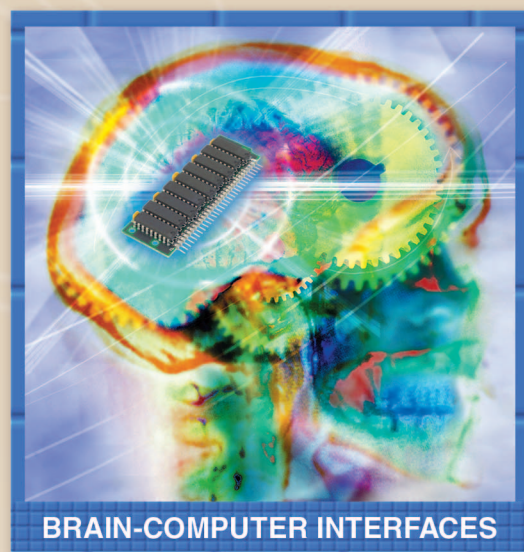


Benjamin Blankertz, Ryota Tomioka, Steven Lemm,
Motoaki Kawanabe, and Klaus-Robert Müller

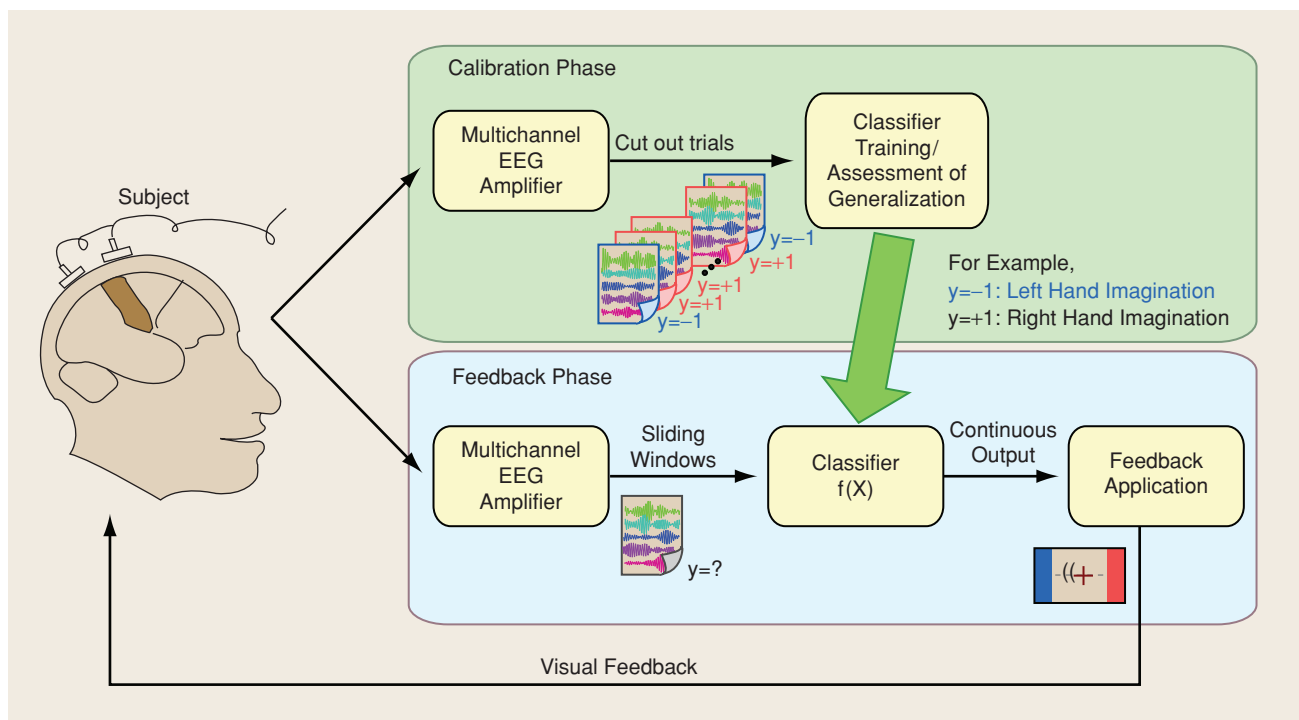


Optimizing Spatial Filters for Robust EEG Single-Trial Analysis

Revealing tricks of the trade

Due to the volume conduction multichannel electroencephalogram (EEG) recordings give a rather blurred image of brain activity. Therefore spatial filters are extremely useful in single-trial analysis in order to improve the signal-to-noise ratio. There are powerful methods from machine learning and signal processing that permit the optimization of spatio-temporal filters for each subject in a data dependent fashion beyond the fixed filters based on the sensor geometry, e.g., Laplacians. Here we elucidate the theoretical background of the common spatial pattern (CSP) algorithm, a popular method in brain-computer interface (BCI) research. Apart from reviewing several variants of the basic algorithm, we reveal tricks of the trade for achieving a powerful CSP performance, briefly elaborate on theoretical aspects of CSP, and demonstrate the application of CSP-type pre-processing in our studies of the Berlin BCI (BBCI) project.

Digital Object Identifier 10.1109/MSP.2007.909009



[FIG1] Overview of the machine-learning-based BCI system. The system runs in two phases. In the calibration phase, we instruct the subjects to perform certain tasks and collect short segments of labeled EEG (trials). We train the classifier based on these examples. In the feedback phase, we take sliding windows from continuous stream of EEG; the classifier outputs a real value that quantifies the likelihood of class membership; we run a feedback application that takes the output of the classifier as an input. Finally the subject receives the feedback on the screen as, e.g., cursor control.

INTRODUCTION

Noninvasive BCI has in the recent years become a highly active research topic in neuroscience, engineering, and signal processing. One of the reasons for this development is the striking advances of BCI systems with respect to usability, information transfer, and robustness for which **modern machine learning and signal processing techniques** have been instrumental [2], [4], [14], and [15]. Invasive BCIs [46], in particular intracranial signals, require completely different signal processing methods and are therefore not discussed here.

This article will review a particularly popular and powerful signal processing technique for EEG-based BCIs called CSP and discusses recent variants of CSP. **Our goal is to provide comprehensive information about CSP and its application.** Thus we address both the BCI expert who is not specialized in signal processing and the BCI novice who is an expert in signal processing.

Consequently, this article will mainly focus on CSP filtering, but we will also briefly discuss BCI paradigms and the neurophysiological background thereof. Finally, we will report on recent results achieved with the BBCI using advanced signal processing and machine learning techniques.

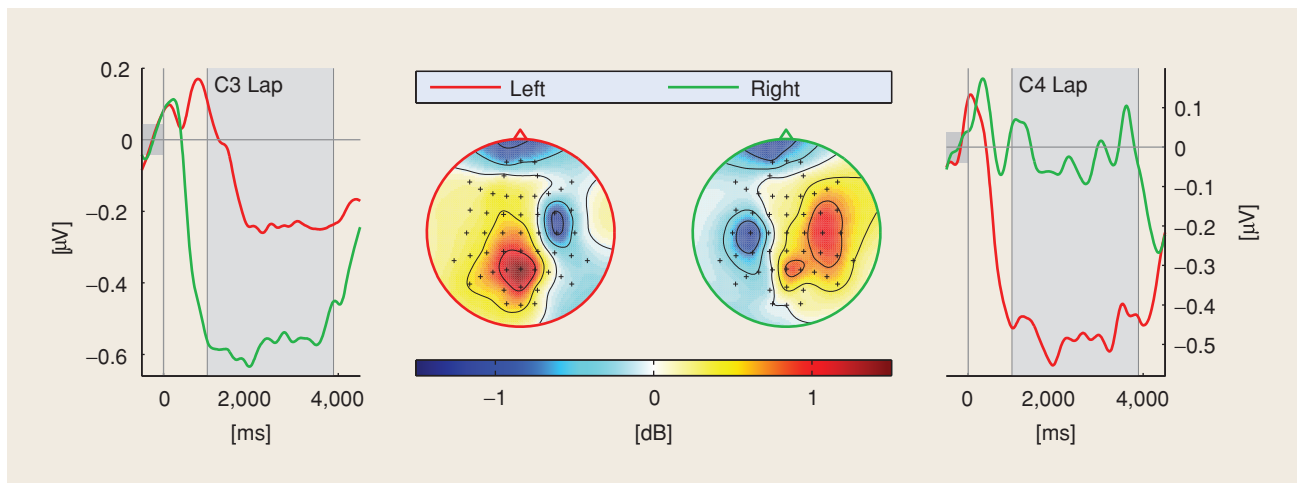
BACKGROUND

OVERVIEW OF A BCI SYSTEM

An overview of a BCI system based on machine learning is shown in Figure 1. The system operates in two phases, namely

the calibration phase and the feedback phase. The feedback phase is the time the users can actually transfer information through their brain activity and control applications; in this phase, the system is composed of the classifier that classifies between different mental states and the user interface that translates the classifier output into control signals, e.g., cursor position or selection from an alphabet. In the calibration phase, we collect examples of EEG signals in order to train the classifier. Here we describe a typical experiment as performed in the BBCI project. We use three types of imaginary movements, namely, left hand (L), right hand (R), and right foot (F) as the mental states to be classified. Other paradigms based on, e.g., modulation of attention to external stimulation can be found in [55]. The subjects are instructed to perform one of the three imaginary movements indicated on the screen for 3.5 s at the interval of 5.5 s. For more effective performance, it is important to instruct the subjects to concentrate on the kinesthetic aspect rather than the visual [37]. We obtain 420 trials of imaginary movement (140 for each class) in a randomized order for each subject (less is sufficient for feedback performance). The data is then used for the training of the classifier and assessment of generalization error by cross-validation. In particular, we compare three pair-wise classifiers and select the combination of two classes that yields the best generalization performance.

After the calibration measurement subjects perform five feedback sessions consisting of 100 runs. Here the output of the



[FIG2] Event-related desynchronization (ERD) during motor imagery of the left and the right hand. Raw EEG signals of one subject have been band-pass filtered between 9–13 Hz. For the time courses, the envelope of the signals has been calculated by Hilbert transform (see e.g., [9]) and averaged over segments of –500–4,500 ms relative to each cue for left or right hand motor imagery. ERD curves are shown for Laplace filtered channels at C3 and C4, i.e., over left and right primary motor cortex. The topographical maps of ERD were obtained by performing the same procedure for all (non-Laplace filtered) channels and averaging across the shaded time interval 1,000 to 4,000 ms.

binary classifier is translated into the horizontal position of a cursor. Subjects are instructed to move the cursor to that one of the two vertical bars at the edges of the screen which was indicated as target by color. The cursor is initially at the center of the screen; it starts to follow the classifier output based on the brain signal 750 ms after the indication of the target. A trial ends when the cursor touches one of the two bars; the bar that the cursor reached is colored green if correct and red otherwise. The next trial starts after 520 ms (see [2], [4], [7] for more details).

