

New BCI approaches: Selective Attention to Auditory and Tactile Stimulus Streams

N. Jeremy Hill,¹ Cornelius Raths²

¹Max Planck Institute for Biological Cybernetics, Tübingen, Germany

²University of Southern California, Los Angeles, USA

with thanks to Ryota Tomioka at FIRST/TU-Berlin



Overview



Long-term goal: develop and optimize new paradigms that are suitable for people in the “completely locked-in” state (CLIS).



Overview



Long-term goal: develop and optimize new paradigms that are suitable for people in the “completely locked-in” state (CLIS).

- Enable the patient to make a binary decision.



Overview



Long-term goal: develop and optimize new paradigms that are suitable for people in the “completely locked-in” state (CLIS).

- Enable the patient to make a binary decision.
- Use signals elicited auditory and tactile stimuli, measureable by EEG.



Overview



Long-term goal: develop and optimize new paradigms that are suitable for people in the “completely locked-in” state (CLIS).

- Enable the patient to make a binary decision.
- Use signals elicited auditory and tactile stimuli, measureable by EEG.
- Use machine-learning algorithms to classify the signals.



Overview



Long-term goal: develop and optimize new paradigms that are suitable for people in the “completely locked-in” state (CLIS).

- Enable the patient to make a binary decision.
- Use signals elicited auditory and tactile stimuli, measureable by EEG.
- Use machine-learning algorithms to classify the signals.

This presentation concerns the results from preliminary experiments with healthy subjects.



Why Non-Visual?





Why Non-Visual?



- In the CLIS state, patients are functionally blind:



Why Non-Visual?



- In the CLIS state, patients are functionally blind:
 - eyes cannot be opened at will;
 - eyes may move involuntarily (often rolling up);
 - lens cannot be refocused or gaze directed;
 - no microsaccades, so images fade out (Troxler effect);
 - no saccades, so no integration of visual scenes: the fovea images a fixed 2 deg. spot, and resolution is very low in most of the visual field;
 - long immobility of the eye often leads to infections;



Why Non-Motor?



- Motor imagery-based BCI shows promising results with normal subjects, and patients with extensive paralysis (Kübler et al 2005, Neurology 10). So far it has not worked with patients in CLIS.

Why?

- Can the patient still imagine movement?



Why Non-Motor?



- Motor imagery-based BCI shows promising results with normal subjects, and patients with extensive paralysis (Kübler et al 2005, Neurology 10). So far it has not worked with patients in CLIS.

Why?

- Can the patient still imagine movement?
- Can the motor and premotor cortex still produce ERD/ERS during motor imagery?



Why Non-Motor?



- Motor imagery-based BCI shows promising results with normal subjects, and patients with extensive paralysis (Kübler et al 2005, Neurology 10). So far it has not worked with patients in CLIS.

Why?

- Can the patient still imagine movement?
- Can the motor and premotor cortex still produce ERD/ERS during motor imagery?
- (...and are these in fact the same question?)



Why Non-Motor?



- Motor imagery-based BCI shows promising results with normal subjects, and patients with extensive paralysis (Kübler et al 2005, Neurology 10). So far it has not worked with patients in CLIS.

Why?

- Can the patient still imagine movement?
- Can the motor and premotor cortex still produce ERD/ERS during motor imagery?
- (...and are these in fact the same question?)
- Are they still intact enough to (relearn to) do so?
 - ★ EEG is still the most attractive technology for clinical BCI.
 - ★ Most of the EEG signal comes from pyramidal neurons.
 - ★ ALS kills the pyramidal neurons of the motor cortex.



Auditory and Tactile BCI



Exogenous (i.e. stimulus-driven) BCI's rely on the conscious direction of the user's *attention*.



Auditory and Tactile BCI



Exogenous (i.e. stimulus-driven) BCI's rely on the conscious direction of the user's *attention*.

For paralysed users, this means *covert* attention.



Auditory and Tactile BCI



Exogenous (i.e. stimulus-driven) BCI's rely on the conscious direction of the user's *attention*.

For paralysed users, this means *covert* attention.

Does covert attention affect auditory ERPs?



Auditory and Tactile BCI



Exogenous (i.e. stimulus-driven) BCI's rely on the conscious direction of the user's *attention*.

For paralysed users, this means *covert* attention.

Does covert attention affect auditory ERPs?

Yes, e.g.:

- Hillyard et al. 1974, *Science* 182.
- Näätänen 1990, *Behavioral and Brain Sciences* 13.
- Schröger and Wolff 1998 *Cognitive Brain Research* 7.



Auditory and Tactile BCI



Exogenous (i.e. stimulus-driven) BCI's rely on the conscious direction of the user's *attention*.

For paralysed users, this means *covert* attention.

Does covert attention affect tactile ERPs?

Yes, e.g.:

- Desmet et al. 1977, *Journal of Physiology* 271.
- García Lorrea 1995, *Psychophysiology* 32.
- Eimer et al. 2003 *Experimental Brain Research* 151.



Auditory and Tactile BCI



Exogenous (i.e. stimulus-driven) BCI's rely on the conscious direction of the user's *attention*.

For paralysed users, this means *covert* attention.

Does covert attention affect tactile ERPs?

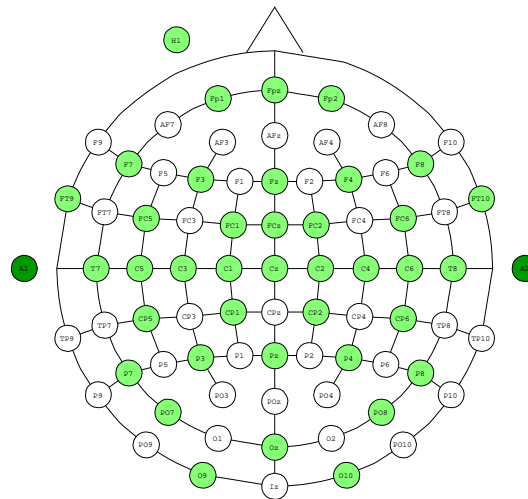
Yes, e.g.:

- Desmet et al. 1977, *Journal of Physiology* 271.
- García Lorrea 1995, *Psychophysiology* 32.
- Eimer et al. 2003 *Experimental Brain Research* 151.

...at least, when you average hundreds of trials. Can we obtain a reliable effect on a timescale suitable for BCI?

