elicited by connected speech of varying intelligibility. BMC Neurosci. 13, 157.

van de Velden, M., Takane, Y., 2012. Generalized canonical correlation analysis with missing values. Comput. Stat. 27, 551–571.

Velden, M.V.D., 2011. On generalized canonical correlation analysis. In: Proc. 58th World Statistical Conference, pp. 758–765.

Vía, J., Santamaría, I., Pérez, J., 2005a. Canonical Correlation Analysis (Cca) Algorithms for Multiple Data Sets: Application to Blind Simo Equalization in *Signal Processing Conference*, 2005 13th European. IEEE, pp. 1–4.

Via, Javier, Santamaria, Ignacio, Pérez, J., 2005b. Canonical correlation analysis (CCA) algorithms for multiple data sets: application to blind SIMO equalization. Signal Process. Conf. 1, 4–7.

Witten, D.M., Tibshirani, R.J., 2009. Extensions of sparse canonical correlation analysis with applications to genomic data. Stat. Appl. Genet. Mol. Biol. 8, 29.

Wong, D.D.E., Fuglsang, S.A., Hjortkjær, J., Ceolini, E., Slaney, M., de Cheveigné, A., 2018. A comparison of regularization methods in forward and backward models for auditory attention decoding. Front. Neurosci. 12, 1–16.

Wronkiewicz, M., Larson, E., Lee, A.K., 2016. Incorporating modern neuroscience findings to improve brain-computer interfaces: tracking auditory attention. J. Neural. Eng. 13, 1–13.

Xu, H., Lorbert, A., Ramadge, P.J., Guntupalli, J.S., Haxby, J.V., 2012. Regularized Hyperalignment of Multi-set Fmri Data in *Statistical Signal Processing Workshop (SSP), 2012 IEEE*. IEEE, pp. 229–232.

Zhang, Q., Borst, J.P., Kass, R.E., Anderson, J.R., 2017. Inter-subject alignment of MEG datasets in a common representational space. Hum. Brain Mapp. 38, 4287–4301.

Zhang, Y., Zhou, G., Zhao, Q., Onishi, A., Jin, J., Wang, X., Cichocki, A., 2011. Multiway canonical correlation analysis for frequency components recognition in ssvep-based bcis. In: International Conference on Neural Information Processing (ICONIP 2011), vol 7062, pp. 287–295.

Zion Golumbic, E.M., Ding, N., Bickel, S., Lakatos, P., Schevon, C.A., McKhann, G.M., Goodman, R.R., Emerson, R., Mehta, A.D., Simon, J.Z., Poeppel, D., Schroeder, C.E., 2013. Mechanisms underlying selective neuronal tracking of attended speech at a "cocktail party". Neuron 77, 980–991.