The performance of the classifier is measured by the accuracy of the prediction in percent. The performance of the overall system is measured by the information transfer rate (ITR) [54] measured in b/min:

$$\text{ITR} = \frac{\text{\# of decisions}}{\text{duration in minutes}} \cdot \left(p \log_2(p) + (1-p) \log_2\left(\frac{1-p}{N-1}\right) + \log_2(N) \right), \quad (1)$$

where p is the accuracy of the subject in making decisions between N targets, e.g., in the feedback explained above, $N = 2$ and p is the accuracy of hitting the correct bars. ITR measures the capacity of a symmetric communication channel that makes mistake with the equal probability $(1-p)/(N-1)$ to all other $N-1$ classes divided by the time required to communicate that amount of information. The ITR depends not only on the accuracy of the classifier but also on the design of the feedback application that translates the classifier output into command. Note that the duration in min refers to the total duration of the run including all inter-trial intervals. In contrast to the accuracy of the decision, the ITR takes different duration of trials and different number of classes into account. The ITR is zero for a ran-

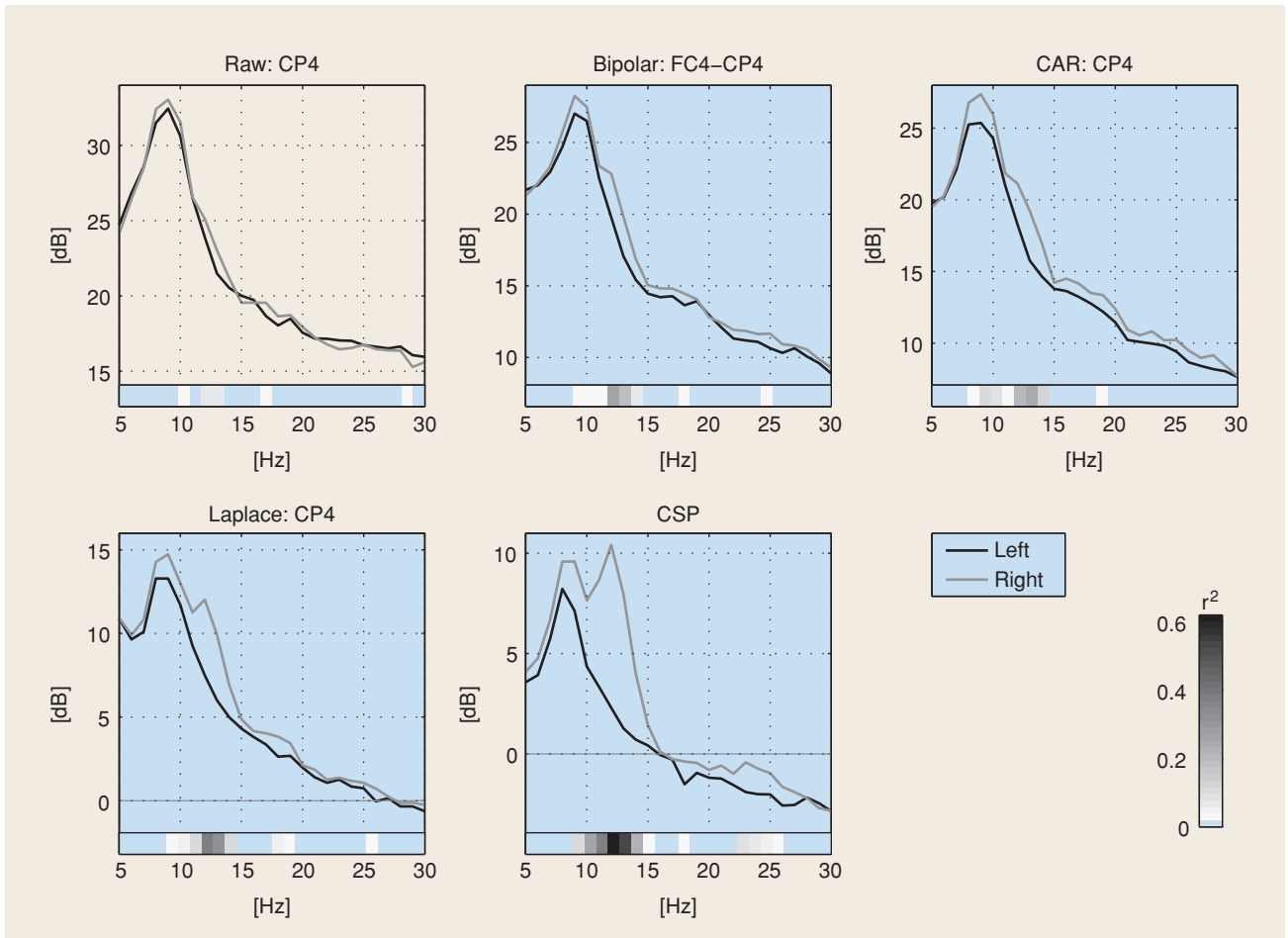
dom classifier i.e., $p = 1/N$. Note that the communication channel model can be generalized to take the nonsymmetric or nonuniform errors into account [44].

Brain activity was recorded from the scalp with multi-channel EEG amplifiers (BrainAmp by Brain Products, Munich, Germany) using 55 Ag/AgCl electrodes in an extended 10–20 system.

NEUROPHYSIOLOGICAL BACKGROUND

Macroscopic brain activity during resting wakefulness comprises distinct “idle” rhythms located over various cortical areas, e.g. the occipital α -rhythm (8–12 Hz) can be measured over the visual cortex [1]. The perirolandic sensorimotor cortices show rhythmic macroscopic EEG oscillations (μ -rhythm, sensorimotor rhythm, SMR) [20], [24] with spectral peak energies of about 9–14 Hz (localized predominantly over the postcentral somatosensory cortex) and around 20 Hz (over the precentral motor cortex). The occipital α -rhythm is quite prominent and can be seen in the raw EEG with the naked eye if the subject closes the eyes (idling of the visual cortex). In contrast the μ -rhythm has a much weaker amplitude and can only be observed after appropriate signal processing. In some subjects no μ -rhythm can be observed in scalp EEG.

Our system is based on the modulation of the SMR. In fact, motor activity, both actual and imagined [25], [42], [45], as well as somatosensory stimulation [38] have been reported to modulate the μ -rhythm. Processing of motor commands or somatosensory stimuli causes an attenuation of the rhythmic activity termed event-related desynchronization (ERD) [42], while an increase in the rhythmic activity is termed event-related synchronization (ERS). For BCIs, the important fact is that the ERD is caused also by imagined movements (healthy users, see Figure 2) and by intended movements in paralyzed patients [30].



[FIG4] Spectra of left versus right hand motor imagery. All plots are calculated from the same dataset but using different spatial filters. The discrimination between the two conditions is quantified by the r^2 -value. CAR stands for common average reference.

The classifier first projects the signal by J spatial filters $\{w_j\}_{j=1}^J \in \mathbb{R}^{C \times J}$; next it takes the logarithm of the power of the projected signal; finally it linearly combines these J dimensional features and adds a bias β_0 . In fact, each projection captures different spatial localization; the modulation of the rhythmic activity is captured by the log-power of the band-pass filtered signal. Note that various extensions are possible (see the section on variants and extensions of the original CSP algorithm). A different experimental paradigm might require the use of nonlinear methods of feature extraction and classification respectively [33]. Direct minimization of discriminative criterion [17] and marginalization of the classifier weight [22] are suggested. On the other hand, methods that are linear in the second order statistics XX^T , i.e., (2) without the log, are discussed in [48], [49] and shown to have some good properties such as convexity.

The coefficients $\{w_j\}_{j=1}^J$ and $\{\beta_j\}_{j=1}^J$ are automatically determined statistically [21] from the training examples i.e., the pairs of trials and labels $\{X_i, y_i\}_{i=1}^n$ we collect in the calibration phase; the label $y \in \{+1, -1\}$ corresponds to, e.g., imaginary movement of left and right hand, respectively, and n is the number of trials.

We use CSP [18], [27] to determine the spatial filter coefficients $\{w_j\}_{j=1}^J$. In the following, we discuss the method in detail and present some recent extensions. The linear weights $\{\beta_j\}_{j=1}^J$ are determined by Fisher's linear discriminant analysis (LDA).

INTRODUCTION TO CSP ANALYSIS

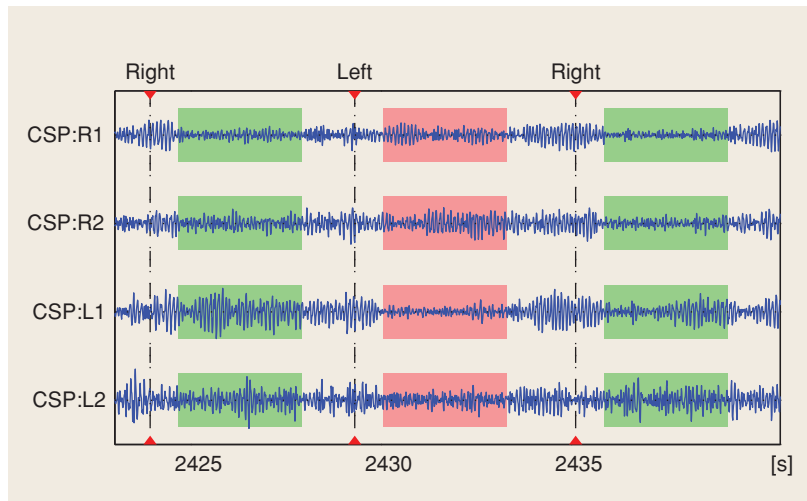
CSP [18], [27] is a technique to analyze multichannel data based on recordings from two classes (conditions). CSP yields a data-driven supervised decomposition of the signal parameterized by a matrix $W \in \mathbb{R}^{C \times C}$ (C being the number of channels) that projects the signal $x(t) \in \mathbb{R}^C$ in the original sensor space to $x_{\text{CSP}}(t) \in \mathbb{R}^C$, which lives in the surrogate sensor space, as follows:

$$x_{\text{CSP}}(t) = W^T x(t).$$

In this article, we call each column vector $w_j \in \mathbb{R}^C$ ($j = 1, \dots, C$) of W a spatial filter or simply a filter; moreover we call each column vector $a_j \in \mathbb{R}^C$ ($j = 1, \dots, C$) of a matrix $A = (W^{-1})^T \in \mathbb{R}^{C \times C}$ a spatial pattern or simply a pattern. In fact, if we think of the signal spanned by A as $x(t) = \sum_{j=1}^C a_j s_j(t)$, each vector a_j characterizes the spatial pattern of the j -th activity; moreover, w_j would filter out all but

the j -th activity because the orthogonality $w_j^\top a_k = \delta_{jk}$ holds, where δ_{jk} is the Kronecker delta ($\delta_{jk} = 1$ for $j = k$ and $= 0$ for $j \neq k$). The matrices A and W are sometimes called the mixing and de-mixing matrix or the forward and backward model [41] in other contexts.

The optimization criterion that is used to determine the CSP filters will be discussed in detail in the subsequent section on technical approaches to CSP analysis. **In a nutshell, CSP filters maximize the variance of the spatially filtered signal under one condition while minimizing it for the other condition.** Since variance of band-pass filtered signals is equal to band-power, **CSP analysis is applied to approximately band-pass filtered signals in order to obtain an effective discrimination of mental states that are characterized by ERD/ERS effects.** Figure 5 shows the result of applying four CSP filters to continuous band-pass filtered EEG data. Intervals of right hand motor imagery are shaded green and show larger variance in the CSP:L1 and CSP:L2 filters, while during left hand motor imagery (shaded red) variance is larger in the CSP:R1 and CSP:R2 filters. See also the visualization of spatial maps of CSP analysis in the section on visualization of the spatial filter coefficients.



[FIG5] Effect of spatial CSP filtering. CSP analysis was performed to obtain four spatial filters that discriminate left from right hand motor imagery. The graph shows continuous band-pass filtered EEG after applying the CSP filters. The resulting signals in filters CSP:L1 and CSP:L2 have larger variance during right hand imagery (segments shaded in green) while signals in filters CSP:R1 and CSP:R2 have larger variance during left hand imagery (segment shaded red).