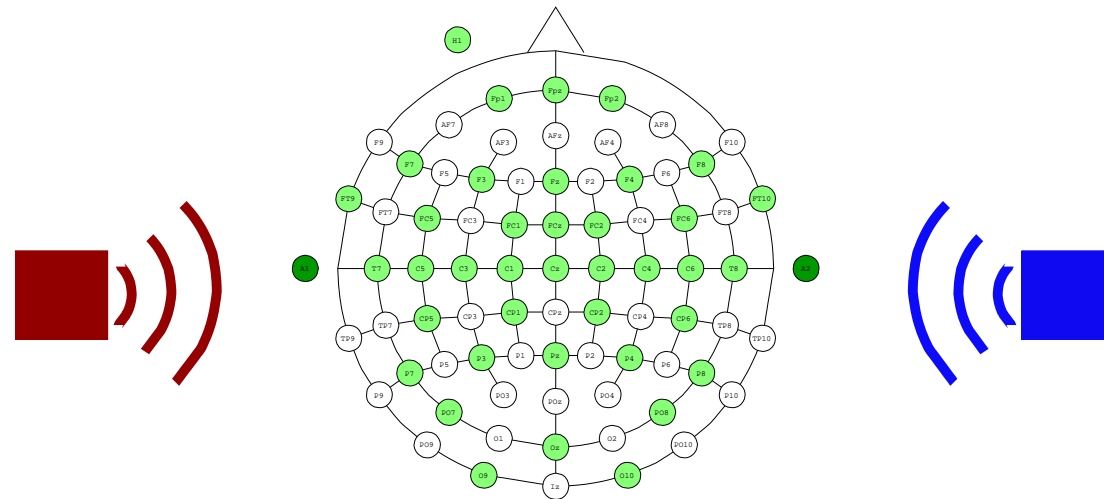


BIOLOGISCHE KYBERNETIK

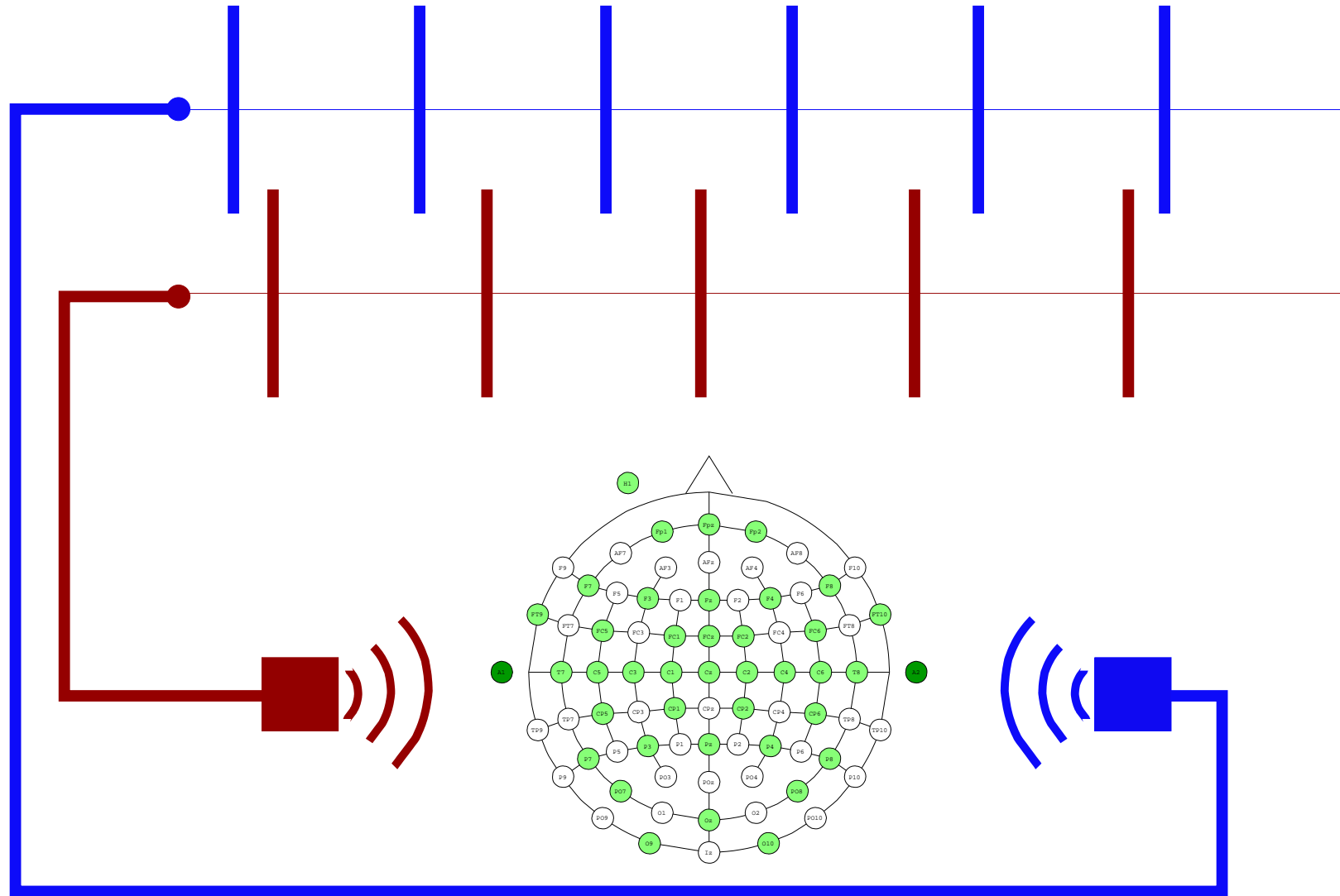




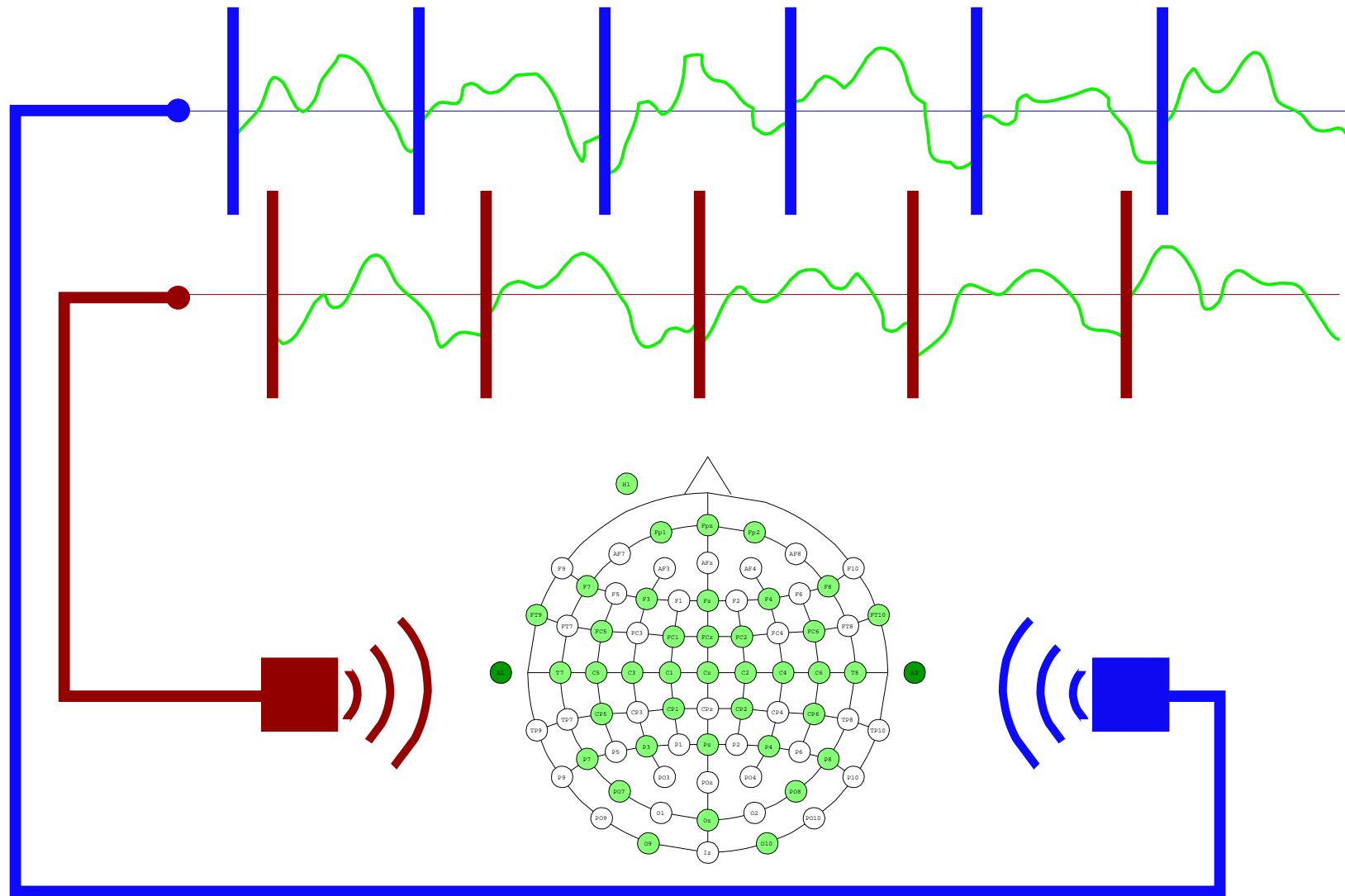
I: Auditory stimulation in EEG



I: Auditory stimulation in EEG



I: Auditory stimulation in EEG

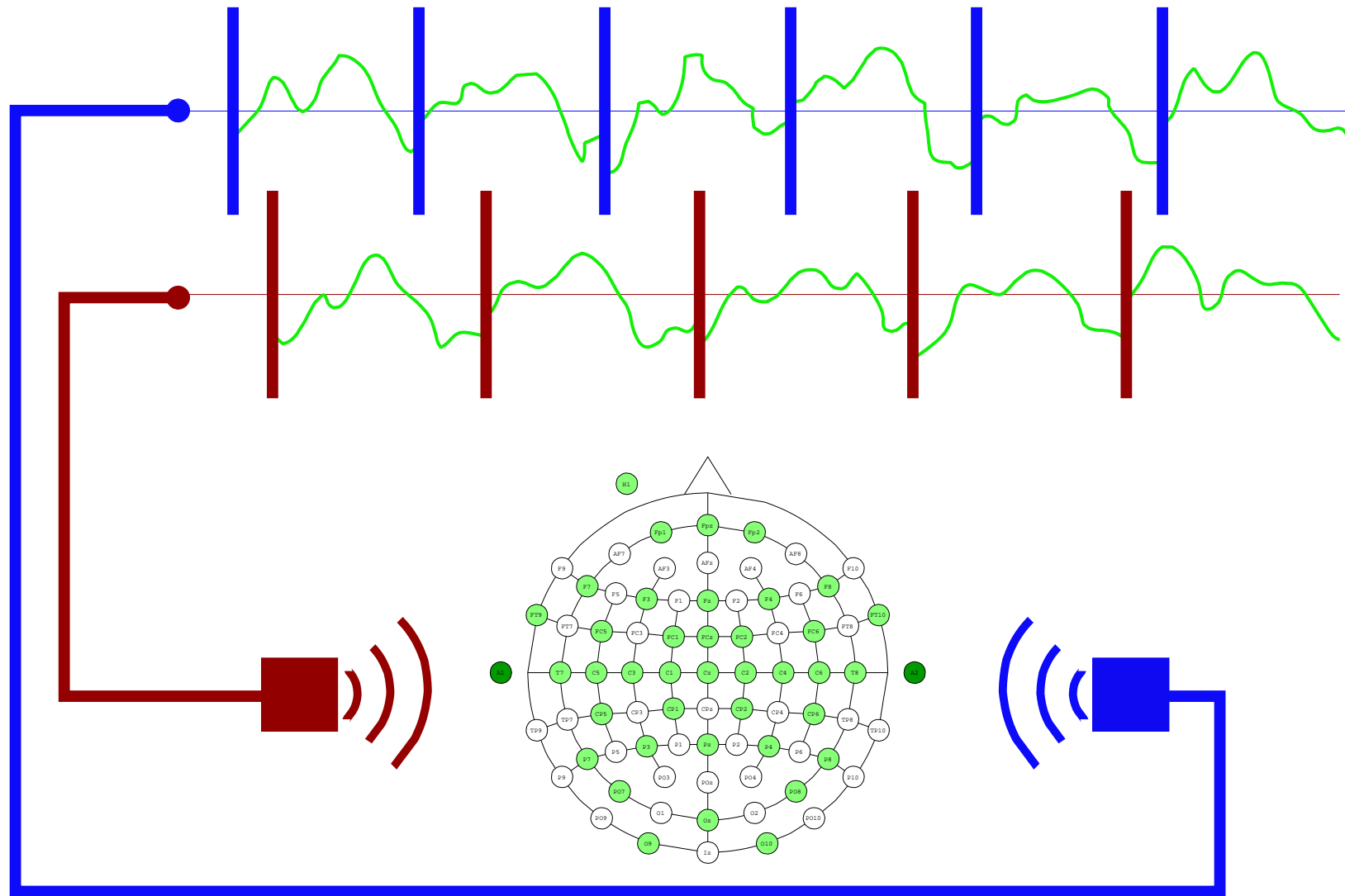




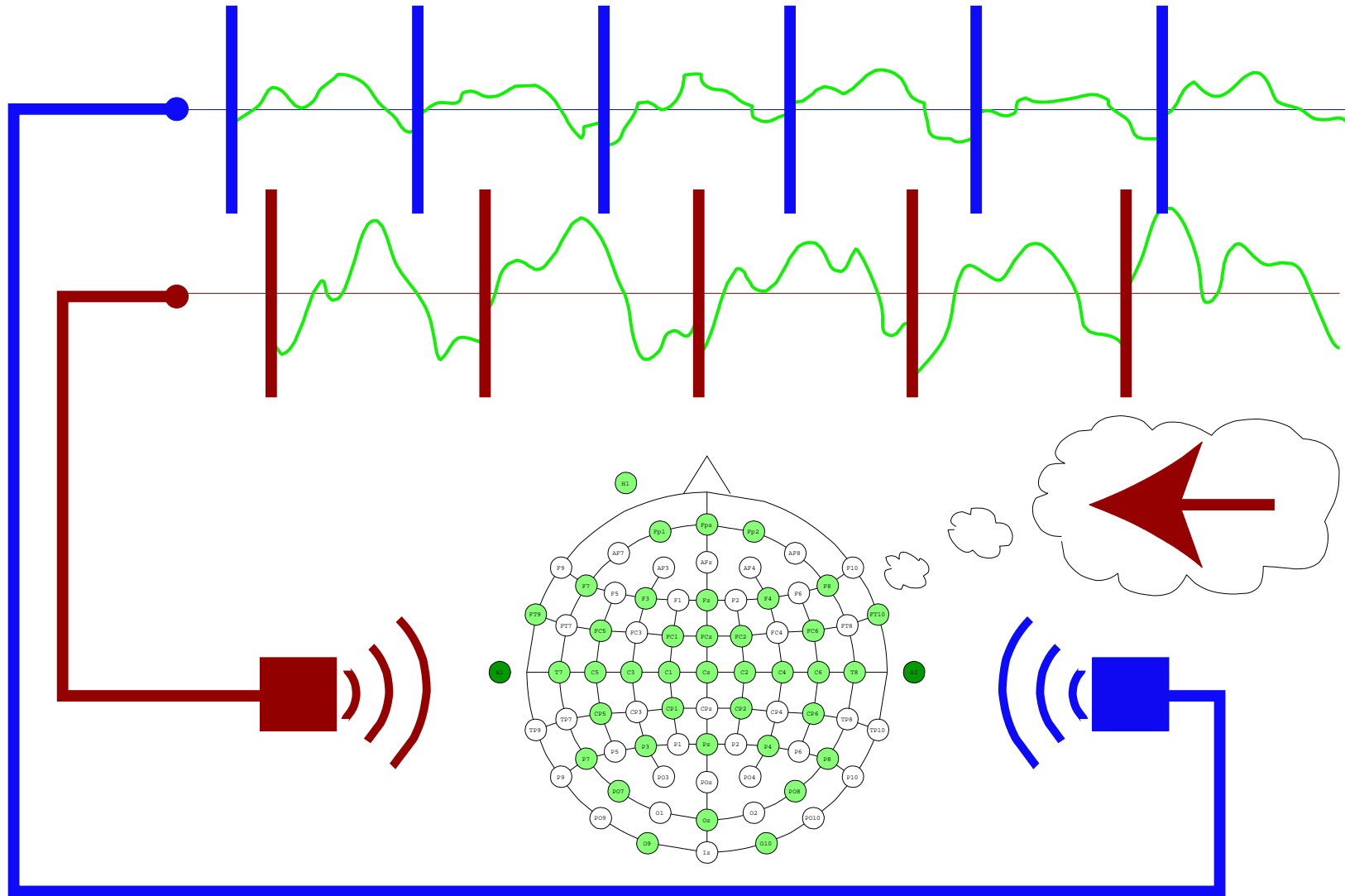
BIOLOGISCHE KYBERNETIK



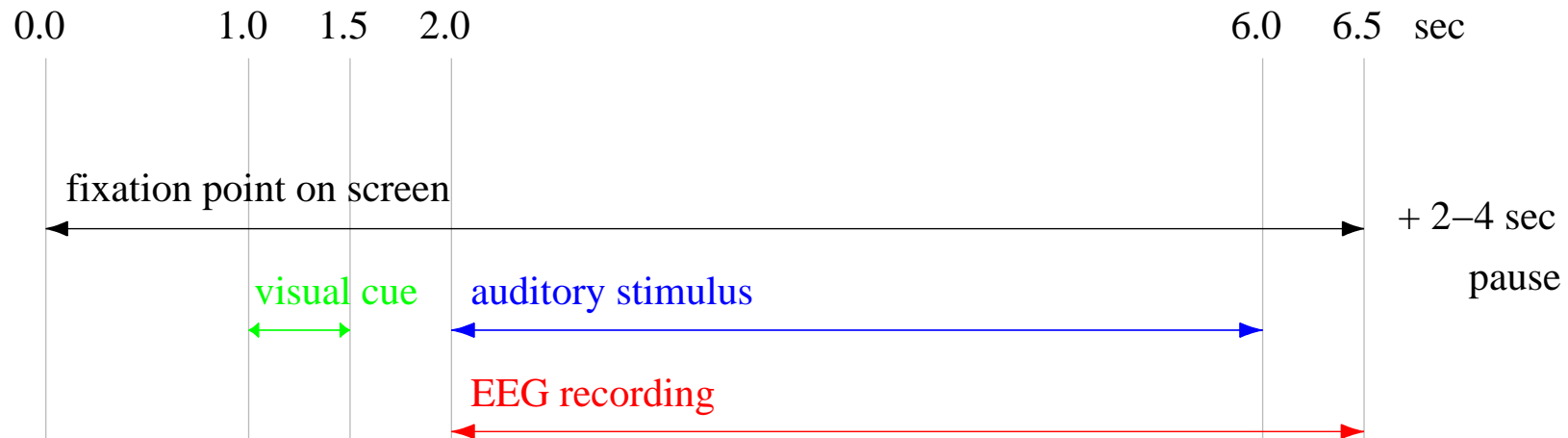
I: Auditory stimulation in EEG



I: Auditory stimulation in EEG



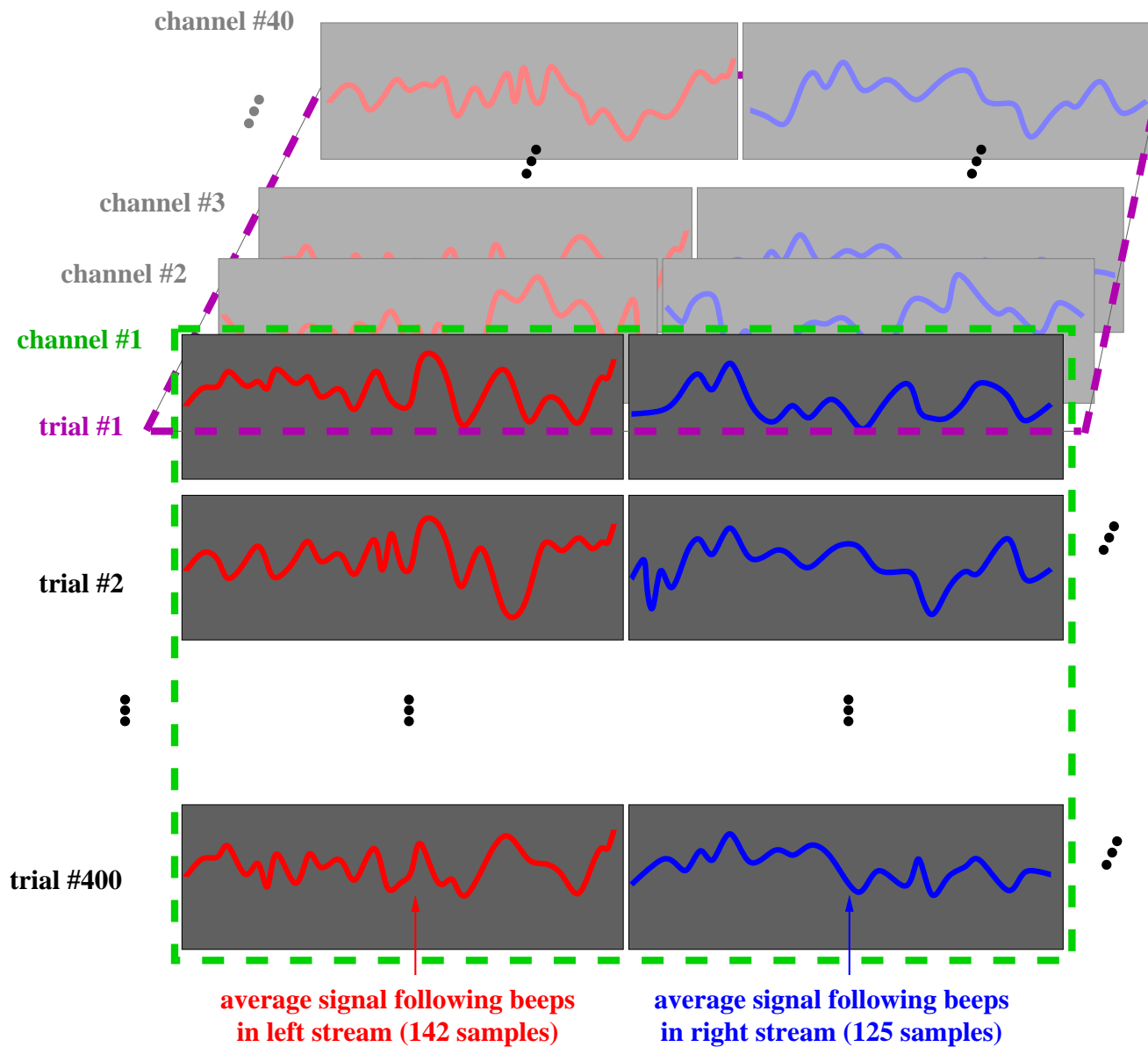
Trial structure



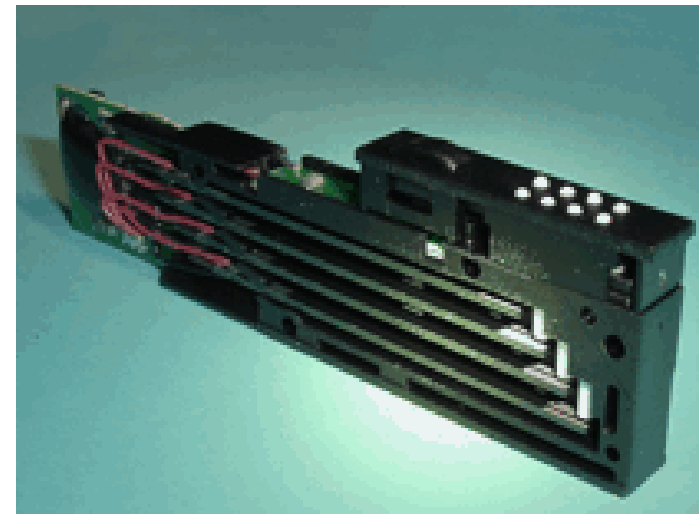
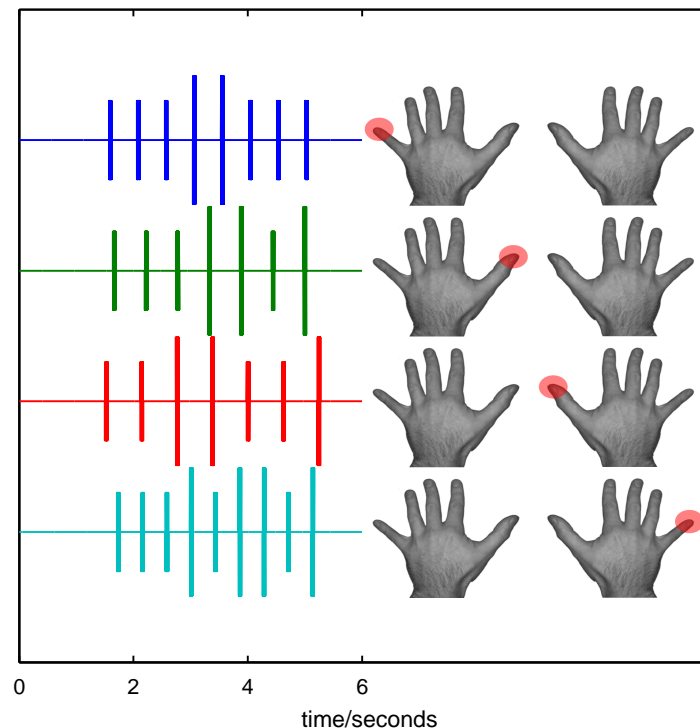
16 subjects, 400 trials each in one 3-hour session.

Cued direction of attention without feedback.

Data structure



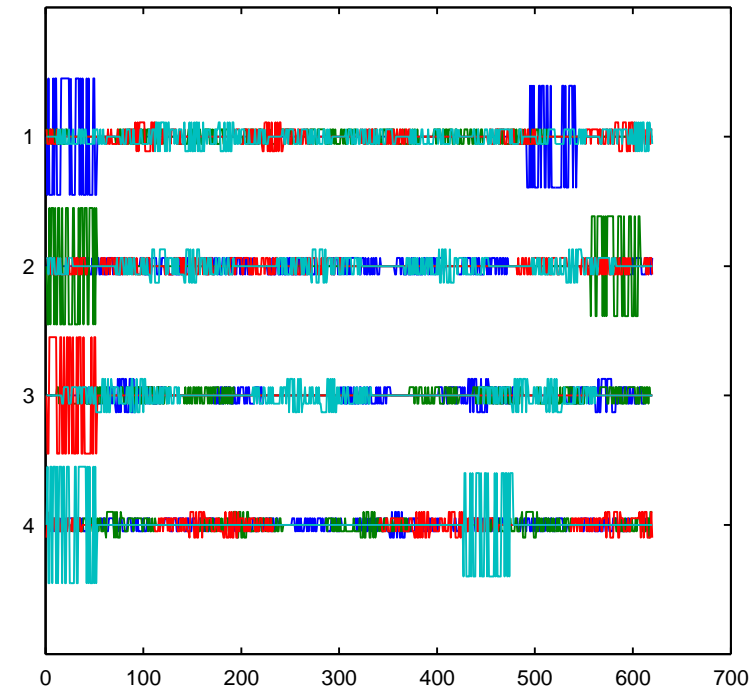
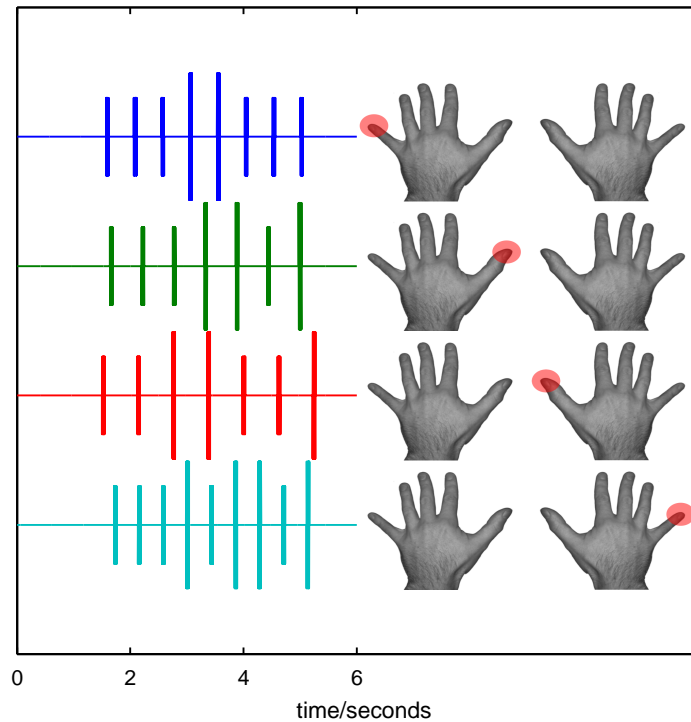
II: Tactile stimulation in MEG



9 subjects, each 200 cued trials without feedback.

5 classes including the no intentional control (NIC) state.

II: Tactile stimulation in MEG



9 subjects, each 200 cued trials without feedback.

5 classes including the no intentional control (NIC) state.



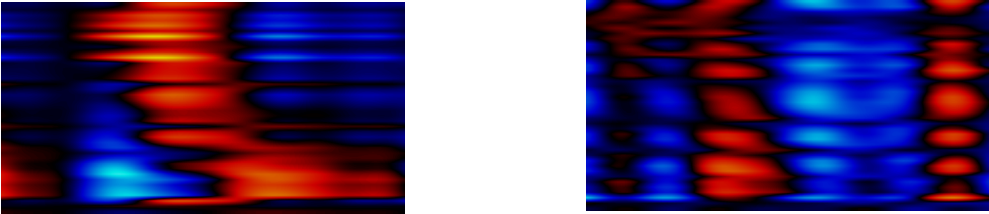
Classification



$$f(X) = \sum_{\text{sensors}} \sum_{\text{time}} \left(\underset{S \times T}{\text{G}} \cdot * \underset{S \times T}{X} \right) + b$$

G is the weight “vector” found by some classifier (SVM, LR, LDA...)

Classification

$$f(X^{(1)}) = \sum_{\text{sensors}} \sum_{\text{time}} \left(\underset{S \times T}{\text{G}} \cdot * \underset{S \times T}{X^{(1)}} \right) + b$$


Classification

$$f(X^{(2)}) = \sum_{\text{sensors}} \sum_{\text{time}} \left(\underset{S \times T}{\text{G}} \cdot * \underset{S \times T}{X^{(2)}} \right) + b$$



Classification



$$f(X^{(3)}) = \sum_{\text{sensors}} \sum_{\text{time}} \left(\underset{S \times T}{\text{G}} \cdot * \underset{S \times T}{X^{(3)}} \right) + b$$



Classification



$$f(X) = \sum_{\text{sensors}} \sum_{\text{time}} \left(\begin{array}{c} \text{Heatmap 1} \\ G \\ S \times T \end{array} \cdot * \begin{array}{c} \text{Heatmap 2} \\ X \\ S \times T \end{array} \right) + b$$



Classification



$$f(X) = \text{trace} \left(\begin{array}{c} \text{Heatmap of } G^T_{T \times S} \\ G^T_{T \times S} \end{array} \begin{array}{c} \text{Heatmap of } X_{S \times T} \\ X_{S \times T} \end{array} \right) + b$$



Classification



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$



Classification



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

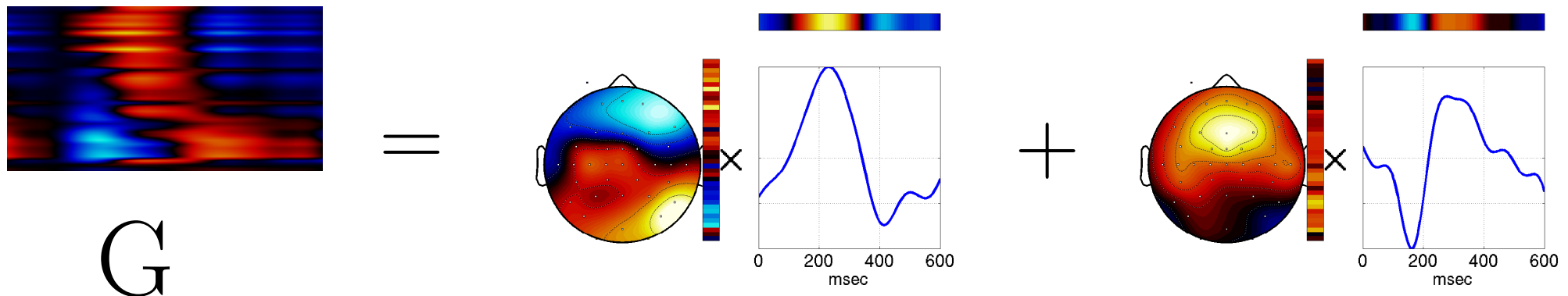
It can be helpful to assume that \mathbf{G} consists of only a small number (let's say 2) of relevant spatio-temporal features—

- whose spatial characteristics are stationary in time;
- whose temporal characteristics are stationary in space:

$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

It can be helpful to assume that \mathbf{G} consists of only a small number (let's say 2) of relevant spatio-temporal features—

- whose spatial characteristics are stationary in time;
- whose temporal characteristics are stationary in space:



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

It can be helpful to assume that \mathbf{G} consists of only a small number (let's say 2) of relevant spatio-temporal features—

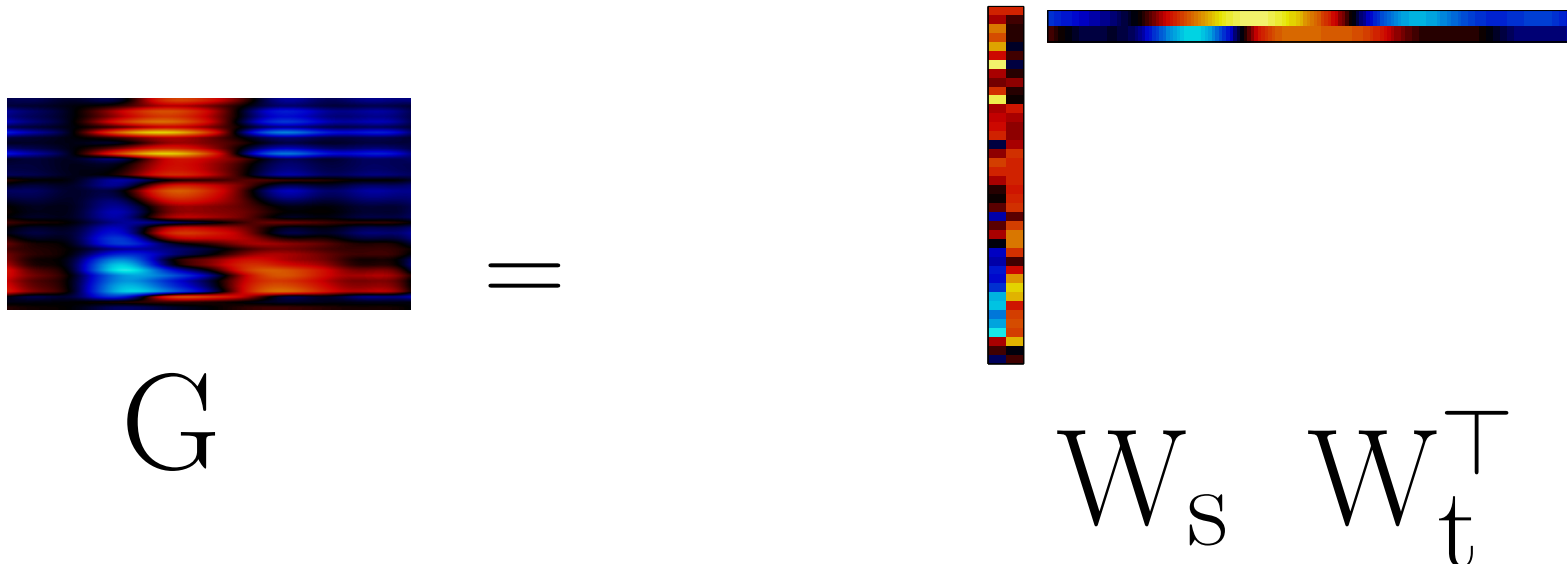
- whose spatial characteristics are stationary in time;
- whose temporal characteristics are stationary in space:

$$\mathbf{G} = \mathbf{w}_{s1} \mathbf{w}_{t1}^\top + \mathbf{w}_{s2} \mathbf{w}_{t2}^\top$$

$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

It can be helpful to assume that \mathbf{G} consists of only a small number (let's say 2) of relevant spatio-temporal features—

- whose spatial characteristics are stationary in time;
- whose temporal characteristics are stationary in space:



The diagram illustrates the decomposition of a spatio-temporal matrix \mathbf{G} into the product of two matrices, \mathbf{W}_s and \mathbf{W}_t^\top . On the left, a square heatmap labeled \mathbf{G} shows a complex pattern of red and blue horizontal bands. To its right is an equals sign. Further right, a vertical heatmap labeled \mathbf{W}_s shows a pattern of red and blue vertical bands. To its right is a horizontal heatmap labeled \mathbf{W}_t^\top showing a pattern of red and blue horizontal bands. The visual representation suggests that the spatial features in \mathbf{G} are captured by \mathbf{W}_s and the temporal features by \mathbf{W}_t^\top .

$$\mathbf{G} = \mathbf{W}_s \mathbf{W}_t^\top$$



Classification



$$\begin{aligned} f(\mathbf{X}) &= \text{tr} [\mathbf{G}^\top \mathbf{X}] \\ &= \text{tr} \left[(\mathbf{W}_s \mathbf{W}_t^\top)^\top \mathbf{X} \right] \end{aligned}$$



Classification



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}]$$



Bilinear Discriminant Analysis



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ 2 \times S & & T \times 2 \end{bmatrix}$$

(because arguments inside a trace operator may be cyclically reordered)



Bilinear Discriminant Analysis



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ 2 \times S & & T \times 2 \end{bmatrix} = \text{tr} \begin{bmatrix} \mathbf{L} \\ 2 \times 2 \end{bmatrix}$$



Bilinear Discriminant Analysis



$$f(X) = \text{tr} [G^T X]$$

$$= \text{tr} [W_t W_s^T X] = \text{tr} \begin{bmatrix} W_s^T & X & W_t \\ 2 \times S & & T \times 2 \end{bmatrix} = \text{tr} \begin{bmatrix} L \\ 2 \times 2 \end{bmatrix}$$

By making G low- instead of full-rank, we are assuming that there is a low-dimensional subspace onto which we can project X without loss (and perhaps with improvement) in performance of $f(X)$. We assert that, for classification, each data point $X^{(i)}$ can be sufficiently represented by a small number of coefficients—two, in our example: $(L_{11}^{(i)}, L_{22}^{(i)})$.



Bilinear Discriminant Analysis



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ 2 \times S & & T \times 2 \end{bmatrix} = \text{tr} \begin{bmatrix} \mathbf{L} \\ 2 \times 2 \end{bmatrix}$$

This implies a basis of spatial and temporal features \mathbf{A}_s and \mathbf{A}_t such that

$$f(\mathbf{X}) = f(\tilde{\mathbf{X}}), \quad \text{where } \tilde{\mathbf{X}} = \begin{array}{ccc} & \ddots & \\ \mathbf{A}_s & \mathbf{L} & \mathbf{A}_t^\top \\ S \times 2 & & 2 \times T \end{array}.$$

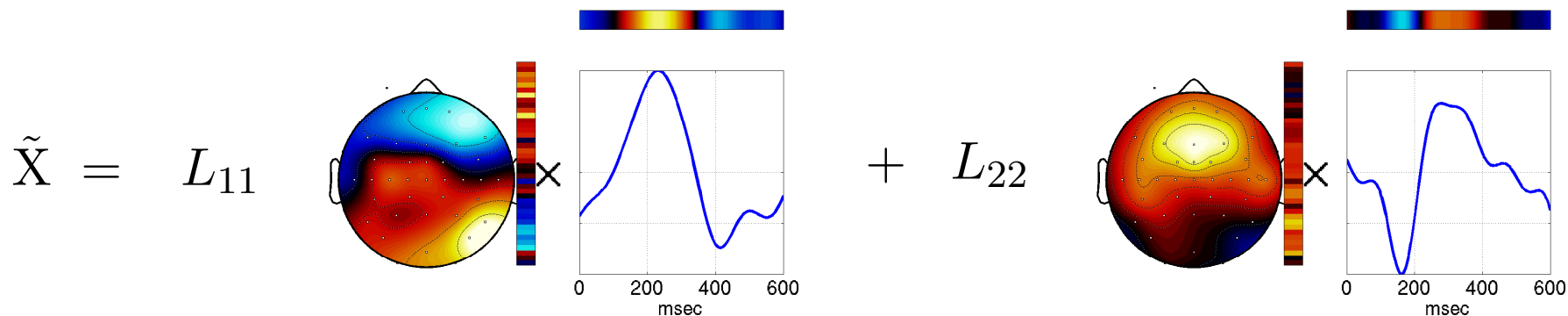
Bilinear Discriminant Analysis

$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ 2 \times S & & T \times 2 \end{bmatrix} = \text{tr} \begin{bmatrix} \mathbf{L} \\ 2 \times 2 \end{bmatrix}$$

This implies a basis of spatial and temporal features \mathbf{A}_s and \mathbf{A}_t such that

$$f(\mathbf{X}) = f(\tilde{\mathbf{X}}), \quad \text{where } \tilde{\mathbf{X}} = \begin{matrix} & \ddots \\ \mathbf{A}_s & \mathbf{L} & \mathbf{A}_t^\top \\ S \times 2 & 2 \times T \end{matrix}.$$





Bilinear Discriminant Analysis



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ 2 \times S & & T \times 2 \end{bmatrix} = \text{tr} \begin{bmatrix} \mathbf{L} \\ 2 \times 2 \end{bmatrix}$$

This implies a basis of spatial and temporal features \mathbf{A}_s and \mathbf{A}_t such that

$$f(\mathbf{X}) = f(\tilde{\mathbf{X}}), \quad \text{where } \tilde{\mathbf{X}} = \begin{matrix} & \ddots & \\ \mathbf{A}_s & \mathbf{L} & \mathbf{A}_t^\top \\ S \times 2 & 2 \times T & \end{matrix}.$$

$$\text{tr} [\mathbf{L}] = \text{tr} \begin{bmatrix} & \ddots & \\ \mathbf{W}_s^\top \mathbf{A}_s & \mathbf{L} & \mathbf{A}_t^\top \mathbf{W}_t \\ & & \end{bmatrix}$$

$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ 2 \times S & & T \times 2 \end{bmatrix} = \text{tr} \begin{bmatrix} \mathbf{L} \\ 2 \times 2 \end{bmatrix}$$

This implies a basis of spatial and temporal features \mathbf{A}_s and \mathbf{A}_t such that

$$f(\mathbf{X}) = f(\tilde{\mathbf{X}}), \quad \text{where } \tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{A}_s & \mathbf{L} & \mathbf{A}_t^\top \\ S \times 2 & 2 \times T & \end{bmatrix}.$$

$$\text{tr} [\mathbf{L}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top \mathbf{A}_s & \mathbf{L} & \mathbf{A}_t^\top \mathbf{W}_t \\ & & \end{bmatrix}, \quad \text{satisfied if } \mathbf{W}_s^\top \mathbf{A}_s = \mathbf{W}_t^\top \mathbf{A}_t = \mathbf{I}_{2 \times 2}$$



Bilinear Discriminant Analysis



$$f(\mathbf{X}) = \text{tr} [\mathbf{G}^\top \mathbf{X}]$$

$$= \text{tr} [\mathbf{W}_t \mathbf{W}_s^\top \mathbf{X}] = \text{tr} \begin{bmatrix} \mathbf{W}_s^\top & \mathbf{X} & \mathbf{W}_t \\ F \times S & T \times F & \end{bmatrix} = \text{tr} \begin{bmatrix} \mathbf{L} \\ F \times F \end{bmatrix}$$

$$f(\mathbf{X}) = f(\tilde{\mathbf{X}}) \quad \text{where} \quad \tilde{\mathbf{X}} = \mathbf{A}_s \mathbf{L} \mathbf{A}_t^\top \quad \text{such that} \quad \mathbf{W}_s^\top \mathbf{A}_s = \mathbf{W}_t^\top \mathbf{A}_t = \mathbf{I}_{F \times F}$$



Bilinear Discriminant Analysis



$$f(X) = \text{tr} [G^T X]$$

$$= \text{tr} [W_t W_s^T X] = \text{tr} \begin{bmatrix} W_s^T & X & W_t \\ F \times S & T \times F & \end{bmatrix} = \text{tr} \begin{bmatrix} L \\ F \times F \end{bmatrix}$$

$$f(X) = f(\tilde{X}) \quad \text{where} \quad \tilde{X} = A_s L A_t^T \quad \text{such that} \quad W_s^T A_s = W_t^T A_t = \underset{F \times F}{I}$$

How do we obtain W_s, W_t ? Simple idea: feed X into a classifier and look at the **singular value decomposition** of the resulting G :

$$G = R_s S^2 R_t^T, \quad \text{where} \quad R_s^T R_s = R_t^T R_t = I$$



Bilinear Discriminant Analysis



$$f(X) = \text{tr} [G^T X]$$

$$= \text{tr} [W_t W_s^T X] = \text{tr} \begin{bmatrix} W_s^T & X & W_t \\ F \times S & T \times F & \end{bmatrix} = \text{tr} \begin{bmatrix} L \\ F \times F \end{bmatrix}$$

$$f(X) = f(\tilde{X}) \quad \text{where} \quad \tilde{X} = A_s L A_t^T \quad \text{such that} \quad W_s^T A_s = W_t^T A_t = \underset{F \times F}{I}$$

How do we obtain W_s, W_t ? Simple idea: feed X into a classifier and look at the **singular value decomposition** of the resulting G :

$$G = R_s S^2 R_t^T, \quad \text{where} \quad R_s^T R_s = R_t^T R_t = I$$

Thus

$$W_s = R_s S, \quad W_t = R_t S,$$

$$A_s = R_s S^{-1}, \quad A_t = R_t S^{-1}.$$



Bilinear Discriminant Analysis



$$f(X) = \text{tr} [G^T X]$$

$$= \text{tr} [W_t W_s^T X] = \text{tr} \begin{bmatrix} W_s^T & X & W_t \\ F \times S & T \times F & \end{bmatrix} = \text{tr} \begin{bmatrix} L \\ F \times F \end{bmatrix}$$

$$f(X) = f(\tilde{X}) \quad \text{where} \quad \tilde{X} = A_s L A_t^T \quad \text{such that} \quad W_s^T A_s = W_t^T A_t = \underset{F \times F}{I}$$

How do we obtain W_s, W_t ? Simple idea: feed X into a classifier and look at the **singular value decomposition** of the resulting G :

$$G = R_s S^2 R_t^T, \quad \text{where} \quad R_s^T R_s = R_t^T R_t = I$$

If you want to reduce the rank, you can throw away columns of R_s and R_t and recompute G .

$$f(X) = \text{tr} [G^T X]$$

$$= \text{tr} [W_t W_s^T X] = \text{tr} \begin{bmatrix} W_s^T & X & W_t \\ F \times S & T \times F & \end{bmatrix} = \text{tr} \begin{bmatrix} L \\ F \times F \end{bmatrix}$$

$$f(X) = f(\tilde{X}) \quad \text{where} \quad \tilde{X} = A_s L A_t^T \quad \text{such that} \quad W_s^T A_s = W_t^T A_t = \underset{F \times F}{I}$$

How do we obtain W_s, W_t ? Simple idea: feed X into a classifier and look at the **singular value decomposition** of the resulting G :

$$G = R_s S^2 R_t^T, \quad \text{where} \quad R_s^T R_s = R_t^T R_t = I$$

Alternatively: obtain a low-rank G in the first place, by using a classifier which is regularized by penalizing the rank of G :

Tomiooka and Aihara (ICML 2007) give a (convex!) formulation for LR.



Bilinear Discriminant Analysis



$$f(X) = \text{tr} [G^T X]$$

$$= \text{tr} [W_t W_s^T X] = \text{tr} \begin{bmatrix} W_s^T & X & W_t \\ F \times S & T \times F & \end{bmatrix} = \text{tr} \begin{bmatrix} L \\ F \times F \end{bmatrix}$$

$$f(X) = f(\tilde{X}) \quad \text{where} \quad \tilde{X} = A_s L A_t^T \quad \text{such that} \quad W_s^T A_s = W_t^T A_t = \underset{F \times F}{I}$$

How do we obtain W_s, W_t ? Simple idea: feed X into a classifier and look at the **singular value decomposition** of the resulting G :

$$G = R_s S^2 R_t^T, \quad \text{where} \quad R_s^T R_s = R_t^T R_t = I$$

...but neither approach performs at its best if G is found directly in the space of X .



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with P such that $E\{P^T X (P^T X)^T\} = I$



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with P such that $E\{P^T X X^T P\} = I$



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with P such that $P^T E\{XX^T\}P = I$



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with P such that $P^T \Sigma P = I$



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with P such that $\Sigma = P^{-T} P^{-1}$



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with (e.g.) $P = \Sigma^{-\frac{1}{2}}$



Spatial Whitening in ERP classification



Spatial whitening (decorrelating) transformations are very common in EEG analysis, as a partial way of undoing the volume conduction effect which causes all sensor outputs to be highly correlated. Examples:

- surface Laplacian filters (approximate)
- ICA
- CSP and friends (esp. for extraction of bandpower features)

Generally: $X_P = P^T X$, with (e.g.) $P = \Sigma^{-\frac{1}{2}}$

Even in linear classification, non-orthonormal transformations of the data affect the classifier's **regularization environment** and can lead to very different results.