TECHNICAL APPROACHES TO CSP ANALYSIS

Let $\Sigma^{(+)} \in \mathbb{R}^{C \times C}$ and $\Sigma^{(-)} \in \mathbb{R}^{C \times C}$ be the estimates of the covariance matrices of the band-pass filtered EEG signal in the two conditions (e.g., left hand imagination and right hand imagination):

$$\Sigma^{(c)} = \frac{1}{|\mathcal{I}_c|} \sum_{i \in \mathcal{I}_c} X_i X_i^\top \quad (c \in \{+, -\}), \quad (3)$$

where \mathcal{I}_c ($c \in \{+, -\}$) is the set of indices corresponding to trials belonging to each condition and $|\mathcal{I}|$ denotes the size of a set

\mathcal{I} . The above expression gives a pooled estimated of covariance in each condition because each X is centered and scaled. Then CSP analysis is given by the simultaneous diagonalization of the two covariance matrices

$$\begin{aligned} W^\top \Sigma^{(+)} W &= \Lambda^{(+)}, \\ W^\top \Sigma^{(-)} W &= \Lambda^{(-)}, \quad (\Lambda^{(c)} \text{ diagonal}), \end{aligned} \quad (4)$$

where the scaling of W is commonly determined such that $\Lambda^{(+)} + \Lambda^{(-)} = I$ [18]. Technically, this can simply be achieved (In MATLAB: $W = \text{eig}(S1, S1 + S2)$) by solving the generalized eigenvalue problem

$$\Sigma^{(+)} w = \lambda \Sigma^{(-)} w. \quad (5)$$

Then (4) is satisfied for W consisting of the generalized eigenvectors w_j ($j = 1, \dots, C$) of (5) (as column vectors) and $\lambda_j^{(c)} = w_j^\top \Sigma^{(c)} w_j$ being the corresponding diagonal elements of $\Lambda^{(c)}$ ($c \in \{+, -\}$), while λ in (5) equals $\lambda_j^{(+)} / \lambda_j^{(-)}$. Note that $\lambda_j^{(c)} \geq 0$ is the variance in condition c in the corresponding surrogate channel and $\lambda_j^{(+)} + \lambda_j^{(-)} = 1$. Hence a large value $\lambda_j^{(+)}$ ($\lambda_j^{(-)}$) close to one indicates that the corresponding spatial filter w_j yields high variance in the positive (negative) condition and low variance in the negative (positive) condition, respectively; this contrast between two classes is useful in the discrimination.

Koles [27] explained that the above decomposition gives a *common* basis of two conditions because the filtered signal $x_{\text{CSP}}(t) = W^\top x(t)$ is uncorrelated in both conditions, which implies ‘independence’ for Gaussian random variables. Figure 6 explains how CSP works in 2-D. CSP maps the samples in Figure 6(a) to those in Figure 6(b); the strong correlation between the original two axes is removed and both distributions are simultaneously de-correlated. Additionally, the two distributions are maximally dissimilar along the new axes. The dashed lines in Figure 6 denote the direction of the CSP projections. Note that the two vectors are not orthogonal to each other; in fact they are rather almost orthogonal to the

direction that the opponent class has the maximum variance.

A generative view on CSP was provided by [40]. Let us consider the following linear mixing model with nonstationary sources:

$$x_c = A s_c, \quad s_c \sim \mathcal{N}(0, \Lambda^{(c)}) \quad (c \in \{+, -\}),$$

where the sources $s_c \in \mathbb{R}^C$ ($c \in \{+, -\}$) are assumed to be uncorrelated Gaussian distributions with covariance matrices $\Lambda^{(c)}$ ($c \in \{+, -\}$) for two conditions respectively. If the empirical estimates $\Sigma^{(c)}$ are reasonably close to the true covariance

matrices $A\Lambda^{(c)}A^\top$, the simultaneous diagonalization gives the maximum likelihood estimator of the backward model $W = (A^{-1})^\top$.

A discriminative view is the following (see the section on variants and extensions of the original CSP algorithm). Let us define S_d and S_c as follows:

$$\begin{aligned} S_d &= \Sigma^{(+)} - \Sigma^{(-)} && \text{: discriminative activity,} \\ S_c &= \Sigma^{(+)} + \Sigma^{(-)} && \text{: common activity,} \end{aligned} \quad (6)$$

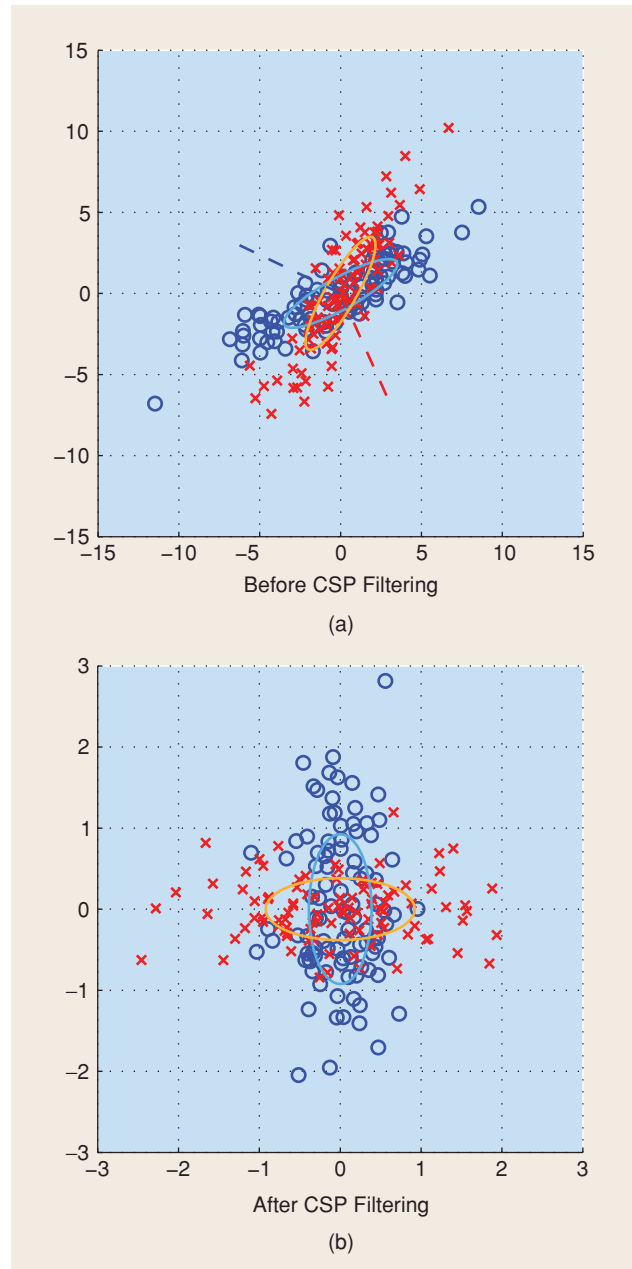
where S_d corresponds to the discriminative activity, i.e., the band-power modulation between two conditions and S_c corresponds to the common activity in the two conditions that we are not interested in. Then a solution to the following maximization problem (Rayleigh coefficient) can be obtained by solving the same generalized eigenvalue problem,

$$\underset{w \in \mathbb{R}^c}{\text{maximize}} \quad \frac{w^\top S_d w}{w^\top S_c w}. \quad (7)$$

It is easy to see that every generalized eigenvector w_j corresponds to a local stationary point with the objective value $\lambda_j^{(+)} - \lambda_j^{(-)}$ (assuming $\lambda_j^{(+)} + \lambda_j^{(-)} = 1$ as above). The large positive (or negative) objective value corresponds to large response in the first (or the second) condition. Therefore, the common practice in a classification setting is to use several eigenvectors from both ends of the eigenvalue spectrum as spatial filters $\{w_j\}_{j=1}^J$ in (2). If the number of components J is too small, the classifier would fail to fully capture the discrimination between two classes (see also the discussion in the section on merits and caveats on the influence of artifacts); on the other hand, the classifier weights $\{\beta_j\}_{j=1}^J$ could severely overfit if J is too large. In practice we find $J = 6$, i.e., three eigenvectors from both ends, often satisfactory. Alternatively one can choose the eigenvectors according to different criterion (see the section on how to select hyperparameters for CSP) or use cross-validation to determine the number of components.

FEEDBACK WITH CSP FILTERS

During BCI feedback, the most recent segment of EEG is processed and translated by the classifier into a control signal, (see Figure 1). This can be done according to (2), where X denotes the band-pass filtered segment of EEG. Due to the linearity of temporal (band-pass) and spatial filtering, these two steps can be interchanged in order. This reduces the computation load (number of signals that are band-pass filtered), since the number of selected CSP filters is typically low (2–6) compared to the number of EEG channels (32–128). Furthermore, it is noteworthy that the length of segment which is used to calculate one time instance of the control signal can be changed during feedback. Shorter segments result in more responsive but also more noisy feedback signal. Longer segments give a smoother control signal, but the delay from intention to control gets longer. This trade-off can be adapted to the aptitude of the



[FIG6] A toy example of CSP filtering in 2-D. Two sets of samples marked by red crosses and blue circles are drawn from two Gaussian distributions. In (a), the distribution of samples before filtering is shown. Two ellipses show the estimated covariances and dashed lines show the direction of CSP projections w_j ($j = 1, 2$). In (b), the distribution of samples after the filtering is shown. Note that both classes are uncorrelated at the same time; the horizontal (vertical) axis gives the largest variance in the red (blue) class and the smallest in the blue (red) class, respectively.

subject and the needs of the application. As a caveat, we remark that for optimal feedback the bias of the classifier $[\beta_0 \text{ in } (2)]$ might need to be adjusted for feedback. Since the mental state of the user is very much different during the feedback phase compared to the calibration phase, also the nontask related

[TABLE 1] RESULTS OF A FEEDBACK STUDY WITH SIX HEALTHY SUBJECTS (IDENTIFICATION CODE IN COLUMN 1). FROM THE THREE CLASSES USED IN THE CALIBRATION MEASUREMENT. THE TWO CHOSEN FOR FEEDBACK ARE INDICATED IN COLUMN 2 (L: LEFT HAND, R: RIGHT HAND, F: RIGHT FOOT). COLUMNS 3 AND 4 COMPARE THE ACCURACY AS CALCULATED BY CROSS-VALIDATION ON THE CALIBRATION DATA WITH THE ACCURACY OBTAINED ONLINE IN THE FEEDBACK APPLICATION “RATE CONTROLLED CURSOR”. THE AVERAGE DURATION \pm STANDARD DEVIATION OF THE FEEDBACK TRIALS IS PROVIDED IN COLUMN 5 (DURATION FROM CUE PRESENTATION TO TARGET HIT). SUBJECTS ARE SORTED ACCORDING TO FEEDBACK ACCURACY. COLUMNS 6 AND 7 REPORT THE INFORMATION TRANSFER RATES (ITR) MEASURED IN B/MIN AS OBTAINED BY SHANNON’S FORMULA, COMPARE (1). HERE THE COMPLETE DURATION OF EACH RUN WAS TAKEN INTO ACCOUNT, I.E., ALSO THE INTER-TRIAL BREAKS FROM TARGET HIT TO THE PRESENTATION OF THE NEXT CUE. THE COLUMN *OVERALL ITR* REPORTS THE AVERAGE ITR OF ALL RUNS (OF 25 TRIALS EACH), WHILE COLUMN *PEAK ITR* REPORTS THE PEAK ITR OF ALL RUNS. FOR SUBJECT *AU* NO REASONABLE CLASSIFIER COULD BE TRAINED (CROSS-VALIDATION ACCURACY BELOW 65% IN THE CALIBRATION DATA), SEE [2] FOR AN ANALYSIS OF THAT SPECIFIC CASE.