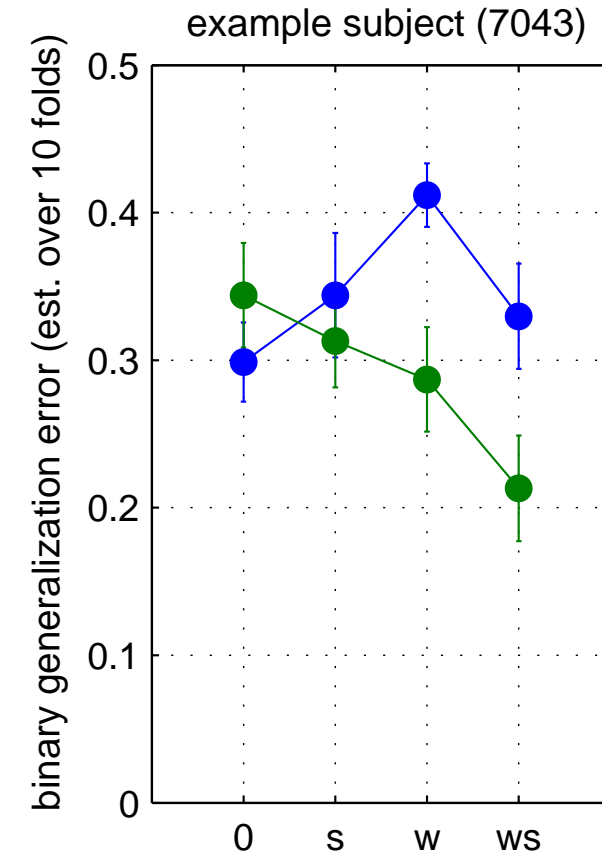
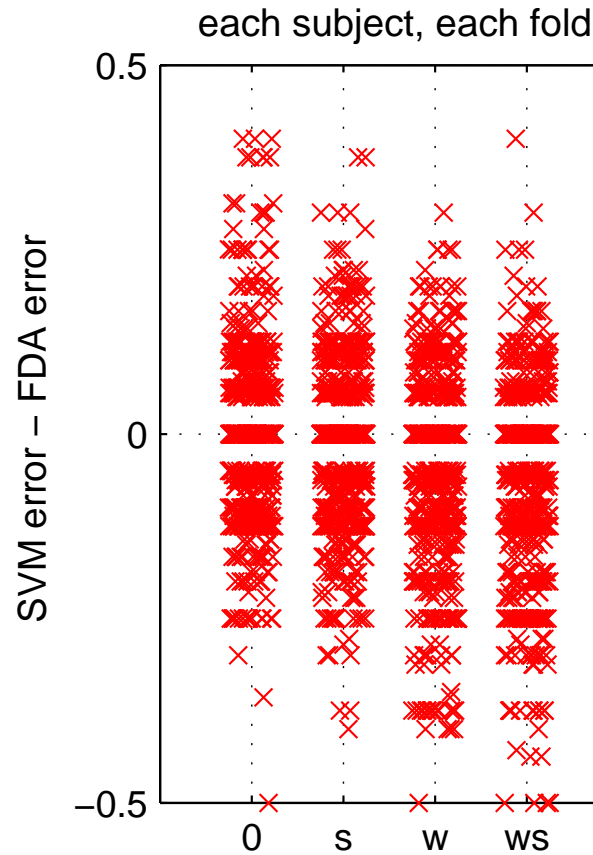
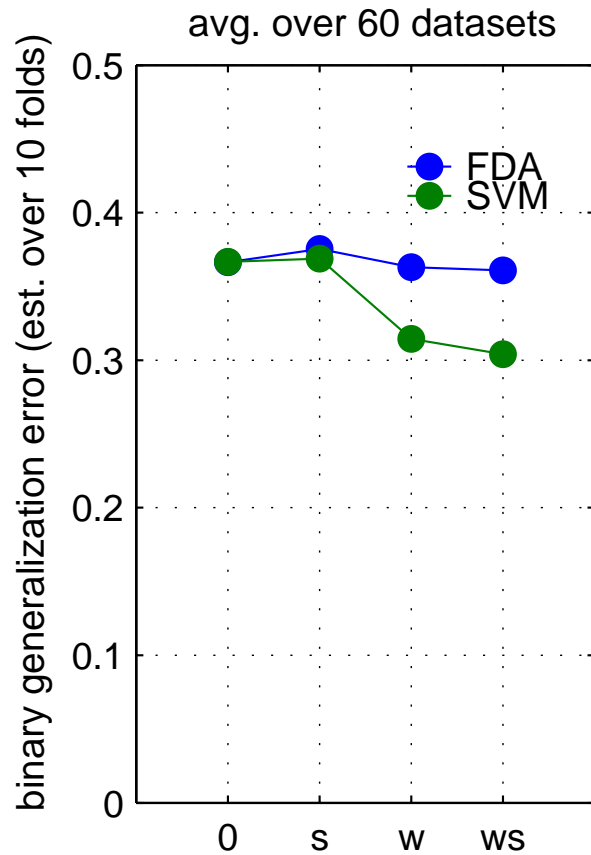


Spatial Whitening in ERP classification



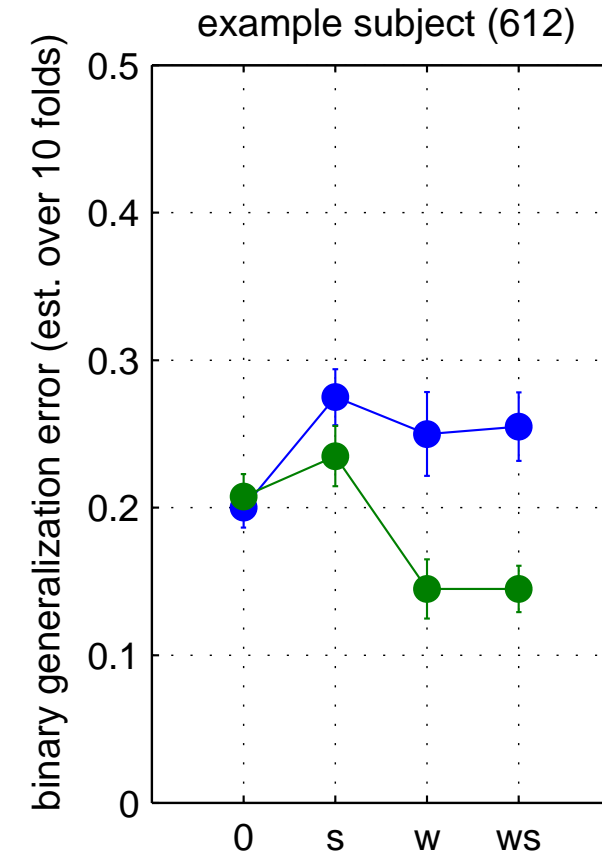
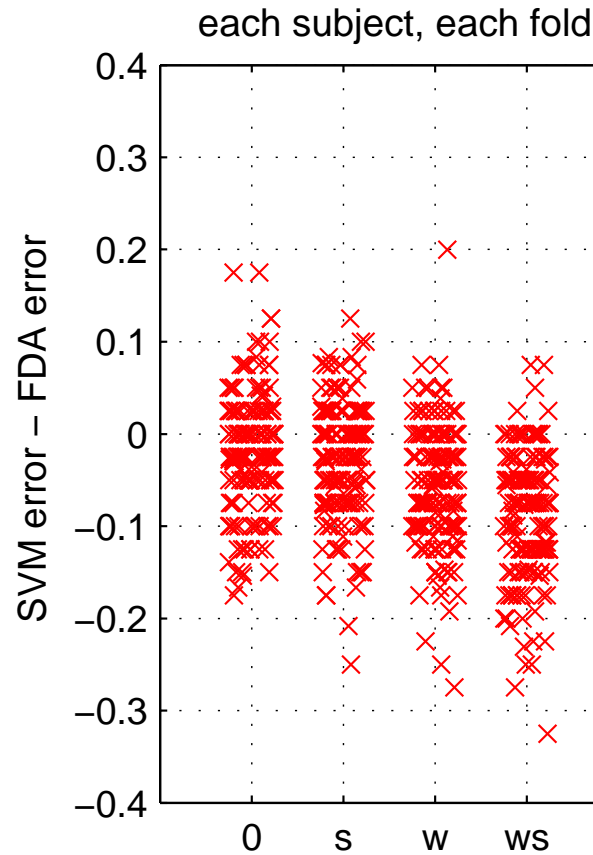
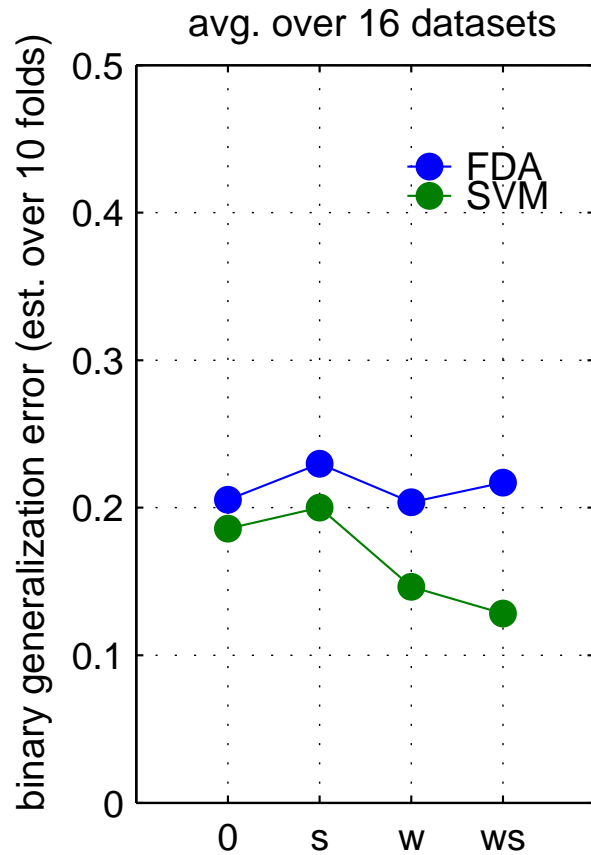
Krusienski et al. (J. Neural Eng. 2006) report in visual ERP classification (in a grid-speller) that Fisher's (unregularized) Linear Discriminant Analysis performs as well as, or better than, Support Vector Machines.

We suggest that this is because no decorrelation was performed. The lack of whitening masked the potential benefits of regularization:



preprocessing (w = whiten, s = center & standardize each trial-by-channel)

Felix Biessmann's visual speller data (10 subjects x 6 stimulus conditions), offline analysis

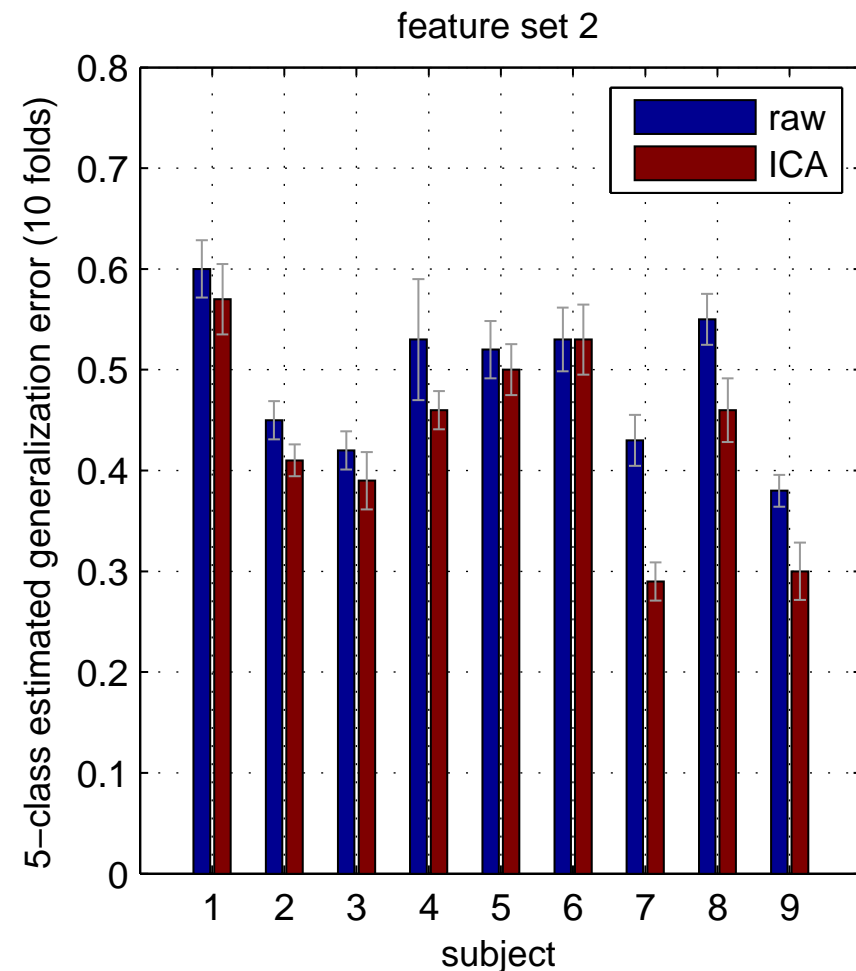
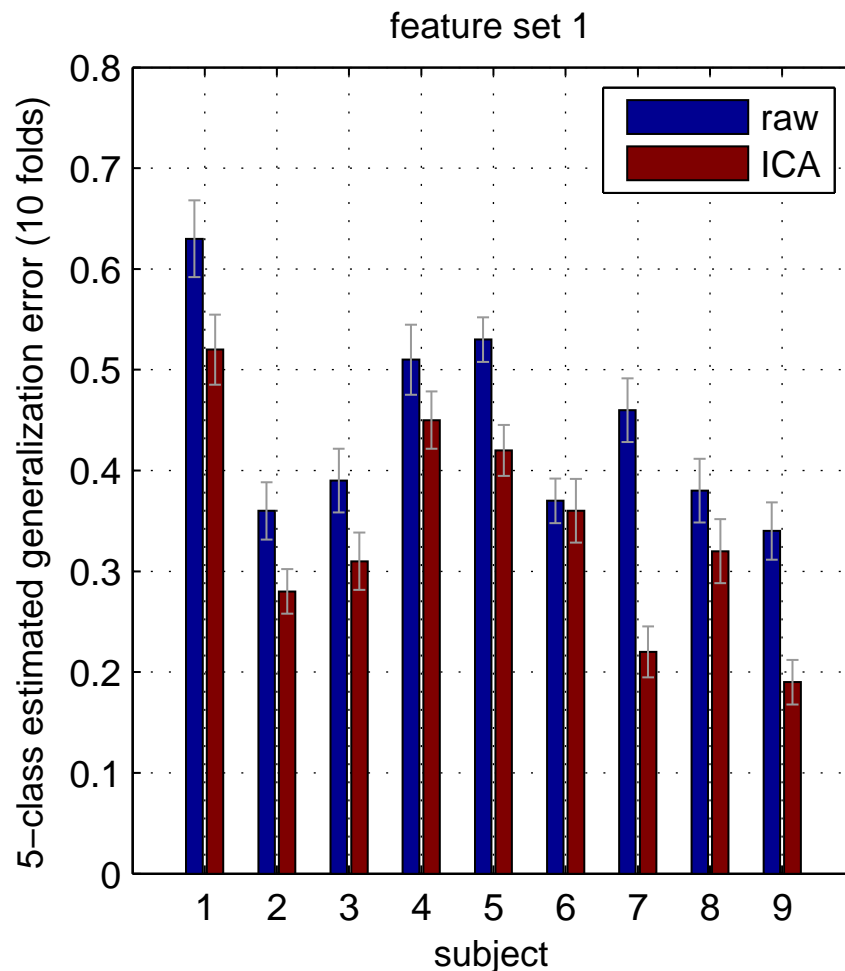


preprocessing (w = whiten, s = center & standardize each trial-by-channel)

auditory ERP data, offline analysis

SVM results on tactile MEG data

(Cornelius Raths' Diploma Thesis 2007, paper in preparation):





Filters and Patterns



Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] &= \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] &= \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$



Filters and Patterns



Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] &= \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] &= \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$

...and decompose G_P instead of G :

$$G_P = R_s \overset{\cdot \cdot \cdot}{S^2} R_t^\top$$

Filters and Patterns

Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] &= \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] &= \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$

...and decompose G_P instead of G :

$$G_P = R_S S^2 R_t^\top$$
$$f(X) = \text{tr} [G_P^\top P_S^\top X P_t] = \text{tr} \begin{bmatrix} \ddots & & \\ S R_S^\top P_S^\top X P_t R_t & \ddots & \\ & & \ddots \end{bmatrix}$$

Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] = \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] = \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$

...and decompose G_P instead of G :

$$G_P = R_S S^2 R_t^\top$$

$$f(X) = \text{tr} [G_P^\top P_S^\top X P_t] = \text{tr} \begin{bmatrix} \ddots & & \\ S R_S^\top P_S^\top X P_t R_t & \ddots & \\ \ddots & & \end{bmatrix}$$

So

$$W_S = P_S R_S S, \quad W_t = P_t R_t S,$$

$$A_S = P_S^{-\top} R_S S^{-1}, \quad A_t = P_t^{-\top} R_t S^{-1}.$$



Preconditioning and Regularization



Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] &= \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] &= \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$

Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] &= \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] &= \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$

the classifier (LR, SVM,...) optimizes G_P instead of G . This does not change $f(X)$, so the loss term is unaffected. However, assuming an L2 regularizer, the classifier's objective becomes:

$$\lambda \text{tr} [G_P^\top G_P] + \sum_i \mathcal{L} \left\{ y^{(i)} f(X^{(i)}) \right\},$$

Learn the classifier weights on the *preconditioned* data $X_P = P_S^\top X P_t$:

$$\begin{aligned} f(X) &= \text{tr} [G_P^\top X_P] &= \text{tr} [G^\top X] \\ &= \text{tr} [G_P^\top P_S^\top X P_t] &= \text{tr} [P_t G_P^\top P_S^\top X] \end{aligned}$$

the classifier (LR, SVM,...) optimizes G_P instead of G . This does not change $f(X)$, so the loss term is unaffected. However, assuming an L2 regularizer, the classifier's objective becomes:

$$\lambda \text{tr} [G_P^\top G_P] + \sum_i \mathcal{L} \left\{ y^{(i)} f(X^{(i)}) \right\},$$

in which further expansion, substitution and reordering gives

$$\text{tr} [G_P^\top G_P] = \text{tr} \begin{bmatrix} S^4 A_t^\top P_t P_t^\top A_t & S^4 A_s^\top P_s P_s^\top A_s \end{bmatrix}.$$



Preconditioning and Prior knowledge



The regularization term can be written as

$$\text{tr} \begin{bmatrix} S^4 A_t^\top \Sigma_t^{-1} A_t & S^4 A_s^\top \Sigma_s^{-1} A_s \end{bmatrix} \quad \text{where} \quad \Sigma_t^{-1} = P_t P_t^\top \quad \text{and} \quad \Sigma_s^{-1} = P_s P_s^\top.$$

It contains terms which look a little like logged Gaussian prior probabilities over the spatial and temporal patterns, (not incorporated in the usual way, though: cross-terms for $F > 1$, and multiplication of spatial and temporal terms).



Preconditioning and Prior knowledge



The regularization term can be written as

$$\text{tr} \begin{bmatrix} S^4 A_t^\top \Sigma_t^{-1} A_t & S^4 A_s^\top \Sigma_s^{-1} A_s \end{bmatrix} \quad \text{where} \quad \Sigma_t^{-1} = P_t P_t^\top \quad \text{and} \quad \Sigma_s^{-1} = P_s P_s^\top.$$

It contains terms which look a little like logged Gaussian prior probabilities over the spatial and temporal patterns, (not incorporated in the usual way, though: cross-terms for $F > 1$, and multiplication of spatial and temporal terms).

Some alternative approaches use direct prior probabilities on the spatial patterns or filters, for example to smooth them (e.g. Dyrholm et al. JMLR 2007, Dyrholm and Parra, Proc. IEEE EMBS 2006, Farquhar et al. Applied Neuroscience Conference 2007).

Similar, though not identical, effects are achieved by suitable choice of Σ_s and Σ_t in the preconditioning approach above, with the advantage that convex formulations can be applied for finding G_P .

The regularization term can be written as

$$\text{tr} \begin{bmatrix} S^4 A_t^\top \Sigma_t^{-1} A_t & S^4 A_s^\top \Sigma_s^{-1} A_s \end{bmatrix} \quad \text{where} \quad \Sigma_t^{-1} = P_t P_t^\top \quad \text{and} \quad \Sigma_s^{-1} = P_s P_s^\top.$$

It contains terms which look a little like logged Gaussian prior probabilities over the spatial and temporal patterns, (not incorporated in the usual way, though: cross-terms for $F > 1$, and multiplication of spatial and temporal terms).

Some alternative approaches use direct prior probabilities on the spatial patterns or filters, for example to smooth them (e.g. Dyrholm et al. JMLR 2007, Dyrholm and Parra, Proc. IEEE EMBS 2006, Farquhar et al. Applied Neuroscience Conference 2007).

Similar, though not identical, effects are achieved by suitable choice of Σ_s and Σ_t in the preconditioning approach above, with the advantage that convex formulations can be applied for finding G_P .

Choosing Σ_s equal to the EEG or MEG sensor covariance matrix seems sensible, since it means the spatial basis functions by which we represent \tilde{X} reflect realistic EEG- or MEG-like volume-conduction properties.

A Σ_t may also be chosen as a prior which smooths the temporal patterns, although to save computation time one can also simply smooth and downsample the data.



The regularization term can be written as

$$\text{tr} \begin{bmatrix} S^4 A_t^\top \Sigma_t^{-1} A_t & S^4 A_s^\top \Sigma_s^{-1} A_s \end{bmatrix} \quad \text{where} \quad \Sigma_t^{-1} = P_t P_t^\top \quad \text{and} \quad \Sigma_s^{-1} = P_s R_s R_s^\top P_s^\top.$$

It contains terms which look a little like logged Gaussian prior probabilities over the spatial and temporal patterns, (not incorporated in the usual way, though: cross-terms for $F > 1$, and multiplication of spatial and temporal terms).

Some alternative approaches use direct prior probabilities on the spatial patterns or filters, for example to smooth them (e.g. Dyrholm et al. JMLR 2007, Dyrholm and Parra, Proc. IEEE EMBS 2006, Farquhar et al. Applied Neuroscience Conference 2007).

Similar, though not identical, effects are achieved by suitable choice of Σ_s and Σ_t in the preconditioning approach above, with the advantage that convex formulations can be applied for finding G_P .

Choosing Σ_s equal to the EEG or MEG sensor covariance matrix seems sensible, since it means the spatial basis functions by which we represent \tilde{X} reflect realistic EEG- or MEG-like volume-conduction properties.

A Σ_t may also be chosen as a prior which smooths the temporal patterns, although to save computation time one can also simply smooth and downsample the data.

The regularization term can be written as

$$\text{tr} \begin{bmatrix} S^4 A_t^\top \Sigma_t^{-1} A_t & S^4 A_s^\top \Sigma_s^{-1} A_s \end{bmatrix} \quad \text{where} \quad \Sigma_t^{-1} = P_t P_t^\top \quad \text{and} \quad \Sigma_s^{-1} = P_s P_s^\top.$$

It contains terms which look a little like logged Gaussian prior probabilities over the spatial and temporal patterns, (not incorporated in the usual way, though: cross-terms for $F > 1$, and multiplication of spatial and temporal terms).

Some alternative approaches use direct prior probabilities on the spatial patterns or filters, for example to smooth them (e.g. Dyrholm et al. JMLR 2007, Dyrholm and Parra, Proc. IEEE EMBS 2006, Farquhar et al. Applied Neuroscience Conference 2007).

Similar, though not identical, effects are achieved by suitable choice of Σ_s and Σ_t in the preconditioning approach above, with the advantage that convex formulations can be applied for finding G_P .

Choosing Σ_s equal to the EEG or MEG sensor covariance matrix seems sensible, since it means the spatial basis functions by which we represent \tilde{X} reflect realistic EEG- or MEG-like volume-conduction properties.

A Σ_t may also be chosen as a prior which smooths the temporal patterns, although to save computation time one can also simply smooth and downsample the data.



Note that an arbitrary invertible $F \times F$ matrix B may be applied after the filters have been found:

$$f(X) = \text{tr} [B \ SR_s^\top P_s^\top \ X \ P_t R_t S \ B^{-1}] .$$



Note that an arbitrary invertible $F \times F$ matrix B may be applied after the filters have been found:

$$f(X) = \text{tr} [B \quad SR_S^\top P_S^\top \quad X \quad P_t R_t S \quad B^{-1}] .$$

This does not change $f(X)$, but it changes the interpretation of the discriminative components:

Now

$$W_S = P_S R_S \begin{matrix} \ddots \\ S \\ \end{matrix} B^\top , \quad W_t = P_t R_t \begin{matrix} \ddots \\ S \\ \end{matrix} B^{-1} ,$$

$$A_S = P_S^{-\top} R_S \begin{matrix} \ddots \\ S^{-1} \\ \end{matrix} B^{-1} , \quad A_t = P_t^{-\top} R_t \begin{matrix} \ddots \\ S^{-1} \\ \end{matrix} B^\top .$$



Note that an arbitrary invertible $F \times F$ matrix B may be applied after the filters have been found:

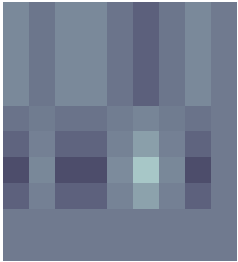
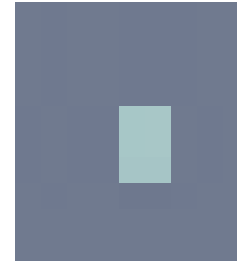
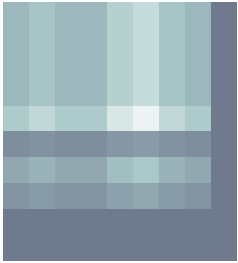
$$f(X) = \text{tr} [B \ SR_S^\top P_S^\top \ X \ P_t R_t S \ B^{-1}] .$$

Various approaches are available:

- set $B = I$ (leave filters and patterns to be determined by the interaction between regularizer and preconditioner)
- optimize B such that the (L_{11}, \dots, L_{FF}) are maximally independent across different instances $L^{(i)}$ (Dyrholm et al., JMLR 2007).
- optimize B according to L1 penalties on A_s and A_t (for sparse basis functions, perhaps confined in space and time) or on W_s and W_t (fewer electrodes to stick).



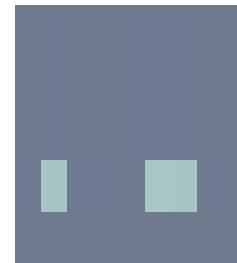
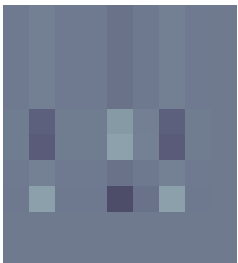
Sparsification



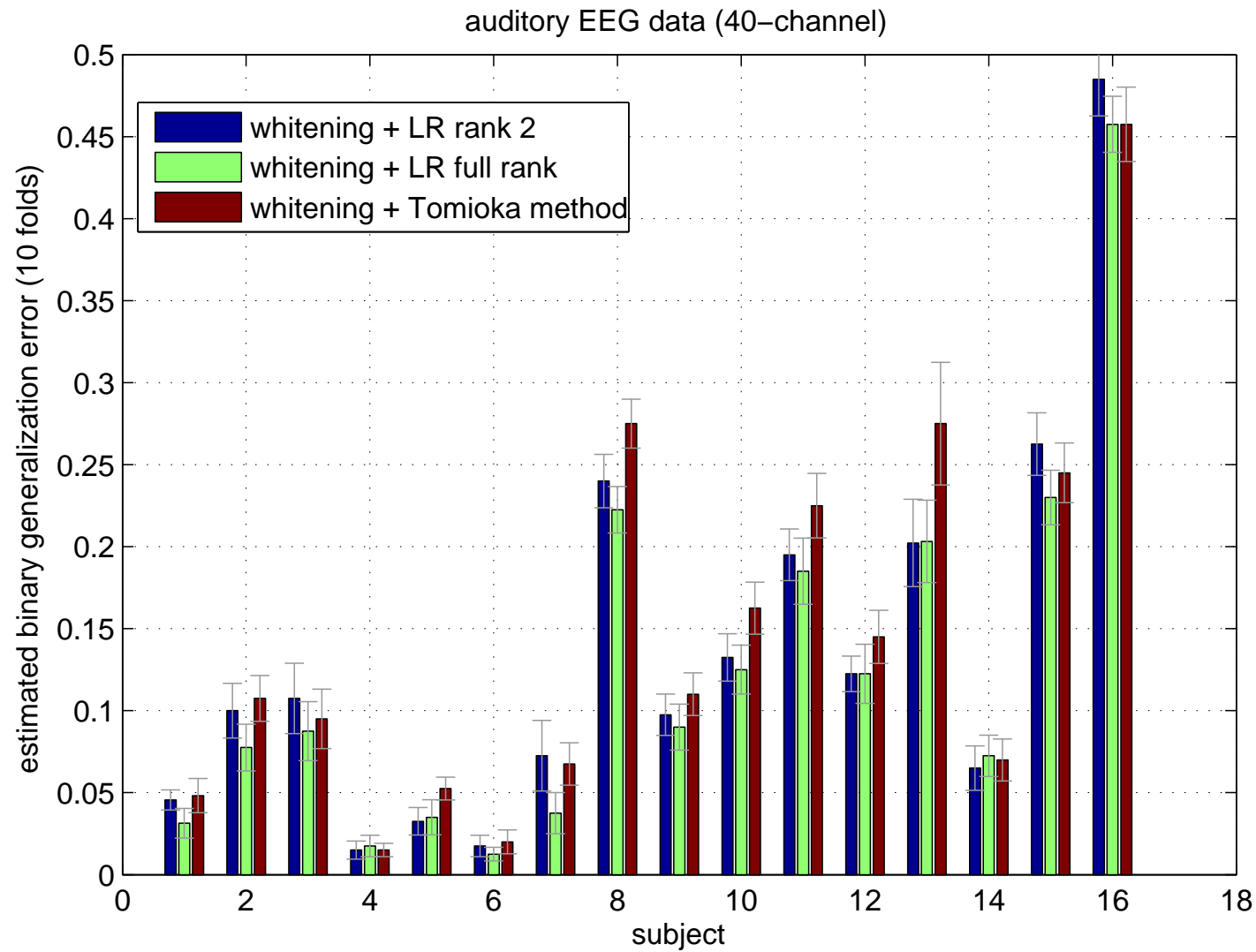
outer products
from SVD of
composite,
 $G = USV^T$



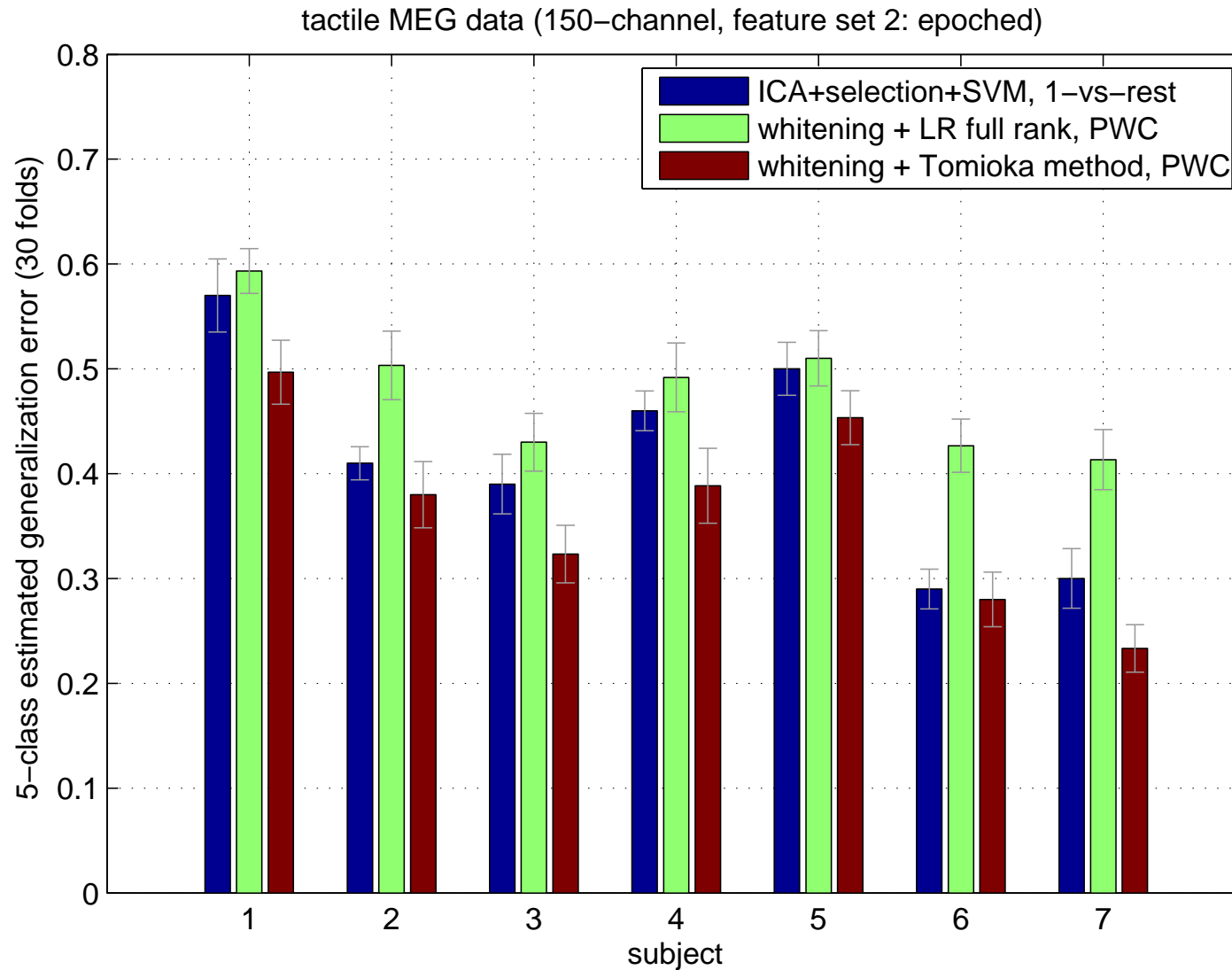
after
minimization of
L1 norms of UB^T
and VB^{-1} , w.r.t
 $F \times F$ matrix B



Does a low-rank constraint help?

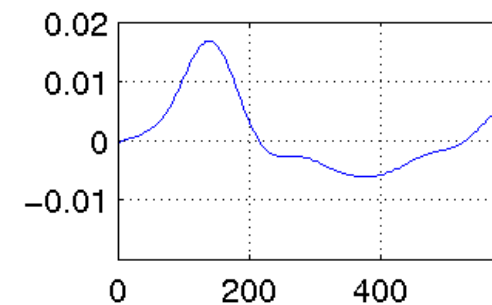
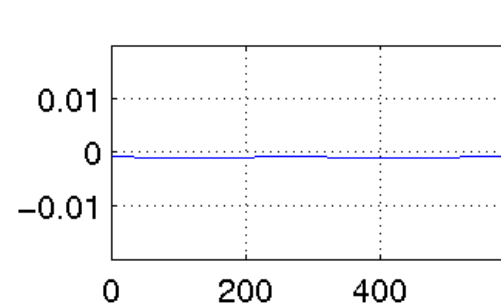
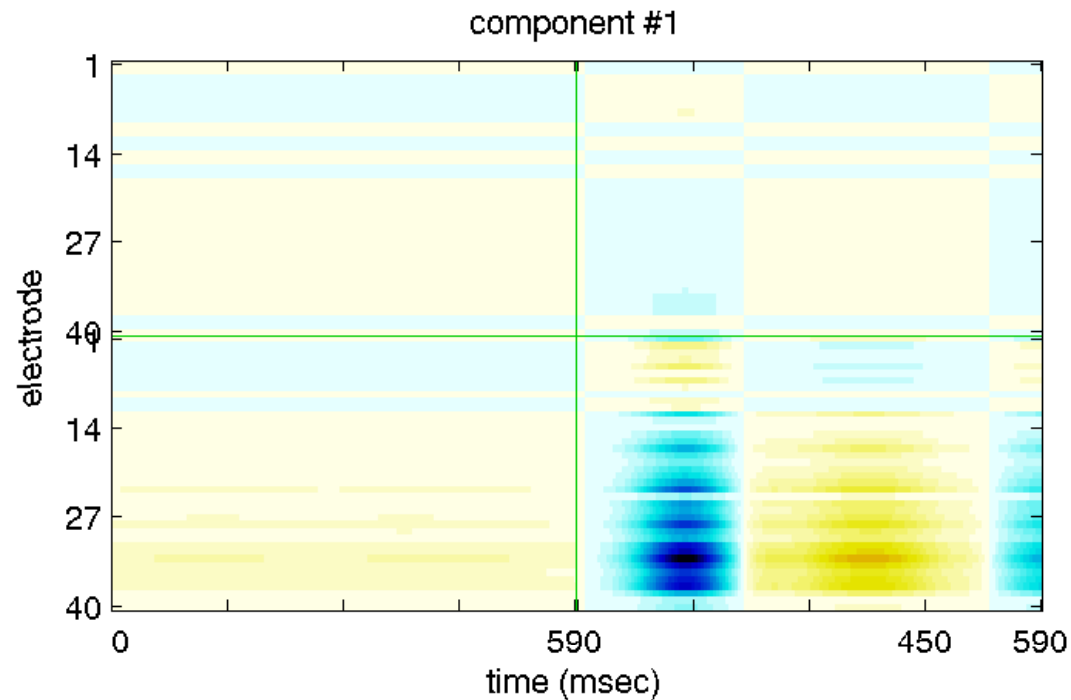
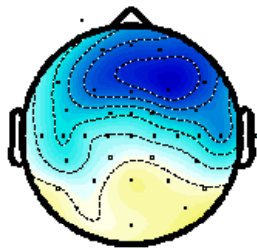


Does a low-rank constraint help?