SUBJECT	CLASSES	CALIBRATION	FEEDBACK			
		ACCURACY [%]	ACCURACY [%]	DURATION [S]	OVERALL ITR [B/MIN]	PEAK ITR [B/MIN]
<i>al</i>	LF	98.0	98.0 \pm 4.3	2.0 \pm 0.9	24.4	35.4
<i>ay</i>	LR	97.6	95.0 \pm 3.3	1.8 \pm 0.8	22.6	31.5
<i>av</i>	LF	78.1	90.5 \pm 10.2	3.5 \pm 2.9	9.0	24.5
<i>aa</i>	LR	78.2	88.5 \pm 8.1	1.5 \pm 0.4	17.4	37.1
<i>aw</i>	RF	95.4	80.5 \pm 5.8	2.6 \pm 1.5	5.9	11.0
<i>au</i>	—	—	—	—	—	—
MEAN		89.5	90.5 \pm 7.6	2.3 \pm 0.8	15.9	27.9

brain activity differs. For a thorough investigation of this issue compare [29], [47].

RESULTS

PERFORMANCE IN TWO BBCI FEEDBACK STUDIES

Here we summarize the results of two feedback studies with healthy subjects. The first was performed to explore the limits of information transfer rates in BCIs system not relying on user training or evoked potentials and the objective of the second was to investigate for what proportion of naive subjects our system could provide successful feedback in the very first session. One of the keys to success in this study was the proper application of CSP analysis. Details can be found in [2], [3], and [7].

Table 1 summarizes performance, in particular the information transfer rates that were obtained in the first study. Note that calibration and feedback accuracy refer to quite different

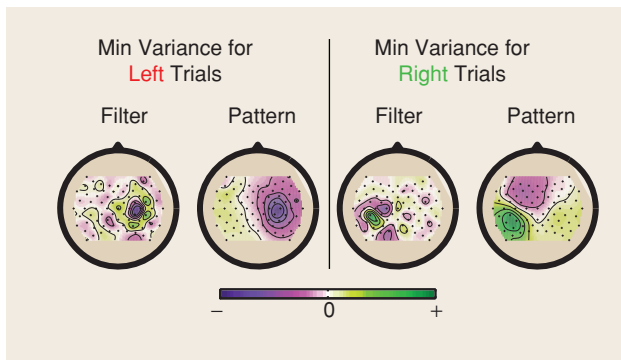
measures. From the calibration measurement, trials of approximately 3 s after each cue presentation have been taken out and the performance of the processing/classification method was validated by cross-validation. The feedback accuracy refers to the actual hitting of the correct target during horizontal cursor control. This involves integration of several classifier outputs to consecutive sliding windows of 300 to 1,000 ms length.

As a test of practical usability, subject *al* operated a mental typewriter based on horizontal cursor control. In a free spelling mode he spelled three German sentences with a total of 135 characters in 30 min, which is a “typing” speed of 4.5 letters per min. Note that the subject corrected all errors using the deletion symbol. For details, see [11]. Recently, using the novel mental typewriter Hex-o-Spell that is based on principles of human-computer interaction, the same subject achieved a typing speed of more than seven letters per min, compare [6], [34].

Table 2 summarizes the performance obtained in the second study. It demonstrates that 12 out of 14 BCI novices were able for control the BCI system in their very first session. In this study, the feedback application was not optimized for fast performance, which results in longer trial duration times.

VISUALIZATION OF THE SPATIAL FILTER COEFFICIENTS

Let us visualize the spatial filter coefficients and the corresponding pattern of activation in the brain and see how they correspond to the neurophysiological understanding of ERD/ERS for motor imagination. Figure 7 displays two pairs of vectors (w_j, a_j) that correspond to the largest and the smallest eigenvalues for one subject topographically mapped onto a scalp and color coded. w_j and a_j are the j -th columns of W and $A = (W^{-1})^T$, respectively. The plot shows the interpolation of the values of the components of vectors w_j and a_j at electrode positions. Note that we use a colormap that has no direct association to signs because the signs of the vectors are irrelevant in our analysis.



[FIG7] Example of CSP analysis. The patterns (a_j) illustrate how the presumed sources project to the scalp. They can be used to verify neurophysiological plausibility. The filters (w_j) are used to project the original signals. Here they resemble the patterns but their intricate weighting is essential to obtain signals that are optimally discriminative with respect to variance. See the introduction to CSP Patterns Analysis for the definition of the terms *filter* and *pattern*.

[TABLE 2] PERFORMANCE RESULTS FOR ALL 14 SUBJECTS OF THE SECOND STUDY. COLUMN 1 SHOWS THE SUBJECT CODE AND COLUMN 2 SHOWS A TWO LETTER CODE WHICH INDICATES THE CLASSES WHICH HAVE BEEN USED FOR FEEDBACK. COLUMN 3 SHOWS THE AVERAGE ACCURACY DURING THE FEEDBACK \pm THE STANDARD ERROR OF INTRA-RUN AVERAGES. THE AVERAGE DURATION \pm STANDARD DEVIATION OF THE FEEDBACK TRIALS IS PROVIDED IN COLUMN 4 (DURATION FROM CUE PRESENTATION TO TARGET HIT). SUBJECTS ARE SORTED ACCORDING TO FEEDBACK ACCURACY. FOR SUBJECT CQ NO REASONABLE CLASSIFIER COULD BE TRAINED.

SUBJECT	CLASSES	CALIBRATION	FEEDBACK	
		ACCURACY [%]	ACCURACY [%]	DURATION [S]
cm	LR	88.9	93.2 \pm 3.9	3.5 \pm 2.7
ct	LR	89.0	91.4 \pm 5.1	2.7 \pm 1.5
cp	LF	93.8	90.3 \pm 4.9	3.1 \pm 1.4
zp	LR	84.7	88.0 \pm 4.8	3.6 \pm 2.1
cs	LR	96.3	87.4 \pm 2.7	3.9 \pm 2.3
cu	LF	82.6	86.5 \pm 2.8	3.3 \pm 2.7
ea	FR	91.6	85.7 \pm 8.5	3.8 \pm 2.2
at	LF	82.3	84.3 \pm 13.1	10.0 \pm 8.3
zr	LF	96.8	80.7 \pm 6.0	3.1 \pm 1.9
co	LF	87.1	75.9 \pm 4.8	4.6 \pm 3.1
eb	LF	81.3	73.1 \pm 5.6	5.9 \pm 4.8
cr	LR	83.3	71.3 \pm 12.6	4.9 \pm 3.7
cn	LF	77.5	53.6 \pm 6.1	3.9 \pm 2.4
cq	—	—	—	—
MEAN		87.3	82.6 \pm 11.4	4.3 \pm 1.9

[TABLE 3] COMPARISON OF CSP-BASED CLASSIFICATION PERFORMANCE WHEN THE HYPERPARAMETERS ARE FIXED A-PRIORI, SELECTED AUTOMATICALLY BY THE PROPOSED HEURISTICS, OR SELECTED MANUALLY. EVALUATION BY A CHRONOLOGICAL SPLIT OF THE CALIBRATION DATA (FIRST HALF FOR TRAINING, SECOND HALF FOR TESTING). NOTE THAT "AUTO" USES ONLY THE FIRST HALF FOR HYPERPARAMETER SELECTION, WHEREAS "MANUAL" USES THE WHOLE CALIBRATION DATA.

SBJ	FIXED	AUTO	MANUAL
zq	2.5	0.5	0.1
zp	11.9	14.8	8.1
zr	0.8	0.2	0.2
cs	9.6	4.1	1.3
at	6.9	6.7	6.7
ct	20.7	8.9	5.2
zk	9.9	6.0	1.5
cm	14.9	6.5	5.0
cm	15.1	6.4	2.1
cm	18.2	18.2	6.9
cm	13.7	8.2	5.0
ea	5.7	1.7	1.6
eb	25.0	27.1	12.1
MEAN	11.9	8.4	4.3

[TABLE 4] COMPARISON OF PERFORMANCE ANALOG TO TABLE 3, BUT WITH EVALUATION BY TRAINING ON THE WHOLE CALIBRATION MEASUREMENT AND TESTING ON THE FEEDBACK DATA (WINDOWS OF 1,000 MS DURING CURSOR MOVEMENT). NOTE THAT THESE ERROR RATES DO NOT REFLECT THE ERRORS IN HITTING THE CORRECT BARS; A SUCCESSFUL TRIAL OFTEN INCLUDES ERRONEOUS INTERMEDIATE STEPS.

SBJ	FIXED	AUTO	MANUAL
zq	17.4	13.1	12.5
zp	24.4	24.6	22.8
zr	25.3	18.6	23.1
cs	26.1	23.0	21.8
at	39.6	34.9	33.6
ct	12.1	31.0	10.9
zk	28.2	27.3	28.8
cm	19.9	8.8	7.4
cm	6.2	2.5	2.0
cm	7.7	6.6	6.1
cm	27.7	7.0	5.9
ea	21.6	20.4	19.1
eb	50.3	42.3	39.1
MEAN	23.6	20.0	17.9

DISCUSSION

DEPENDENCE OF LINEAR SPATIAL FILTERING PRIOR TO CSP

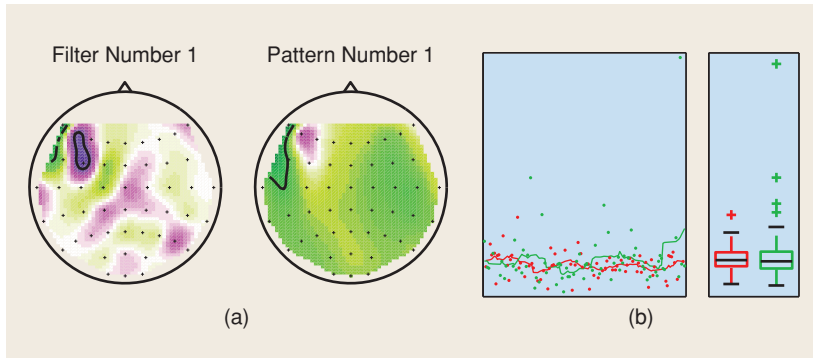
The question arises whether the results of CSP-based classification can be enhanced by preprocessing the data with a linear spatial filter [like principle component analysis (PCA), independent component analysis (ICA) or re-referencing like Laplace filtering]. The question is difficult to answer in general, but two facts can be derived. Let $B \in \mathbb{R}^{C \times C_0}$ be the matrix representing an arbitrary linear spatial filter while using notions X_i , $\Sigma^{(+)}$, $\Sigma^{(-)}$, S_d , and S_c as noted in the section on technical approaches to CSP patterns analysis. Denoting all variables corresponding to the B -filtered signals by $\tilde{\cdot}$, the signals are $\tilde{X} = B^T X$. This implies $\tilde{\Sigma}^{(+)} = B^T \Sigma^{(+)} B$,

$\tilde{\Sigma}^{(-)} = B^T \Sigma^{(-)} B$, $\tilde{S}_d = B^T S_d B$, and $\tilde{S}_c = B^T S_c B$. The filter matrices calculated by CSP are denoted by W and \tilde{W} .

1) If matrix B is invertible, the classification results will exactly be identical, regardless of applying filter B before calculating CSP or not: Let us consider the CSP solution characterized by simultaneous diagonalization of $\Sigma^{(+)}$ and $\Sigma^{(-)}$ in (4) with constraint $\Lambda^{(+)} + \Lambda^{(-)} = I$. This implies

$$\begin{aligned} (B^{-1}W)^T \tilde{\Sigma}^{(+)} B^{-1}W &= \Lambda^{(+)} \\ (B^{-1}W)^T \tilde{\Sigma}^{(-)} B^{-1}W &= I - \Lambda^{(+)}, \end{aligned}$$

which means that $B^{-1}W$ is a solution to the simultaneous diagonalization of $\tilde{\Sigma}^{(+)}$ and $\tilde{\Sigma}^{(-)}$. Since the solution is unique up to the sign of the columns, we obtain



[FIG8] (a) CSP filter/pattern corresponding to the ‘best’ eigenvalue in the data set of subject *cr*. This CSP solution is highly influenced by one single-trial in which channel FC3 has a very high variance. (b) shows the variance of all single-trials of the training data (x-axis: number of trial in chronological order, y-axis: log variance of the trial in the CSP surrogate channel; green: left hand imagery, red: right hand imagery). The trial which caused the distorted filter can be identified as the point in the upper right corner. Note that the class-specific box-plots in (b) show no difference in median of the variances (black line).

$$\tilde{W}D = B^{-1}W \quad \text{with diagonal } D: \quad (D)_{j,j} = \text{sign} \left(w_j^T B \tilde{w}_j \right).$$

Accordingly, the filtered signals are identical up to the sign: $W^T X = D \tilde{W}^T B^T X = D \tilde{W}^T \tilde{X}$, so the features, the classifier and the classification performance do not change.

2) If matrix B is not invertible, the objective of CSP analysis (on the training data) can only get worse. This can easily be seen in terms of the objective of the CSP-maximization in the formulation of the Rayleigh coefficient, (7). Then the following holds

$$\begin{aligned} \max_{\tilde{w} \in \mathbb{R}^{C_0}} \frac{\tilde{w}^T \tilde{S}_d \tilde{w}}{\tilde{w}^T \tilde{S}_c \tilde{w}} &= \max_{\tilde{w} \in \mathbb{R}^{C_0}} \frac{\tilde{w}^T B^T S_d B \tilde{w}}{\tilde{w}^T B^T S_c B \tilde{w}} \\ &\leq \max_{w \in \mathbb{R}^C} \frac{w^T S_d w}{w^T S_c w} \end{aligned}$$

since every term on the left hand side of the inequality is covered on the right hand side for $w = B \tilde{w}$. That means, the CSP-optimum for the unfiltered signals (right hand side) is greater than or equal to the CSP-optimum for the signals filtered by B (left hand side). However, this result holds only for the training data, i.e., it may be affected by overfitting effects. If the prefiltering reduces artifacts, it is well possible that the generalization performance of CSP improves. On the other hand, the prefiltering could also discard discriminative information which would be detrimental for performance.

MERITS AND CAVEATS

The CSP technique is very successfully used in online BCI systems [2], [19]. Also in the BCI Competition III many of the successful methods involved CSP type spatial filtering [8]. Apart from the above results, an advantage of CSP is the interpretability of its solutions. Far from being a black-box method, the result

of the CSP optimization procedure can visualized as scalp topographies (filters and patterns). These maps can be used to check plausibility and to investigate neurophysiological properties, compare the section on visualization of the spatial filter coefficients and also Figure 8.

It is important to point out that CSP is not a source separation or localization method. On the contrary, each filter is optimized for two effects: maximization of variance for one class while minimizing variance for the other class. Let us consider, e.g., a filter that maximizes variance for class foot and minimizes it for right: A strong focus on the left hemispherical motor area (corresponding to the right hand) can have two plausible reasons. It can either originate from an ERD during right hand imagery, or from an ERS

during foot imagery (hand areas are more relaxed if concentration focuses on the foot, therefore the idle rhythm may increase; lateral inhibition [36], [43]). Or it can be a mixture of both effects. For the discrimination task, this mixing effect is irrelevant. However, this limitation has to be kept in mind for neurophysiological interpretation.

Several parameters have to be selected before CSP can be used: the band-pass filter and the time intervals (typically a fixed time interval relative to all stimuli/responses) and the subset of CSP filters that are to be used. Often some general settings are used (frequency band 7–30 Hz ([35]), time interval starting 1,000 ms after cue, two or three filters from each side of the eigenvalue spectrum). But there is report that on-line performance can be much enhanced by subject-specific settings [2]. In Appendix A, we give a heuristic procedure for selection of CSP hyperparameters and demonstrate its favorable impact on classification. A practical example where parameters are selected manually is given in [15].

In addition, one should keep in mind that the discriminative criterion (6) tells only the separation of the mean power of two classes. The mean separation might be insufficient to tell the discrimination of samples around the decision boundary. Moreover, the mean might be sensitive to outliers. Artifacts, such as blinking and other muscle movements can dominate over EEG signals giving excessive power in some channels. If the artifact happens to be unevenly distributed in two conditions (due to its rareness), one CSP filter will likely to capture it with very high eigenvalue. Taking one specific data set from our database as an example, the CSP filter/pattern corresponding to the best eigenvalue shown in Figure 8 is mainly caused by one single trial. This is obviously a highly undesirable effect. But it has to be noted that the impact on classification is not as severe as it may seem on the first sight; typically the feature corresponding to such an artifact CSP filter component gets a near-zero weight in the classification step and is thereby neglected.

Finally, we would like to remark that the evaluation of CSP-based algorithms needs to take into account that this technique uses label information. This means that CSP filters may only be calculated from training data (of course the resulting filters need then to be applied also to the test set). In a cross-validation, CSP filters have to be calculated repeatedly on the training set within each fold/repetition. Otherwise severe underestimation of the generalization error may occur.

APPLICATION OF CSP TO SOURCE PROJECTION

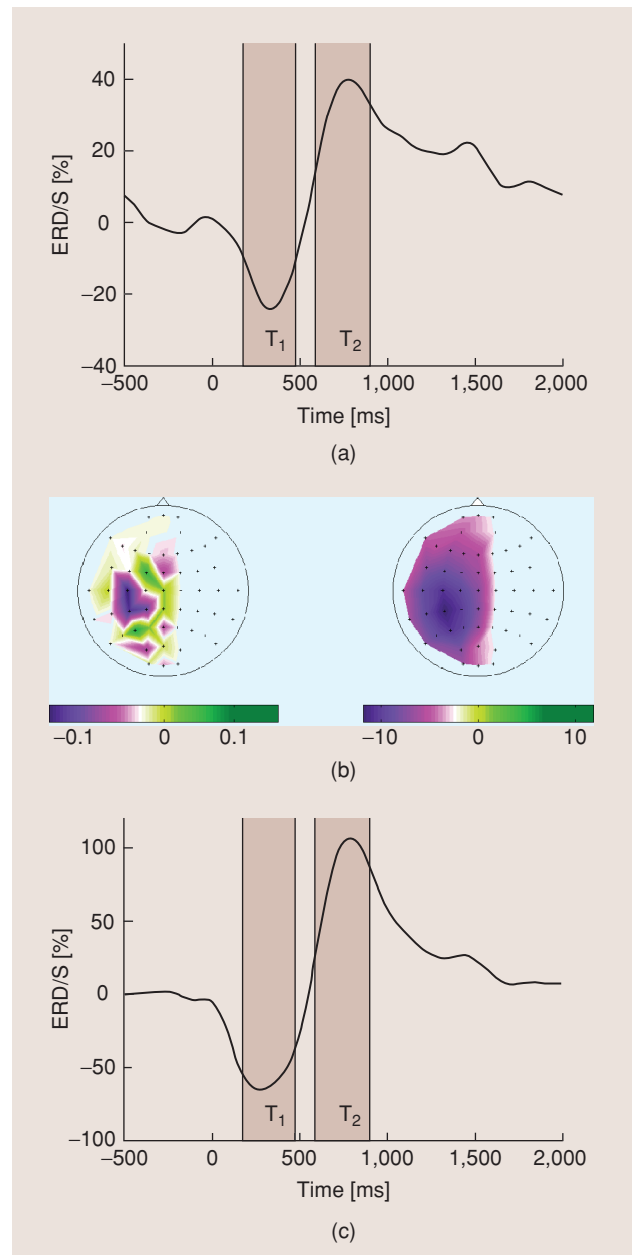
Here we report a novel application of CSP with a different flavor than above. Instead of single trial classification of mental states, CSP is used in the analysis of event-related modulations of brain rhythms. We show that CSP can be used to enhance the signal of interest while suppressing the background activity.

Conventionally event-related (de-)synchronization is defined as the relative difference in signal power of a certain frequency band, between two conditions, for instance a prestimulus or reference period and an immediate post-stimulus period [42]:

$$\text{ERD}(t) = \frac{\text{Power}(t) - \text{Reference power}}{\text{Reference power}}.$$

Thus ERD and ERS describe the relative power modulation of the ongoing activity, induced by a certain stimulus or event. Typically the sensor (possibly after Laplace filtering) that exhibit the strongest ERD/ERS effect at a certain frequency band is used for the analysis. Nevertheless, the CSP technique can help to further improve on the signal-to-noise ratio by optimizing the spatial filters focusing on rhythmic cortical generators that undergo the rhythmic perturbation.

We briefly outline how the CSP algorithm can be used for this purpose in an illustrative example of somatosensory stimulation. In particular, we use single trial EEG recordings of electrical stimulations of the median nerve at the right wrist. Such somatosensory stimulation typically causes modulations of the μ -rhythm, yielding a sequence of ERD followed by a rebound (ERS), overshooting the pre-event baseline level. Figure 9(a) depicts the time course of the averaged ERD/ERS for the α -band at approximately 10 HZ obtained from the best sensor. Based on this averaged band power modulations, we determine two disjoint temporal intervals T_1 and T_2 , associated with the desynchronization and the hyper-synchronization phase, respectively. These two intervals serve as the opposed conditions (classes) in the conventional CSP framework. We estimate covariance matrices $\Sigma^{(+)}$ and $\Sigma^{(-)}$ as in (3) pooling covariance matrices in the two intervals separately. Solving the CSP problem according to (5), yields a set of spatial filters. The filter that minimized the variance for the desynchronization period, while simultaneously maximizing those of the synchronization period constitutes the optimal spatial projection onto the cortical generator under consideration, i.e., onto the contralateral μ -rhythm. Here we restrict our CSP analysis only to the hemisphere that is contralateral to the stimulation in



[FIG9] Illustration of an improved source projection using the CSP technique. (a) The time course of the averaged band-power (10 HZ) at the channel (CP3) with the most prominent ERD/ERS following a median nerve stimulation at the right wrist. The gray-shaded areas indicate the two selected virtual classes for the CSP-algorithm, where T_1 corresponds to the ERD phase, while T_2 reflects the ERS interval. (b) Depicts the CSP-filter that minimizes the variance for T_1 , along with the projection of the corresponding source to the scalp. See main text for the reason to constrain the filter to the left hemisphere. (c) Time course of the averaged band-power of the projected signal. Note that this source projection procedure has yielded ERD and ERS that are much more accentuated as they have almost tripled in magnitude.

order to obtain unilateral spatial filter that has no cross talk with the other hemisphere. Figure 9 depicts the obtained spatial CSP filter, along the time course of ERD/ERS of the projected signal.

Note, in case the modulation of rhythmic activity comprises only of an ERD or an ERS response, the same approach can be used by simply contrasting a pre-stimulus reference interval against the period of modulation. In other words, CSP should be thought as a general tool for contrasting different brain states that yields a spatial filter solution that can be used to enhance the signal-to-noise ratio and can be interpreted from the physiological viewpoint.

CSP FILTERS MAXIMIZE THE VARIANCE OF THE SPATIALLY FILTERED SIGNAL UNDER ONE CONDITION WHILE MINIMIZING IT FOR THE OTHER CONDITION.