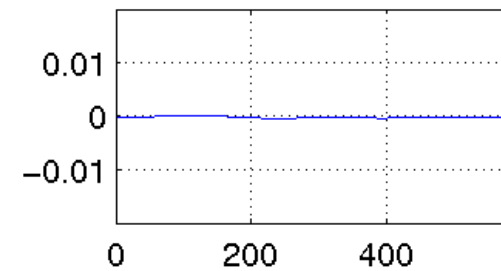
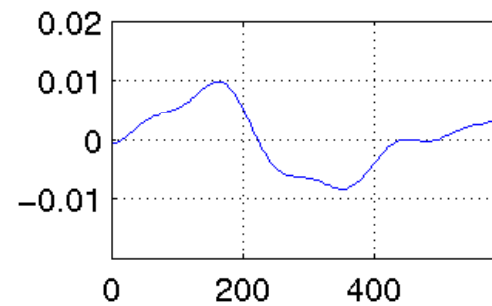
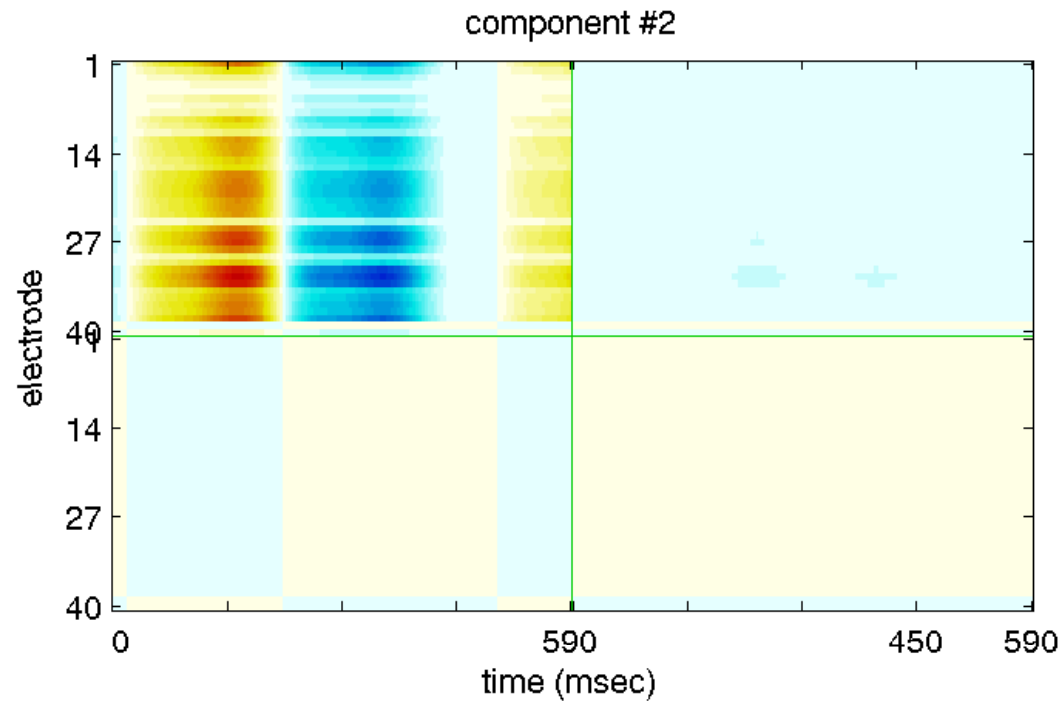
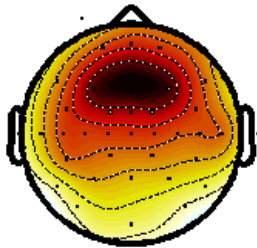
Example decomposition (auditory)

Auditory EEG data:



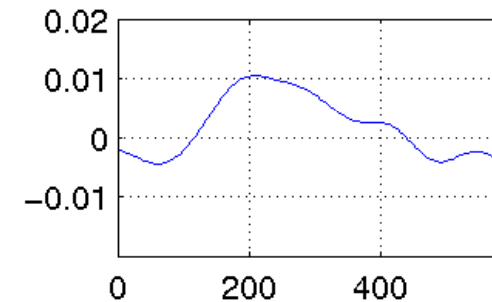
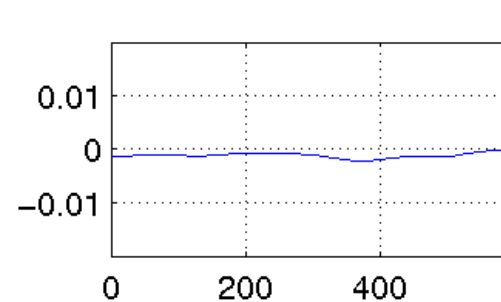
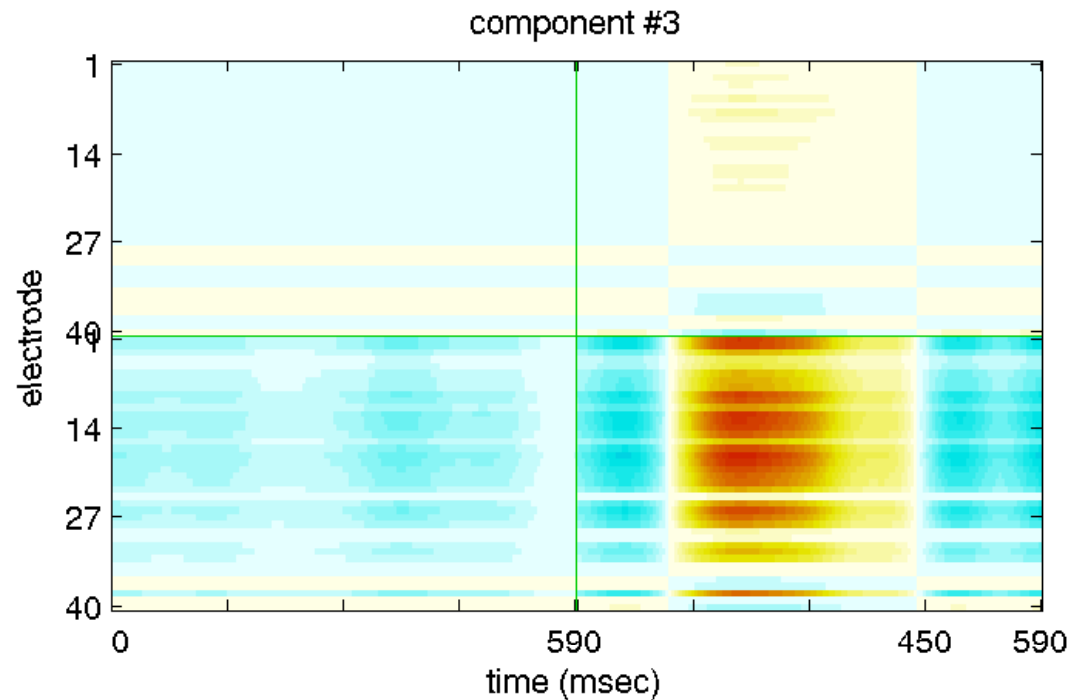
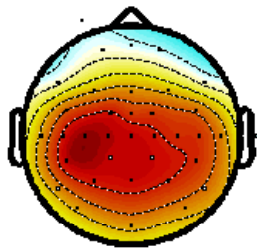
Example decomposition (auditory)

Auditory EEG data:



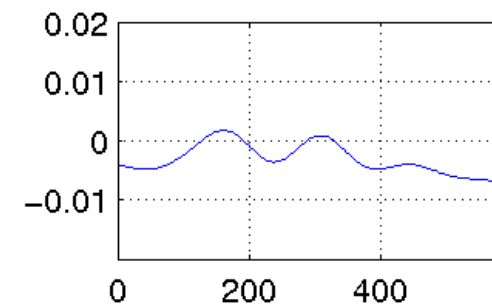
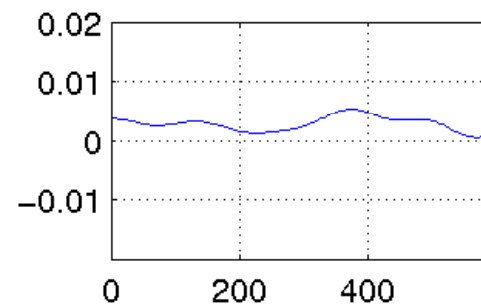
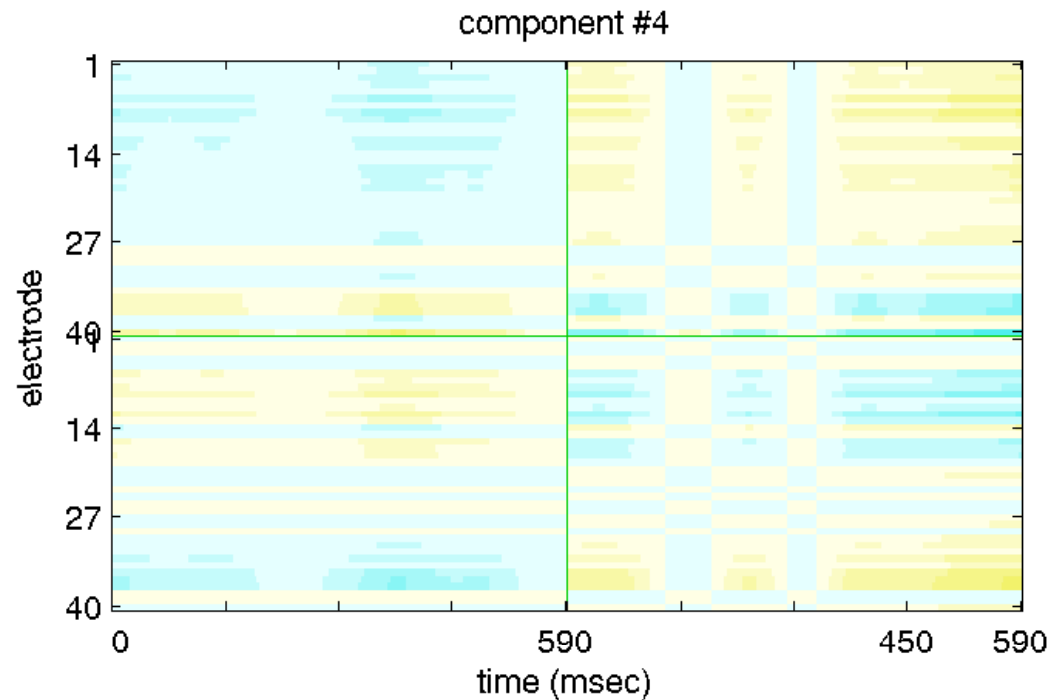
Example decomposition (auditory)

Auditory EEG data:



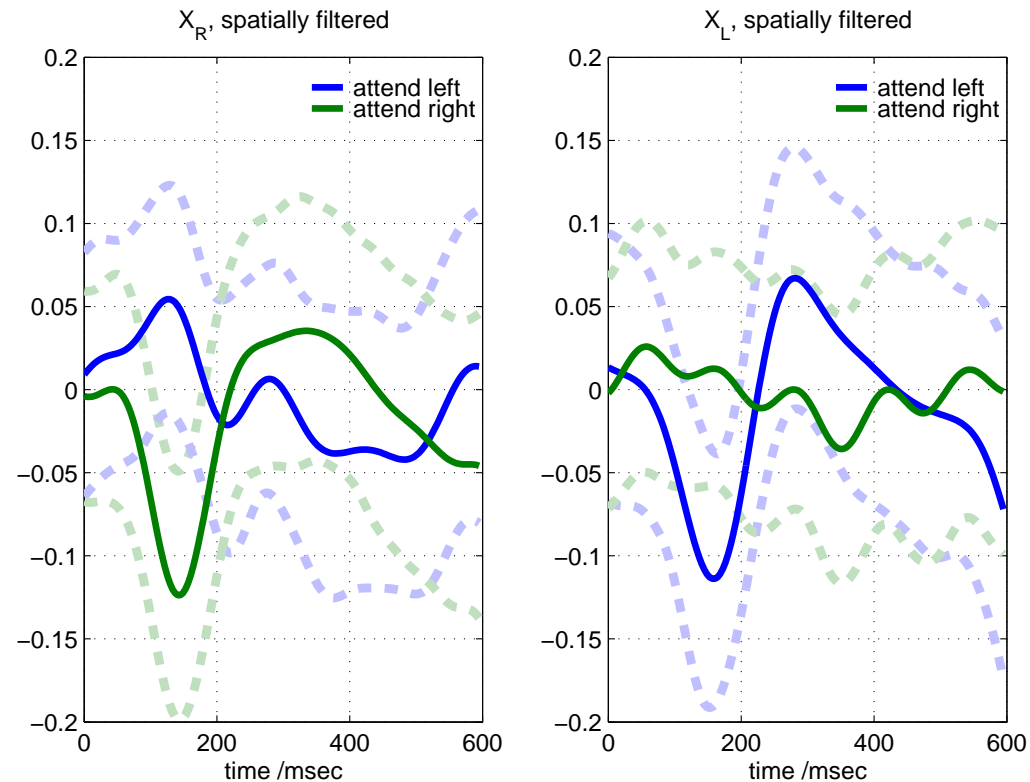
Example decomposition (auditory)

Auditory EEG data:

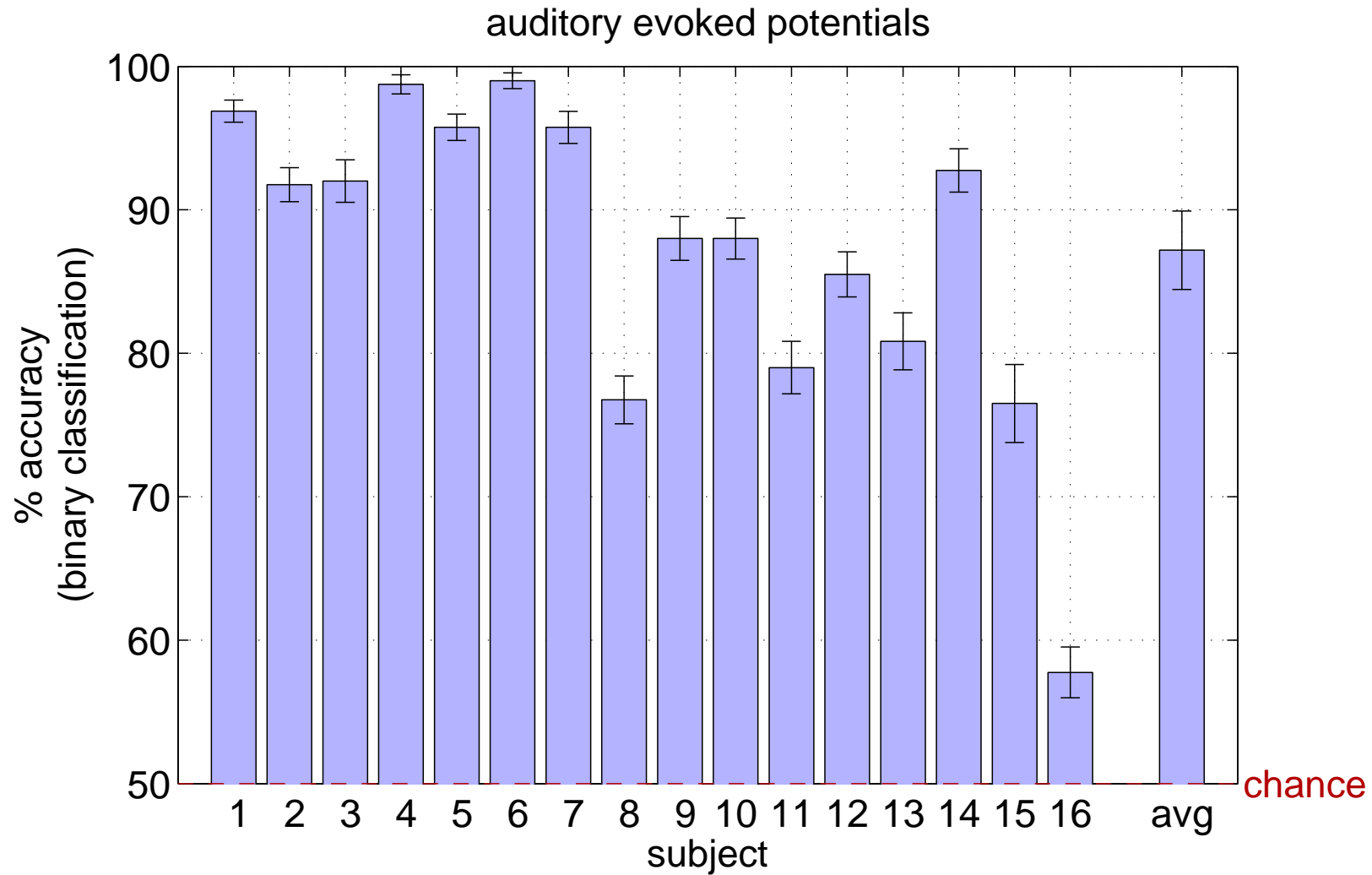


Example decomposition (auditory)

Auditory EEG data:

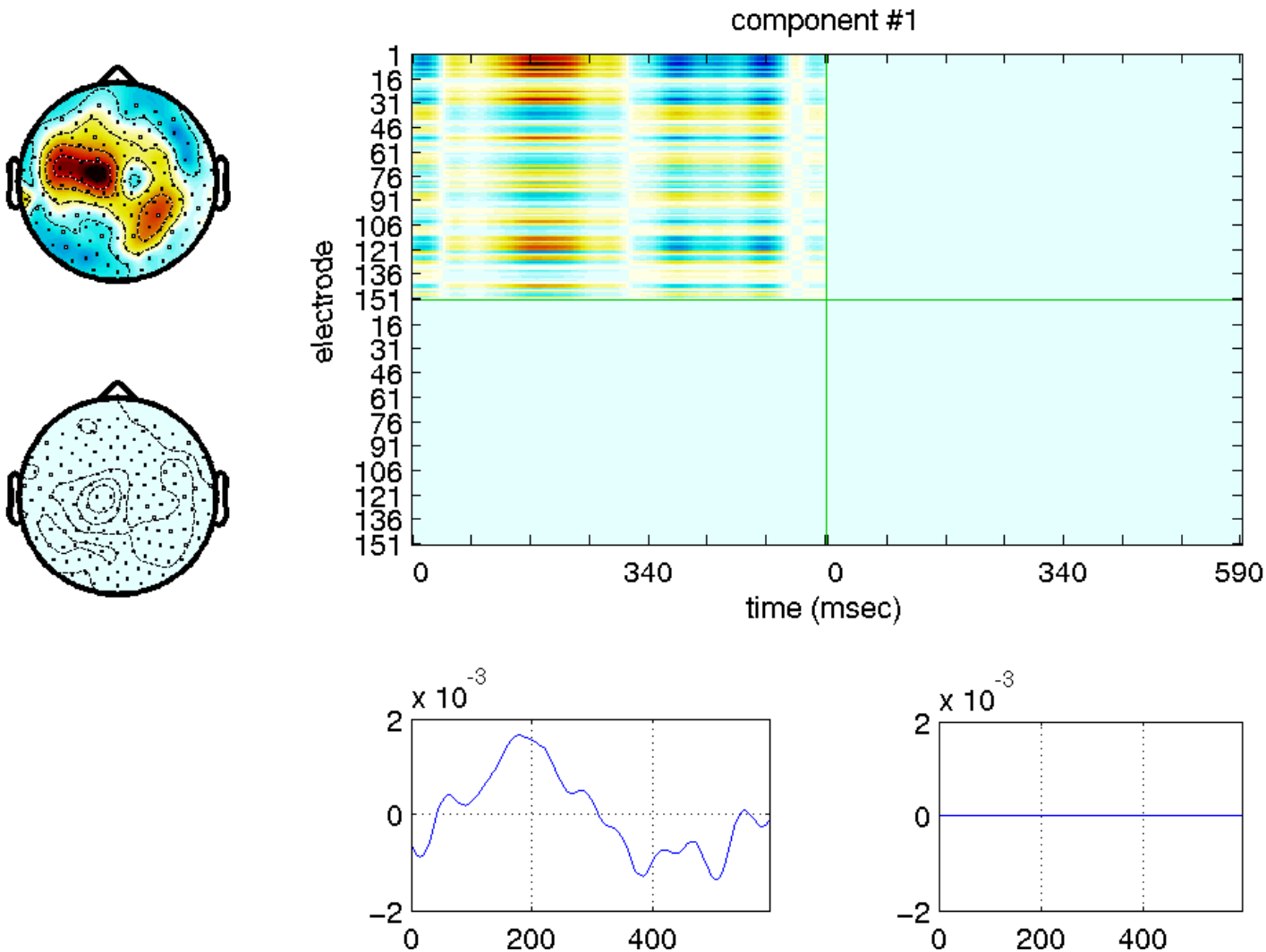


Classification performance (auditory)