VARIANTS AND EXTENSIONS OF THE ORIGINAL CSP ALGORITHM

MULTICLASS

In its original form, CSP is restricted to binary problems. A general way to extend this algorithm to the multiclass case is to apply CSP to a set of binary subproblems (all binary pairs or, preferably, in a one-versus-rest scheme). A more direct approach by approximate simultaneous diagonalization was proposed in [12].

AUTOMATIC SELECTION OF SPECTRAL FILTER

The common spatio-spectral pattern (CSSP) algorithm [31] solves the standard CSP problem on the EEG time series augmented by delayed copies of the original signal, thereby obtaining simultaneously optimized spatial filters in conjunction with simple frequency filters. More specifically, CSP is applied to the original x concatenated with its off τ ms delayed version $x(t - \tau)$. This amounts to an optimization in an extended spatial domain, where the delayed signals are treated as new channels $\tilde{x}(t) = (x(t)^\top, x(t - \tau)^\top)^\top$. Consequently this yields spatial projections $\tilde{w} = (w^{(0)\top}, w^{(\tau)\top})^\top$, that correspond to vectors in this extended spatial domain. Any spatial projection in state space can be expressed as a combination of a pure spatial and spectral filter applied to the original data x , as follows:

$$\begin{aligned} \tilde{w}^\top \tilde{x}(t) &= \sum_{c=1}^C w_c^{(0)} x_c(t) + w_c^{(\tau)} x_c(t - \tau) \\ &= \sum_{c=1}^C \gamma_c \left(\frac{w_c^{(0)}}{\gamma_c} x_c(t) + \frac{w_c^{(\tau)}}{\gamma_c} x_c(t - \tau) \right), \end{aligned} \quad (8)$$

where $\{\gamma_c\}_{c=1}^C$ defines a pure spatial filter, whereas

$$\left(\frac{w_c^{(0)}}{\gamma_c}, \overbrace{0, \dots, 0}^{\tau-1}, \frac{w_c^{(\tau)}}{\gamma_c} \right)$$

defines a finite impulse response (FIR) filter at each electrode c . Accordingly this technique automatically neglects or emphasizes specific frequency bands at each electrode position in a way

that is optimal for the discrimination of two given classes of signals. Note that individual temporal filters are determined for each input channel.

The common sparse spectral spatial pattern (CSSSP) algorithm [13] eludes the problem of manually selecting the frequency band in a different way. Here a temporal FIR filter is optimized simultaneously with a spatial filter. In contrast to CSSP only one temporal filter is

used, but this filter can be of higher complexity. In order to control the complexity of the temporal filter, a regularization scheme is introduced which favors sparse solutions for the FIR coefficients. Although some values of the regularization parameter seem to give good results in most cases, for optimal performance a model selection has to be performed.

In [50] an iterative method (SPEC-CSP) is proposed which alternates between spatial filter optimization in the CSP sense and the optimization of a spectral weighting. As result, one obtains a spatial decomposition and a temporal filter with are jointly optimized for the given classification problem.

CONNECTION TO A DISCRIMINATIVE MODEL

Here we show how CSP analysis is related to a discriminative model. This connection is of theoretical interest in itself, and can also be used to further elaborate new variants of CSP. See [48], [49] for related models.

The quantity $S_d = \Sigma^{(+)} - \Sigma^{(-)}$ in (6) can be interpreted as the empirical average $\hat{\mathbb{E}}_{X,y}[yXX^\top]$ of the sufficient statistics yXX^\top of a linear logistic regression model:

$$P(y|X, V, b) = \frac{\exp(yf(X; V, b))}{Z(X, V, b)}$$

$$f(X; V, b) = \text{Tr}[V^\top XX^\top] + b,$$

where $y \in \{+1, -1\}$ is the label corresponding to two classes, $V \in \mathbb{R}^{C \times C}$ is the regression coefficient, b is the bias, and $Z(X, V, b) = e^{f(X; V, b)} + e^{-f(X; V, b)}$. In fact, given a set of trials and labels $\{X_i, y_i\}$ the log-likelihood of the above problem can be written as follows:

$$\begin{aligned} \log \prod_{i=1}^n P(y_i | X_i, V, b) &= \text{Tr} \left[V^\top \left(\sum_{i=1}^n y_i X_i X_i^\top \right) \right] \\ &\quad + b \sum_{i=1}^n y_i - \sum_{i=1}^n \log Z(X_i, V, b) \\ &= \frac{n}{2} \text{Tr} [V^\top S_d] - \sum_{i=1}^n \log Z(X_i, V, b), \end{aligned}$$

where for simplicity we assumed that each condition contains equal number ($n/2$) of trials. Unfortunately, because of the log-normalization $Z(X, V, b)$ term, the maximum likelihood prob-

lem cannot be solved as simple as the simultaneous diagonalization. One can upper bound the $\log Z(X, V, b)$ under the following condition:

$$\sum_{i=1}^n |\text{Tr}[V^\top X_i X_i^\top]| \leq 1,$$

and maximize the lower bound of the likelihood as follows:

$$\begin{aligned} & \underset{V \in \mathbb{R}^{c \times c}}{\text{maximize}} && \frac{n}{2} \text{Tr}[V^\top S_d], \\ & \text{subject to} && \sum_{i=1}^n |\text{Tr}[V^\top X_i X_i^\top]| \leq 1. \end{aligned}$$

Indeed this yields the first generalized eigenvector of the CSP problem (5) when V is rank = 1 matrix $V = ww^\top$.

REGULARIZING CSP

In practical BCI applications, the smaller the number of electrodes, the smaller the effort and time to set up the cap and also the smaller the stress of patients would be. CSP analysis can be used to determine where the electrodes should be positioned; therefore it would be still useful for experiments with a small number of electrodes. In [16], ℓ_1 regularization on the CSP filter coefficients was proposed to enforce a sparse solution; that is, many filter coefficients become numerically zero at the optimum. Therefore, it provides a clean way of selecting the number and the positions of electrodes. Their results have shown that the number of electrodes can be reduced to 10–20 without significant drop in the performance.

ADVANCED TECHNIQUES TOWARDS REDUCING CALIBRATION DATA

Because there exists substantial day-to-day variability in EEG data, the calibration session (15–35 min) is conventionally carried out every time before day-long experiments even for an experienced subject. Thus, in order to increase the usability of BCI systems, it is desirable to make use of previous recordings so that we can reduce the calibration measurement as small as possible (compare also data set IVa of the BCI competition III, [8]). For experienced BCI users whose EEG data were recorded more than once, [28] proposed a procedure to utilize results from the past recordings. They extracted prototypical filters by a clustering algorithm from the data recorded before and use them as an additional prior information for the current new session learning problem.

Recently [32] proposed an extended EM algorithm, where the extraction and classification of CSP features are performed jointly and iteratively. This method can be applied to the cases where either only a small number of calibration measurements (semisupervised) or even no labeled trials (unsupervised) are available. Basically, their algorithm repeats the following steps until a stable result is obtained: (i) constructing an expanded training data which consists of calibration trials with observed

labels and a part of unlabeled (feedback) data with labels estimated by the current classifier, (ii) reextracting the CSP feature and updating the classifier based on the current data sets. They analyzed the data IVa of BCI competition III [8] and reported that because of the iterative reextraction of the CSP features, they could achieve satisfactory performance from only 30 labeled and 120 unlabeled data or even from 150 unlabeled trials (off-line analysis). Note that only results of selected subjects of the competition data set IVa were reported. Although there was no experimental result presented, it was claimed that the extended EM procedure can also adapt to nonstationarity in EEG signals.

DEALING WITH THE NONSTATIONARY OF EEG SIGNALS

Another practical issue is nonstationarity in EEG data. There are various suggestions how to handle the nonstationarity in BCI systems [10], [26], [51], and [52]. With respect to CSP-based BCIs, the result of [29], [47] was that a simple adaptation of the classifier bias can compensate nonstationarity astonishingly well. Further changes like retraining LDA and recalculating CSP contributed only slightly or sometimes increased the error rate.

The question whether the CSP filter W or the pattern A should generalize to a new recording was raised by [23]. From a source separation point of view, the j -th column w_j of the filter W tries to capture the j -th source denoted by the j -th column a_j of the pattern A while trying to suppress all other sources that are irrelevant to the motor-imagination task. Therefore, if the disturbances change while the relevant source remains unchanged the optimal filter should adaptively change to cancel out the new disturbances while still capturing the relevant source. In [23] the fixed spatial pattern (FSP) approach was proposed; that is to keep the spatial pattern of the relevant source, i.e., subset of the columns of A unchanged while changing the spatial filter adaptively in a new recording. The true labels (i.e., the actual intension of a subject) are not required when the FSP is applied because only the irrelevant sources, which are assumed to be common to two classes, are re-estimated.

A novel approach to make CSP more robust to nonstationarities during BCI feedback was proposed in [5]. In this article, a short measurement of nontask related disturbances is used to enforce spatial filters which are invariant against those disturbances. In invariant CSP (iCSP) the covariance matrix of the disturbance is added to the denominator in the Rayleigh coefficient representation of CSP, compare (7).

CONCLUDING DISCUSSION

We have reviewed a spatial filtering technique that often finds its successful use in BCI: CSP. The method is based on the second order statistics of the signal between electrodes and the solution is obtained by solving a generalized eigenvalue problem. We have shown a generative and a discriminative interpretation of the method. We have applied the method to two motor imagination based BCI studies. In the

first study, we have reported the peak information transfer rate from one subject of 35.4 b/min. In the second study, we have shown that 12 out of 14 naive subjects could perform BCI control on their first BCI experiments. We have pointed out not only the advantage of the method, such as low computation cost and interpretability but also some caveats such as model selection and pre-processing issues or deterioration under outliers. We showed subsequently that CSP can be extended and robustified in order to alleviate these critical aspects. In this review, we have focused our attention to applications of CSP for single trial EEG analysis in the context of BCI. Note however that CSP-filtering and extensions thereof can be applied to extract general discriminative spatio-temporal structure from multivariate data streams beyond EEG. Future work will continue the quest to develop novel spatio-temporal filtering methods that allow more accurate and interpretable classification even for nonstationary, noisy, interacting data sources. Special attention will be placed on the construction of probabilistically interpretable nonlinear modeling that allows the integration of feature extraction and classification steps within a one step procedure in the spirit of, e.g., [17], [22], [48], and [49].

ACKNOWLEDGMENTS

This work was supported in part by grants of the Bundesministerium für Bildung und Forschung (BMBF), FKZ 01IBE01A (BCI III) and 01GQ0415 (BCCNB-A4), by MEXT, Grant-in-Aid for JSPS fellows, 17-11866 and by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778. This publication only reflects the authors' views.

We thank Vadim Nikulin, Gabriel Curio, Guido Dornhege, Matthias Krauledat, and Michael Tangermann for helpful discussions.

APPENDIX

HOW TO SELECT HYPERPARAMETERS FOR CSP

Here we give a heuristic procedure to automatically select all parameters that are needed for successful CSP application. There is no claim whatsoever that these heuristics are close to optimal or natural in any sense. However, we have found them practically working and evaluate them here in comparison to the general setting and to manual selection by the experimenter.

SELECTION OF A FREQUENCY BAND

We provide our heuristic for the selection of a discriminative frequency band in pseudo code, see Algorithm 1. The EEG trials X should be spatially filtered by a Laplacian or bipolar filter. In our experience, the algorithm works best if only few channels are used. A good choice is, e.g., to choose $C = \{c_1, c_2, c_3\}$ with c_i being one from each area of the left hand, right hand and feet with $\max \sqrt{\sum_f (\text{score}_c(f))^2}$.

Algorithm 1 Selection of a discriminative frequency band.

Let $X_{(c,i)}$ denote trial i at channel c with label y_i and let C denote the set of channels.

- 1) $\text{dB}_c(f, i) \leftarrow \log \text{band-power of } X_{(c,i)} \text{ at frequency } f \text{ (from 5 to 35Hz)}$
- 2) $\text{score}_c(f) \leftarrow \text{corrcoef}(\text{dB}_c(f, i), y_i)_i$
- 3) $f_{\max} \leftarrow \arg\max_f \sum_{c \in C} \text{score}_c(f)$
- 4) $\text{score}_c^*(f) \leftarrow \begin{cases} \text{score}_c(f) & \text{if } \text{score}_c(f_{\max}) > 0 \\ -\text{score}_c(f) & \text{otherwise} \end{cases}$
- 5) $\text{fscore}(f) \leftarrow \sum_{c \in C} \text{score}_c^*(f)$
- 6) $f_{\max}^* \leftarrow \arg\max_f \text{fscore}(f)$
- 7) $f_0 \leftarrow f_{\max}^* \quad f_1 \leftarrow f_{\max}^*$
- 8) **while** $\text{fscore}(f_0 - 1) \geq \text{fscore}(f_{\max}^*) * 0.05$ **do**
- 9) $f_0 \leftarrow f_0 - 1$
- 10) **while** $\text{fscore}(f_1 + 1) \geq \text{fscore}(f_{\max}^*) * 0.05$ **do**
- 11) $f_1 \leftarrow f_1 + 1$
- 12) **return** frequency band $[f_0, f_1]$

SELECTION OF A TIME INTERVAL.

The heuristic selection of a time interval proceeds similar to the selection of the frequency band, (see e.g., 2).

Algorithm 2 Selection of a discriminative time interval.

Let $X_{(c,i)(t)}$ denote time sample t of trial i at channel c with label y_i and let C denote the set of channels.

- 1) $\text{env}_c(t, i) \leftarrow \text{envelope of } X_{(c,i)(t)}, \text{ calculated by Hilbert transform (e.g. [9]) and smoothed}$
- 2) $\text{score}_c(t) \leftarrow \text{corrcoef}(\text{env}_c(t, i), y_i)_i$
- 3) $t_{\max} \leftarrow \arg\max_t \sum_{c \in C} |\text{score}_c(t)|$
- 4) $\text{score}_c^*(t) \leftarrow \begin{cases} \text{score}_c(t) & \text{if } \sum_{t_{\max}-100\text{ms} < t' < t_{\max}+100\text{ms}} \text{score}_c(t') > 0 \\ -\text{score}_c(t) & \text{otherwise} \end{cases}$
- 5) $\text{tscore}(t) \leftarrow \sum_{c \in C} \text{score}_c^*(t)$
- 6) $t_{\max}^* \leftarrow \arg\max_t \text{tscore}(t)$
- 7) $\text{thresh} \leftarrow 0.8 * \sum_t \text{tscore}^+(t)$ (with $f^+(x) = f(x)$ if $f(x) > 0$ and $= 0$ otherwise)
- 8) $t_0 \leftarrow t_{\max}^*; \quad t_1 \leftarrow t_{\max}^*$
- 9) **while** $\sum_{t_0 \leq t \leq t_1} \text{tscore}(t) < \text{thresh}$ **do**
- 10) **if** $\sum_{t < t_0} \text{tscore}^*(t) > \sum_{t > t_1} \text{tscore}^*(t)$ **then**
- 11) $t_0 \leftarrow t_0 - 1$
- 12) **else**
- 13) $t_1 \leftarrow t_1 + 1$
- 14) **return** time interval $[t_0, t_1]$

SELECTION OF A SUBSET OF FILTERS

The classical measure for the selection of CSP filters is based on the eigenvalues in (5). Each eigenvalue is the relative variance of the signal filtered with the corresponding spatial filter (variance in one condition divided by the sum of variances in both condi-

tions). This measure is not robust to outliers because it is based on simply pooling the covariance matrices in each condition (3). In fact, one single trial with very high variance can have a strong impact on the CSP solution (see also Figure 8). A simple way to circumvent this problem is to calculate the variance of the filtered signal within each trial and then calculate the corresponding ratio of medians:

$$\text{score}(w_j) = \frac{\text{med}_j^{(+)}}{\text{med}_j^{(+)} + \text{med}_j^{(-)}}$$

where,

$$\text{med}_j^{(c)} = \text{median}_{i \in \mathcal{I}_c} (w_j^\top X_i X_i^\top w_j) \quad (c \in \{+, -\}).$$

As with eigenvalues, a ‘ratio-of-medians’ score near 1 or near 0 indicates good discriminability of the corresponding spatial filter. These scores are more robust with respect to outliers than the eigenvalue score, e.g., the filter shown in Figure 8 would get a minor (i.e., near 0.5) ratio-of-medians score.

EVALUATION OF HEURISTIC SELECTION PROCEDURE

Here we compare the impact of individually choosing the hyperparameters for CSP-based classification. We compare the method “fixed” which uses a broad frequency band 7–30 Hz and the time window 1,000–3,500 ms post stimulus. The method “auto” uses the heuristics presented in this section to select frequency band and time interval. In “manual” we use the settings that were chosen by an experienced experimenter by hand for the actual feedback (see [15] for a practical example with manual selection). Note there is a substantial improvement of performance in most of the data sets. Interestingly in one feedback data set (subject ct) the “auto” method performs badly, although the selected parameters were reasonable.

AUTHORS

Benjamin Blankertz (blanker@cs.tu-berlin.de) received the Diploma degree in mathematics 1994 and the Ph.D. in mathematical logic in 1997, both from University of Muenster, Germany. He conducted studies in computational models for perception of music and computer-aided music analysis. Since 2000, he is with the Intelligent Data Analysis (IDA) group at Fraunhofer FIRST, and since 2007 with Technical University (TU) of Berlin working in the BBCI project. His scientific interests are in the fields of machine learning, analysis of biomedical data, and psychoacoustics.

Ryota Tomioka (ryota.tomioka@first.fraunhofer.de) received his bachelor’s degree from the University of Tokyo, School of Engineering in 2003, and master’s degree from the Graduate School of Frontier Sciences, 2005. He was awarded a Japan Society for the Promotion of Science (JSPS) fellowship for young scientists in 2005, and he works on machine learning and signal processing techniques for BCI at the Fraunhofer Institute FIRST in Berlin.

Steven Lemm (steven.lemm@first.fraunhofer.de) received the Diploma degree in mathematics in 2002 from the Berlin University of Technology, Berlin. He conducted studies in probabilistic modeling of computer networks and communication protocols in order to optimize the expected network throughput. Since 1998, he has been with the IDA group at the Fraunhofer-Institute FIRST, and also as of 2002, he has been with the Neurophysics Group at the Department of Neurology of the Campus Benjamin Franklin, Charité University Medicine, Berlin. His main scientific interests are in the field of machine learning, with a special focus on the analysis of biomedical data.

Motoaki Kawanabe (motoaki.kawanabe@first.fraunhofer.de) graduated from the University of Tokyo in 1990. He received the Ph.D. on semiparametric statistics and information geometry in 1995. He later worked as an associate professor at the statistics laboratory of the same department. In November 2000, he moved to the IDA group of Fraunhofer Institute FIRST. His research interests are robust statistics, signal processing, and kernel methods in machine learning. He started collaboration with the BBCI team in 2005, and in particular has worked on nonstationarity of EEG signals.

Klaus-Robert Müller (krm@cs.tu-berlin.de) received the Diplom degree in mathematical physics 1989 and the Ph.D. in theoretical computer science in 1992, both from University of Karlsruhe, Germany. From 1992–1994 he worked as a postdoctoral fellow at GMD FIRST, in Berlin where he started to build up the IDA group. From 1994–1995 he was a European Community STP Research Fellow at University of Tokyo. In 1999, he received the Olympus Prize, awarded by the German Pattern Recognition Society, DAGM. In 2003, he became full professor at University of Potsdam in Germany. He became chair of the machine learning department at TU Berlin in 2006. He was honored in 2006 with the SEL Alcatel Communication Award. His research areas include statistical learning theory for neural networks, support vector machines, and ensemble learning techniques. His present application interests are expanded to the analysis of biomedical data, most recently to brain computer interfacing and genomic data analysis.

REFERENCES

- [1] H. Berger, “Über das Elektroenkephalogramm des Menschen,” *Arch. Psychiat. Nervenkr.*, vol. 99, no. 6, pp. 555–574, 1933.
- [2] B. Blankertz, G. Dornhege, M. Krauledat, G. Curio, and K.-R. Müller, “The non-invasive Berlin Brain-Computer Interface: Fast acquisition of effective performance in untrained subjects,” *NeuroImage*, vol. 37, no. 2, pp. 539–550, 2007.
- [3] B. Blankertz, G. Dornhege, M. Krauledat, K.-R. Müller, V. Kunzmann, F. Losch, and G. Curio, “The Berlin Brain-Computer Interface: EEG-based communication without subject training,” *IEEE Trans. Neural Sys. Rehab. Eng.*, vol. 14, no. 2, pp. 147–152, 2006.
- [4] B. Blankertz, G. Dornhege, S. Lemm, M. Krauledat, G. Curio, and K.-R. Müller, “The Berlin Brain-Computer Interface: Machine learning based detection of user specific brain states,” *J. Univ. Comput. Sci.*, vol. 12, no. 6, pp. 581–607, 2006.
- [5] B. Blankertz, M. Kawanabe, R. Tomioka, F. Hohlefeld, V. Nikulin, and K.-R. Müller, “Invariant common spatial patterns: Alleviating nonstationarities in brain-computer interfacing,” in *Advances in Neural Information Processing Systems 20*, Cambridge, MA: MIT Press, 2008.
- [6] B. Blankertz, M. Krauledat, G. Dornhege, J. Williamson, R. Murray-Smith, and K.-R. Müller, “A note on brain actuated spelling with the Berlin Brain-Computer Interface,” in *Universal Access in HCI*, C. Stephanidis, Ed. (Part II, *HCI*), 2007, vol. 4555 of *LNCS* Berlin: Springer-Verlag, vol. 4555, pp. 759–768, 2007.