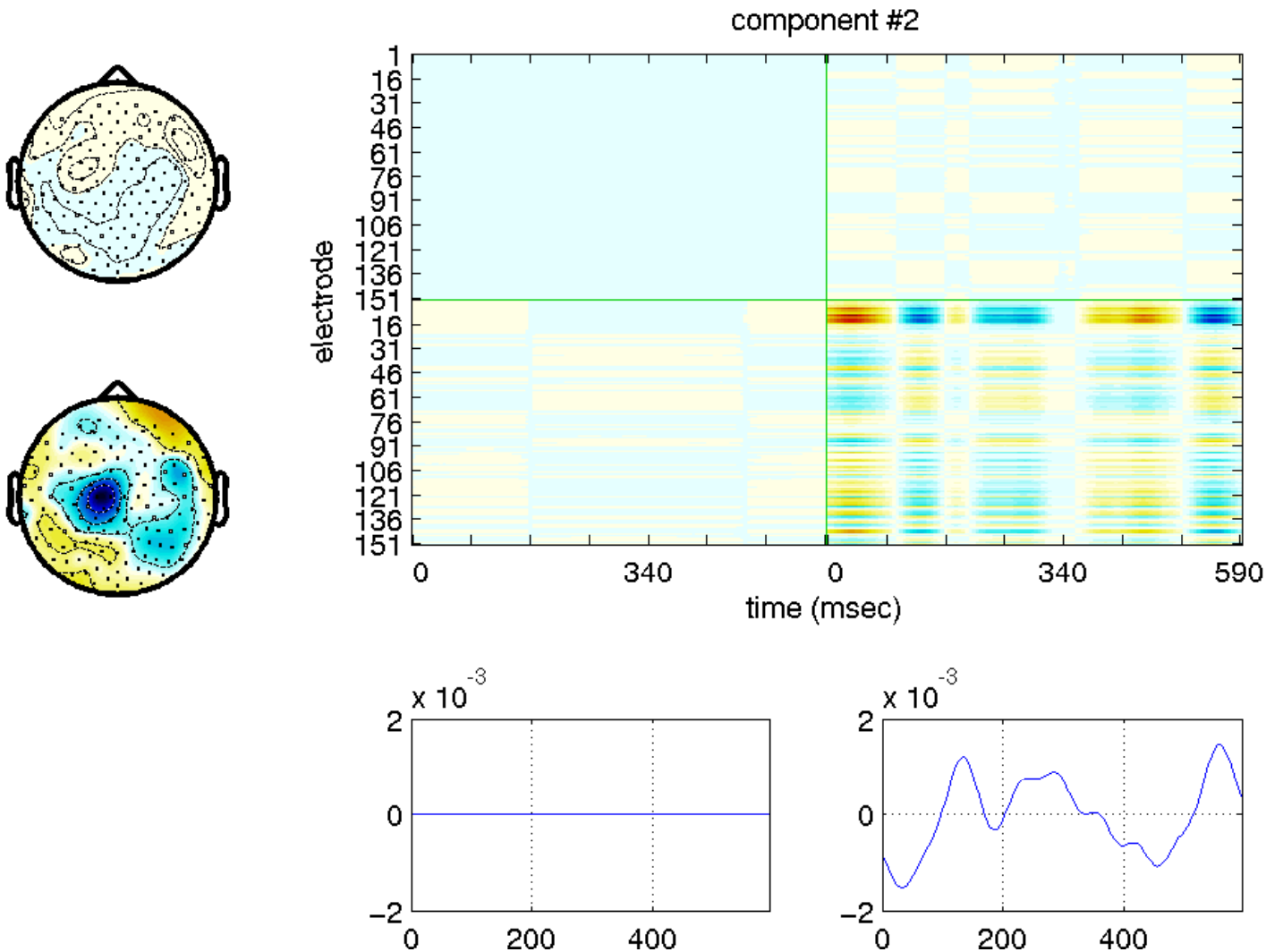
Example decomposition (tactile)

Subset of tactile MEG data (left little finger versus right little finger):



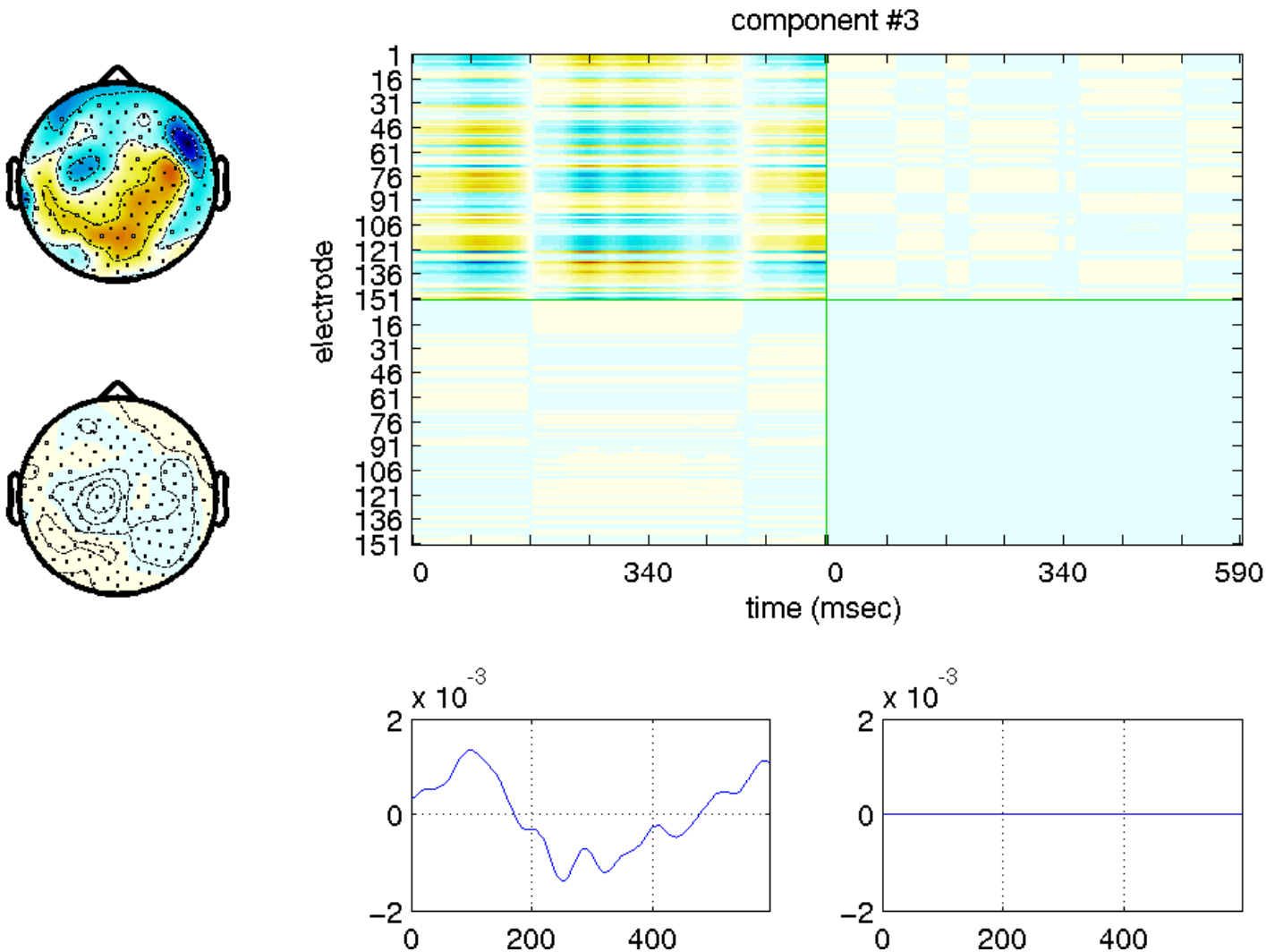
Example decomposition (tactile)

Subset of tactile MEG data (left little finger versus right little finger):



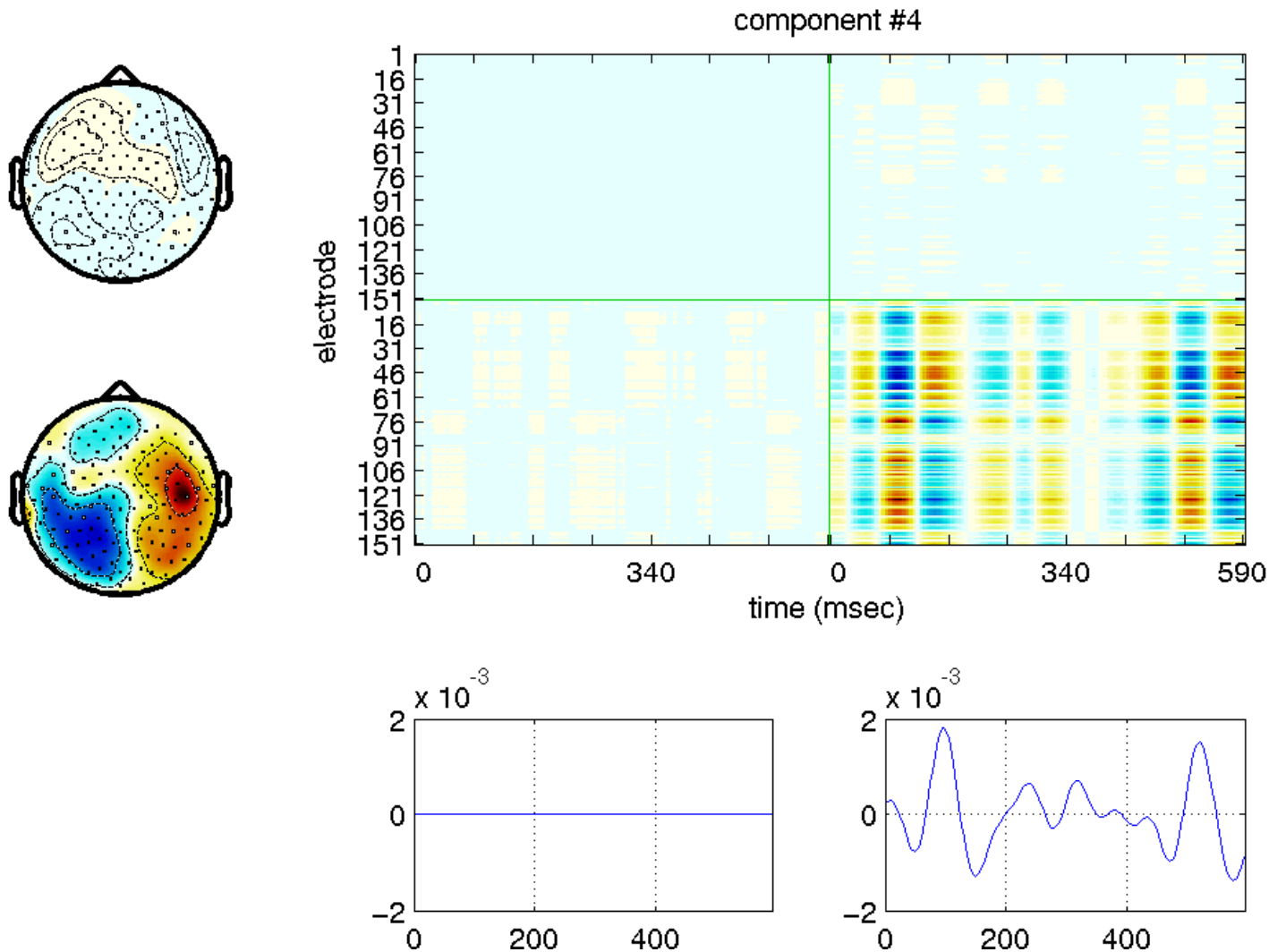
Example decomposition (tactile)

Subset of tactile MEG data (left little finger versus right little finger):

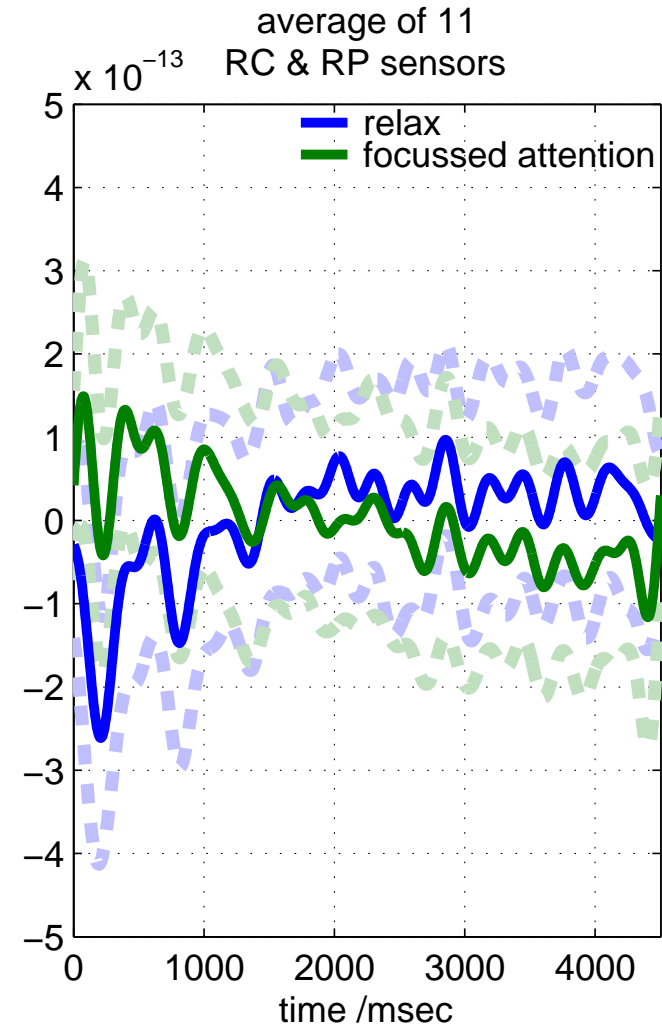
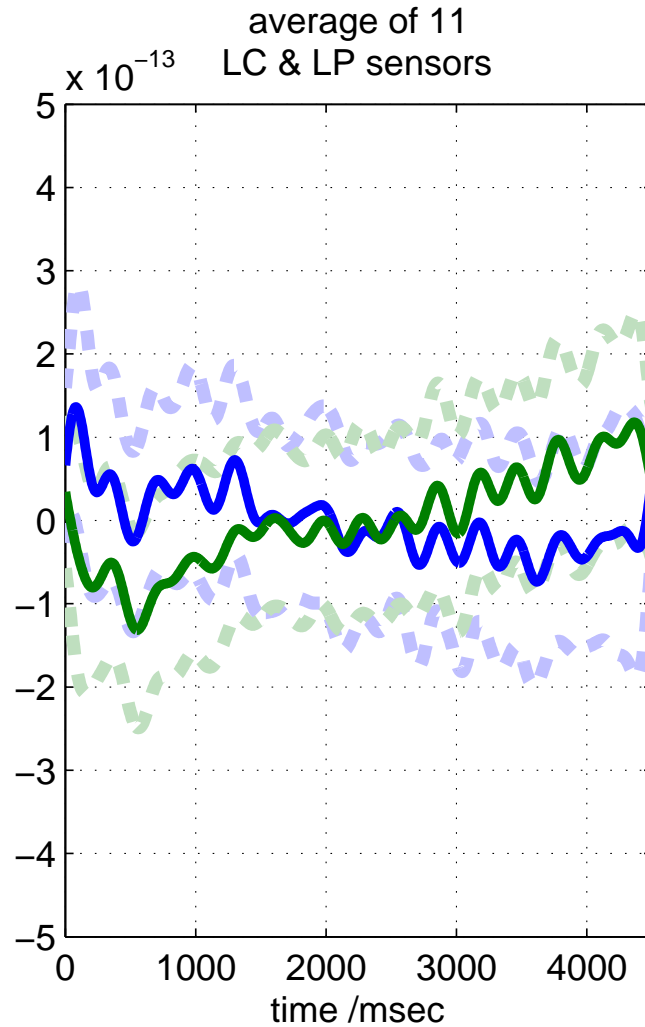


Example decomposition (tactile)

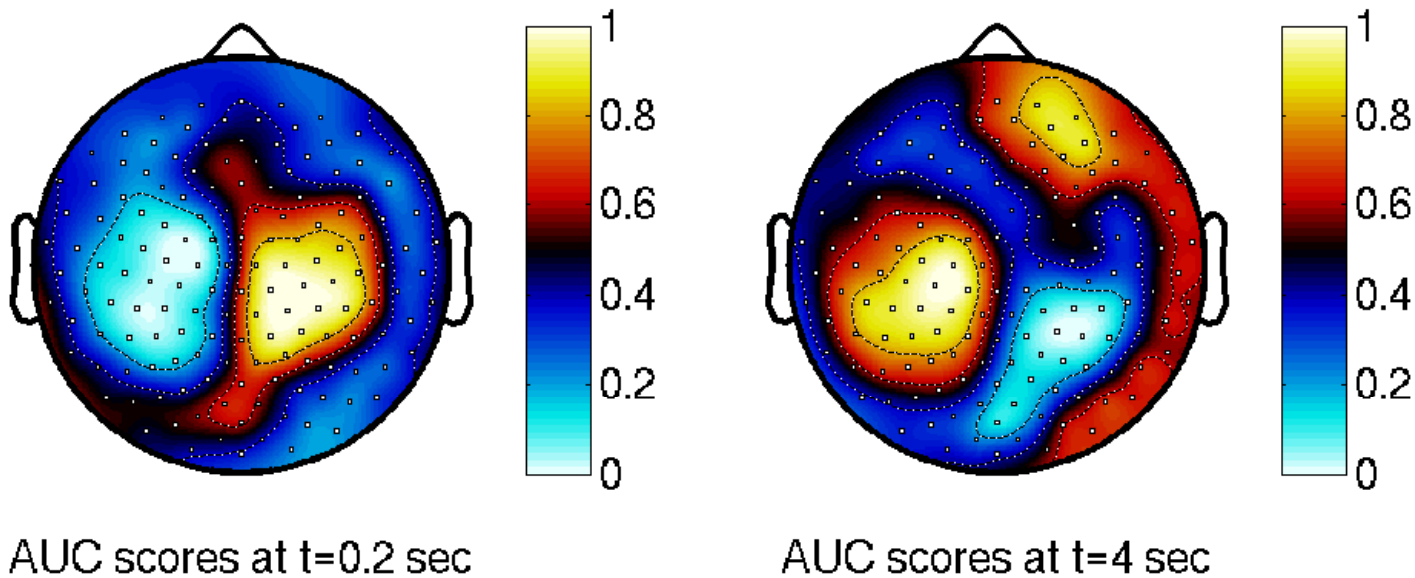
Subset of tactile MEG data (left little finger versus right little finger):