- [7] B. Blankertz, F. Losch, M. Krauledat, G. Dornhege, G. Curio, and K.-R. Müller, "The Berlin Brain-Computer Interface: Accurate performance from first-session in BCI-naïve subjects," submitted for publication.
- [8] B. Blankertz, K.-R. Müller, D. Krusienski, G. Schalk, J.R. Wolpaw, A. Schlögl, G. Pfurtscheller, J. del R. Millán, M. Schröder, and N. Birbaumer, "The BCI competition III: Validating alternative approaches to actual BCI problems," *IEEE Trans. Neural Syst. Rehab. Eng.*, vol. 14, no. 2, pp. 153–159, 2006.
- [9] R.N. Bracewell, *The Fourier Transform and Its Applications*, 3rd ed. New York: McGraw-Hill, 1999.
- [10] J. del R. Millán, "On the need for on-line learning in brain-computer interfaces," in *Proc. Int. Joint Conference Neural Networks*, Budapest, Hungary, July 2004. IDIAP-RR 03-30, pp. 2877–2882.
- [11] G. Dornhege, "Increasing information transfer rates for brain-computer interfacing. Ph.D. dissertation, Dept. Comput. Sci., Univ. Potsdam, Germany, 2006.
- [12] G. Dornhege, B. Blankertz, G. Curio, and K.-R. Müller, "Boosting bit rates in non-invasive EEG single-trial classifications by feature combination and multi-class paradigms," *IEEE Trans. Biomed. Eng.*, vol. 51, no. 6, pp. 993–1002, June 2004.
- [13] G. Dornhege, B. Blankertz, M. Krauledat, F. Losch, G. Curio, and K.-R. Müller, "Optimizing spatio-temporal filters for improving brain-computer interfacing," in *Advances in Neural Inf. Proc. Systems (NIPS 05)*, vol. 18, pp. 315–322, Cambridge, MA, 2006. MIT Press.
- [14] G. Dornhege, José del R. Millán, Thilo Hinterberger, Dennis McFarland, and K.-R. Müller, Eds. *Toward Brain-Computer Interfacing*. Cambridge, MA: MIT Press, 2007.
- [15] G. Dornhege, M. Krauledat, K.-R. Müller, and B. Blankertz, "General signal processing and machine learning tools for BCI," in *Toward Brain-Computer Interfacing*. Cambridge, MA: MIT Press, 2007, pp. 207–233.
- [16] J. Farquhar, N.J. Hill, T.N. Lal, and B. Schölkopf, "Regularised CSP for sensor selection in BCI," in *Proc. 3rd Int. Brain-Computer Interface Workshop Training Course 2006*, Verlag der Technischen Universität Graz, Graz, Austria, pp. 14–15.
- [17] J. Farquhar, J. Hill, and B. Schölkopf, "Learning optimal EEG features across time, frequency and space," in *NIPS 2006 Workshop Current Trends Brain-Computer Interfacing*, Whistler, Canada, 2006 [Online]. Available: http://www.kyb.mpg.de/publications/attachments/nips2006_4262%5B0%5D.pdf
- [18] K. Fukunaga, *Introduction to Statistical Pattern Recognition*, 2nd ed. New York: Academic, 1990.
- [19] C. Guger, H. Ramoser, and G. Pfurtscheller, "Real-time EEG analysis with subject-specific spatial patterns for a Brain Computer Interface (BCI)," *IEEE Trans. Neural Syst. Rehab. Eng.*, vol. 8, no. 4, pp. 447–456, 2000.
- [20] R. Hari and R. Salmelin, "Human cortical oscillations: A neuromagnetic view through the skull," *Trends Neurosci.*, vol. 20, no. 1, pp. 44–49, 1997.
- [21] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference and Prediction* (Springer Series in Statistics). New York: Springer-Verlag, 2001.
- [22] J. Hill and J. Farquhar, "An evidence-based approach to optimizing feature extraction in EEG signal classification," Max Planck Inst. Biolog. Cybernet., Tech. Report, 2007, unpublished.
- [23] N.J. Hill, J. Farquhar, T.N. Lal, and B. Schölkopf, "Time-dependent demixing of task-relevant EEG signals," in *Proc. 3rd Int. Brain-Computer Interface Workshop Training Course 2006*, Verlag der Technischen Universität Graz, Graz, Austria, pp. 20–21.
- [24] H. Jasper and H.L. Andrews, "Normal differentiation of occipital and precentral regions in man," *Arch. Neurol. Psych. (Chicago)*, vol. 39, pp. 96–115, Jan. 1938.
- [25] H. Jasper and W. Penfield, "Electrocorticograms in man: Effect of voluntary movement upon the electrical activity of the precentral gyrus," *Arch. Psychiatric Zeitschrift Neurol.*, vol. 183, pp. 163–174, 1949.
- [26] M. Kawanabe, M. Krauledat, and B. Blankertz, "A Bayesian approach for adaptive BCI classification," in *Proc. 3rd Int. Brain-Computer Interface Workshop Training Course 2006*, pp. 54–55. Verlag der Technischen Universität Graz, Graz, Austria, 2006.
- [27] Z.J. Koles, "The quantitative extraction and topographic mapping of the abnormal components in the clinical EEG," *Electroencephalogr. Clin. Neurophysiol.*, vol. 79, no. 6, pp. 440–447, 1991.
- [28] M. Krauledat, M. Schröder, B. Blankertz, and K.-R. Müller, "Reducing calibration time for brain-computer interfaces: A clustering approach," in *Advances in Neural Information Processing Systems 19*, B. Schölkopf, J. Platt, and T. Hoffman, Eds. Cambridge, MA: MIT Press, 2007, pp. 753–760.
- [29] M. Krauledat, P. Shenoy, B. Blankertz, R.P.N. Rao, and K.-R. Müller, "Adaptation in CSP-based BCI systems," in *Toward Brain-Computer Interfacing*. Cambridge, MA: MIT Press, 2007, pp. 305–309.
- [30] A. Kübler, F. Nijboer, J. Mellinger, T.M. Vaughan, H. Pawelzik, G. Schalk, D.J. McFarland, N. Birbaumer, and J.R. Wolpaw, "Patients with ALS can use sensorimotor rhythms to operate a brain-computer interface," *Neurology*, vol. 64, no. 10, pp. 1775–1777, 2005.
- [31] S. Lemm, B. Blankertz, G. Curio, and K.-R. Müller, "Spatio-spectral filters for improving classification of single trial EEG," *IEEE Trans. Biomed. Eng.*, vol. 52, no. 9, pp. 1541–1548, 2005.
- [32] Y. Li and C. Guan, "An extended EM algorithm for joint feature extraction and classification in brain-computer interfaces," *Neural Comput.*, vol. 18, no. 11, pp. 2730–2761, 2006.
- [33] K.-R. Müller, S. Mika, G. Rätsch, K. Tsuda, and B. Schölkopf, "An introduction to kernel-based learning algorithms," *IEEE Trans. Neural Networks*, vol. 12, no. 2, pp. 181–201, May 2001.
- [34] K.-R. Müller and B. Blankertz, "Toward noninvasive brain-computer interfaces," *IEEE Signal Processing Mag.*, vol. 23, no. 5, pp. 125–128, Sept. 2006.
- [35] J. Müller-Gerking, G. Pfurtscheller, and H. Flyvbjerg, "Designing optimal spatial filters for single-trial EEG classification in a movement task," *Clin. Neurophysiol.*, vol. 110, no. 5, pp. 787–798, 1999.
- [36] C. Neuper and G. Pfurtscheller, "Event-related dynamics of cortical rhythms: frequency-specific features and functional correlates," *Int. J. Psychophysiol.*, vol. 43, no. 1, pp. 41–58, 2001.
- [37] C. Neuper, R. Scherer, M. Reiner, and G. Pfurtscheller, "Imagery of motor actions: Differential effects of kinesthetic and visual-motor mode of imagery in single-trial EEG," *Brain Res. Cogn. Brain Res.*, vol. 25, no. 3, pp. 668–677, 2005.
- [38] V. Nikouline, K. Linkenkaer-Hansen, H. Wikström, M. Kesäniemi, E. Antonova, R. Ilmoniemi, and J. Huttunen, "Dynamics of mu-rhythm suppression caused by median nerve stimulation: A magnetoencephalographic study in human subjects," *Neurosci. Lett.*, vol. 294, no. 3, pp. 163–166, 2000.
- [39] P.L. Nunez, R. Srinivasan, A.F. Westdorp, R.S. Wijesinghe, D.M. Tucker, R.B. Silberstein, and P.J. Cadusch, "EEG coherence I: Statistics, reference electrode, volume conduction, Laplacians, cortical imaging, and interpretation at multiple scales," *Electroencephalogr. Clin. Neurophysiol.*, vol. 103, no. 5, pp. 499–515, 1997.
- [40] L. Parra and P. Sajda, "Blind source separation via generalized eigenvalue decomposition," *J. Mach. Learn. Res.*, vol. 4, pp. 7–8, pp. 1261–1269, 2003.
- [41] L.C. Parra, C.D. Spence, A.D. Gerson, and P. Sajda, "Recipes for the linear analysis of EEG," *NeuroImage*, vol. 28, no. 2, pp. 326–341, 2005.
- [42] G. Pfurtscheller and A. Aranbizar, "Evaluation of event-related desynchronization preceding and following voluntary self-paced movement," *Electroencephalogr. Clin. Neurophysiol.*, vol. 46, no. 2, pp. 138–146, 1979.
- [43] G. Pfurtscheller, C. Brunner, A. Schlögl, and F.H. Lopes da Silva, "Mu rhythm (de)synchronization and EEG single-trial classification of different motor imagery tasks," *NeuroImage*, vol. 31, no. 1, pp. 153–159, 2006.
- [44] A. Schlögl, J. Kronegg, J. Huggins, and S.G. Mason, "Evaluation criteria for BCI research," in *Towards Brain-Computer Interfacing*, G. Dornhege, Jose del R. Millán, Thilo Hinterberger, Dennis McFarland, and K.-R. Müller, Eds. Cambridge, MA: MIT Press, 2007, pp. 297–312.
- [45] A. Schnitzler, S. Salenius, R. Salmelin, V. Joumäki, and R. Hari, "Involvement of primary motor cortex in motor imagery: A neuromagnetic study," *NeuroImage*, vol. 6, no. 3, pp. 201–208, 1997.
- [46] S.H. Scott, "Converting thoughts into action," *Nature*, vol. 442, no. 7099, pp. 141–142, 2006.
- [47] P. Shenoy, M. Krauledat, B. Blankertz, R.P.N. Rao, and K.-R. Müller, "Towards adaptive classification for BCI," *J. Neural Eng.*, vol. 3, no. 1, pp. R13–R23, 2006.
- [48] R. Tomioka and K. Aihara, "Classifying matrices with a spectral regularization," in *ICML '07: Proc. 24th Int. Conf. Machine Learning*, ACM Press, 2007, pp. 895–902.
- [49] R. Tomioka, K. Aihara, and K.-R. Müller, "Logistic regression for single trial EEG classification," in *Advances in Neural Information Processing Systems 19*, B. Schölkopf, J. Platt, and T. Hoffman, Eds. Cambridge, MA: MIT Press, 2007, pp. 1377–1384.
- [50] R. Tomioka, G. Dornhege, K. Aihara, and K.-R. Müller, "An iterative algorithm for spatio-temporal filter optimization," in *Proc. 3rd Int. Brain-Computer Interface Workshop Training Course 2006*, Verlag der Technischen Universität Graz, Graz, Austria, 2006, pp. 22–23.
- [51] C. Vidaurre, A. Schlögl, R. Cabeza, R. Scherer, and G. Pfurtscheller, "A fully on-line adaptive BCI," *IEEE Trans. Biomed. Eng.*, vol. 53, no. 6, 2006.
- [52] C. Vidaurre, A. Schlögl, R. Cabeza, R. Scherer, and G. Pfurtscheller, "Study of on-line adaptive discriminant analysis for EEG-based brain computer interfaces," *IEEE Trans. Biomed. Eng.*, vol. 54, no. 3, pp. 550–556, 2007.
- [53] J.R. Wolpaw and D.J. McFarland, "Control of a two-dimensional movement signal by a noninvasive brain-computer interface in humans," *Proc. Natl. Acad. Sci. USA*, vol. 101, no. 51, pp. 17849–17854, 2004.
- [54] J.R. Wolpaw, D.J. McFarland, and T.M. Vaughan, "Brain-computer interface research at the Wadsworth Center," *IEEE Trans. Rehab. Eng.*, vol. 8, no. 2, pp. 222–226, 2000.
- [55] J.R. Wolpaw, N. Birbaumer, D.J. McFarland, G. Pfurtscheller, and T.M. Vaughan, "Brain-computer interfaces for communication and control," *Clin. Neurophysiol.*, vol. 113, no. 6, pp. 767–791, 2002.