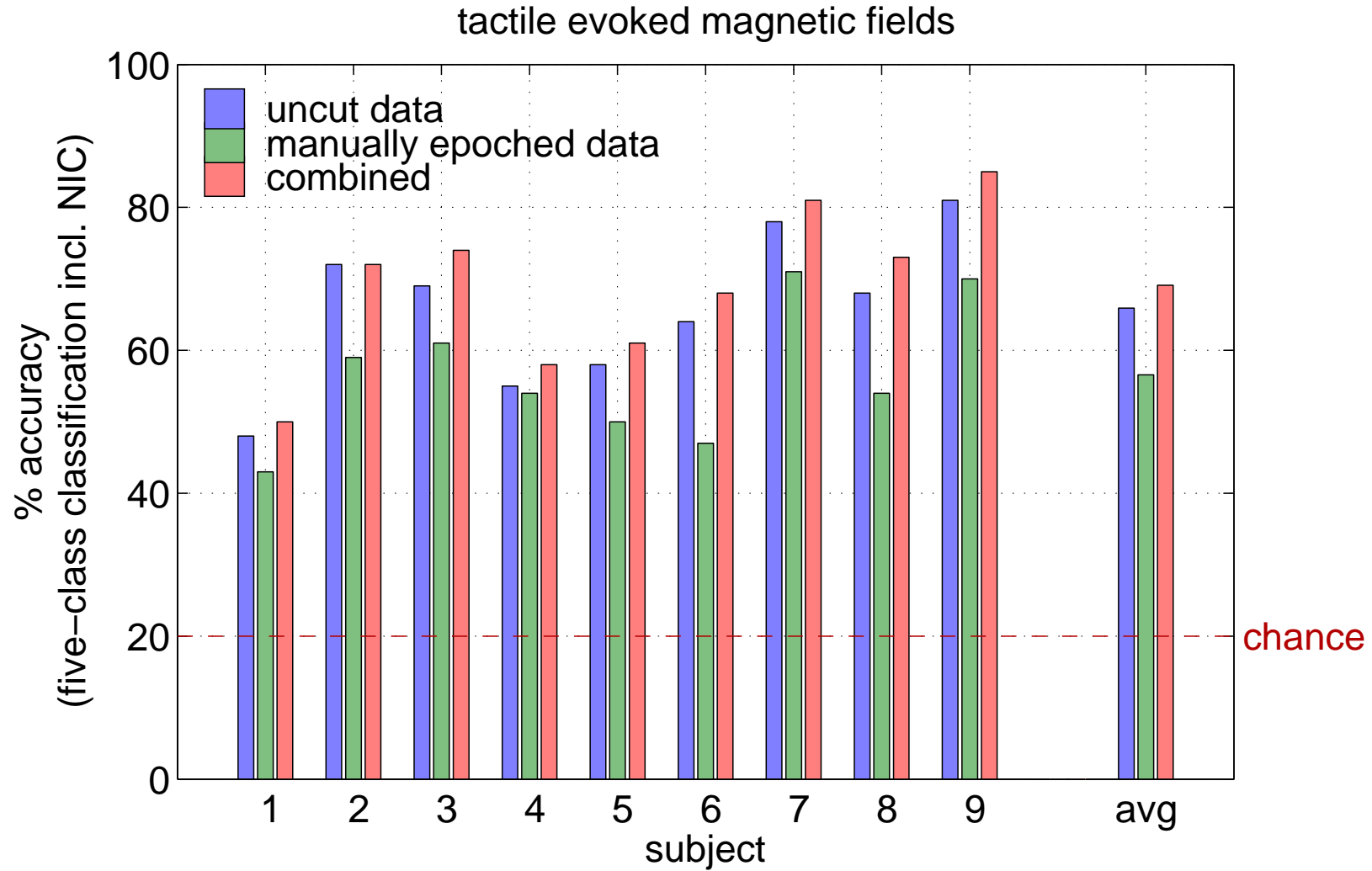
Slow waves indicating attention



Slow waves indicating attention

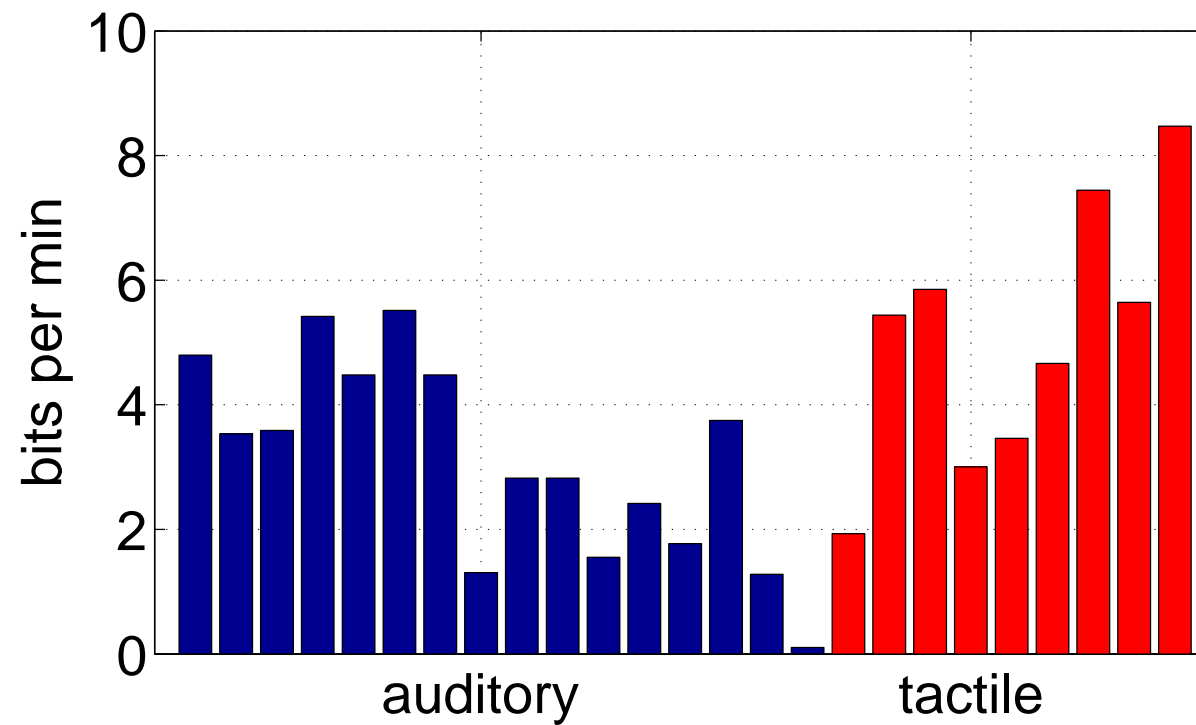


Classification performance



Bit rates

Information transfer rate
(Wolpaw's definition)
assuming 6 trials/minute





Conclusions



- It is possible to classify attention modulation of evoked responses to auditory and tactile stimuli, from 4.5-second signal segments.
- The stimuli are frequent (not an “oddball” paradigm).
- The classifier tends to rely more heavily on early (100–200 msec) components, although P300-like components are also useful.
- In addition, slow waves may help us to distinguish the control/no control problem.



Conclusions

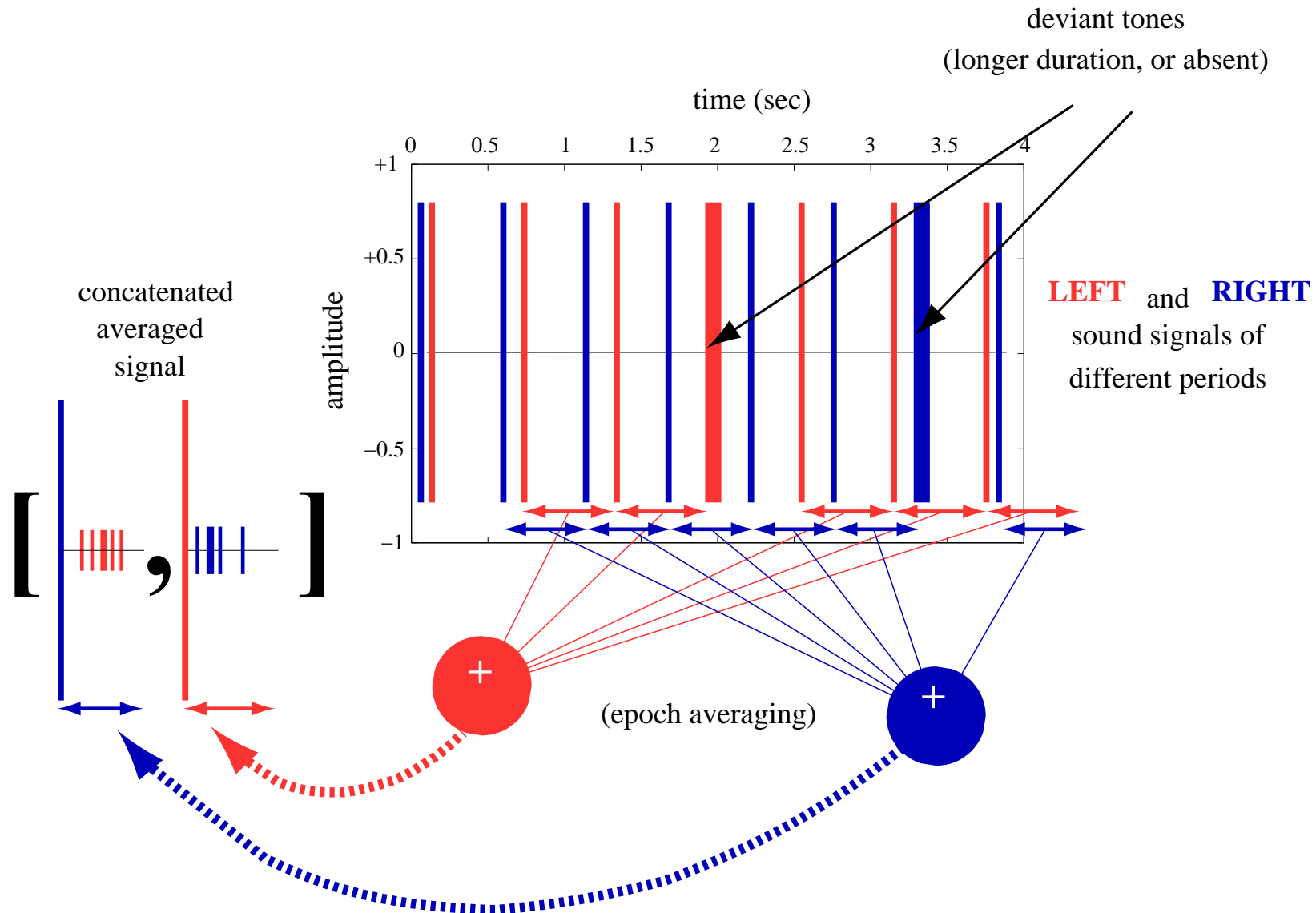


- It is possible to classify attention modulation of evoked responses to auditory and tactile stimuli, from 4.5-second signal segments.
- The stimuli are frequent (not an “oddball” paradigm).
- The classifier tends to rely more heavily on early (100–200 msec) components, although P300-like components are also useful.
- In addition, slow waves may help us to distinguish the control/no control problem.
- How well will this work as an online BCI?
- How much can the stimuli be speeded up?
- How high can the number of classes go, in the tactile experiment?
- How will speed and number of classes interact?

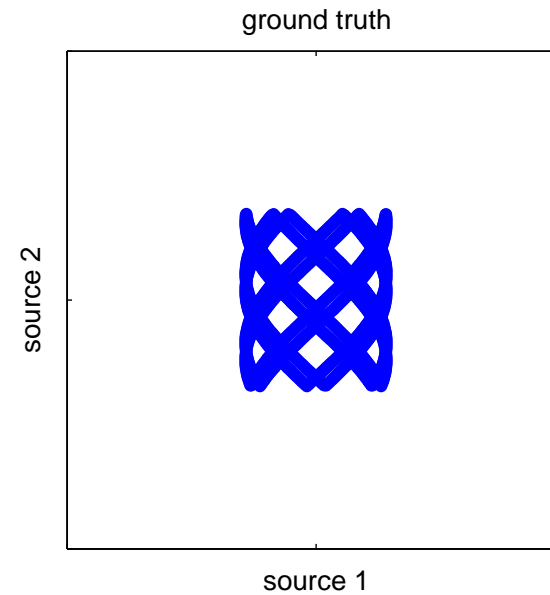
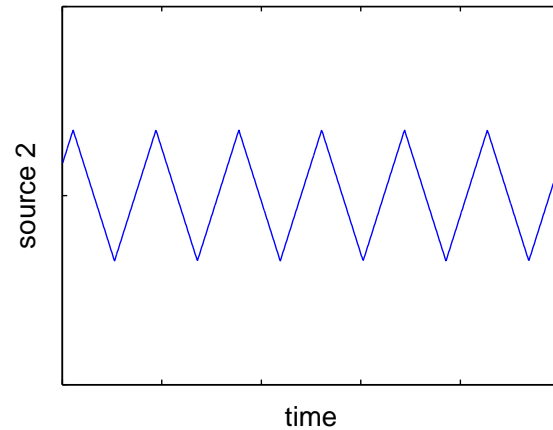
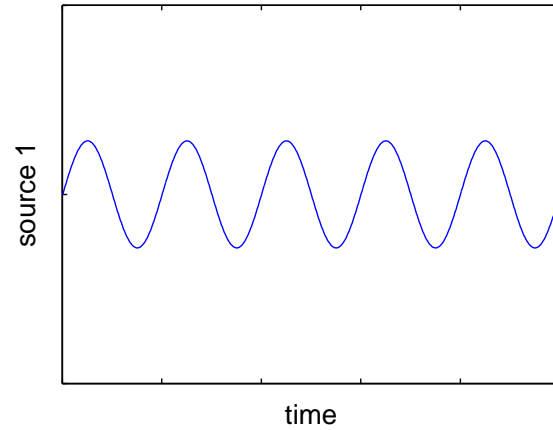


Thank you for your attention.

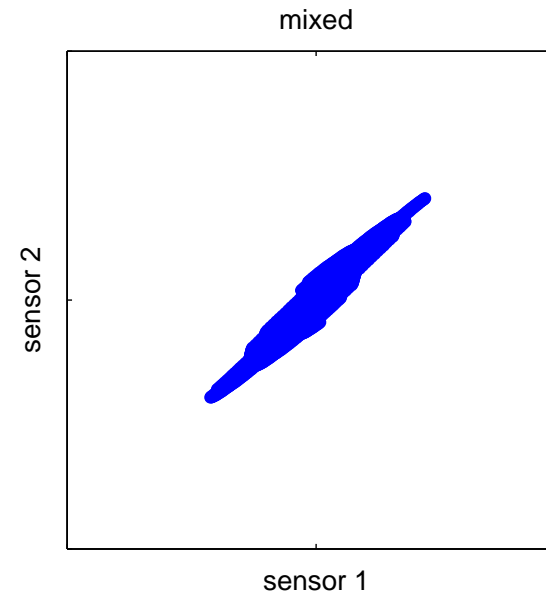
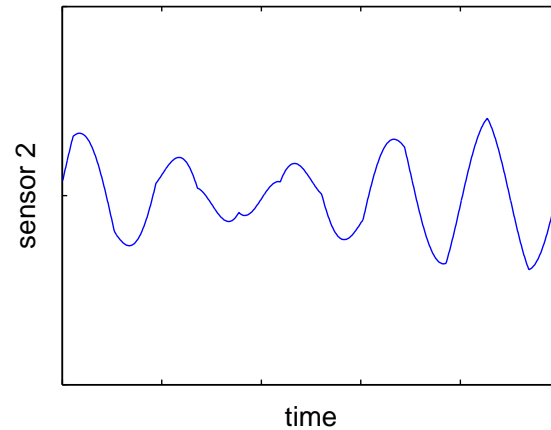
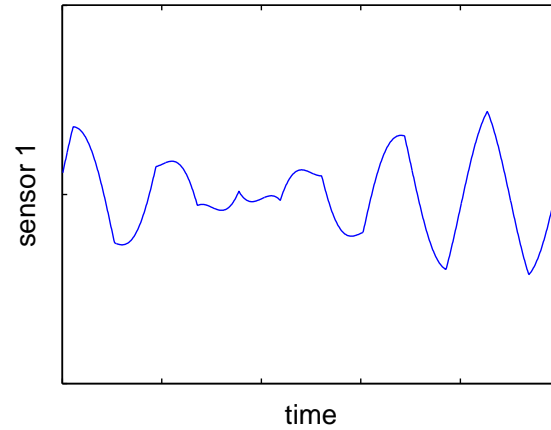
Auditory stimulus design



Whitening and rotation

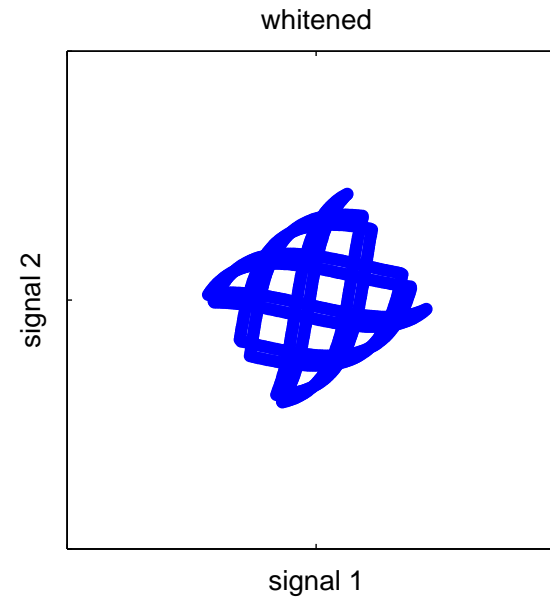
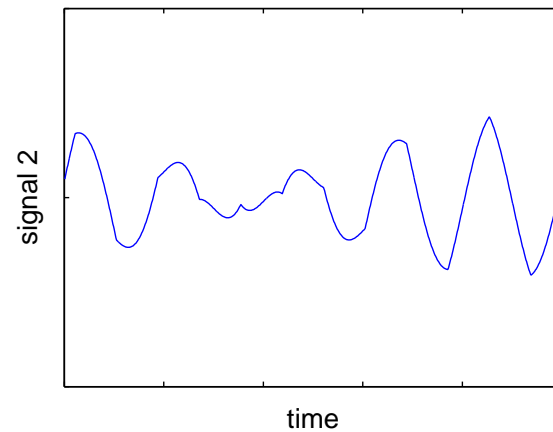
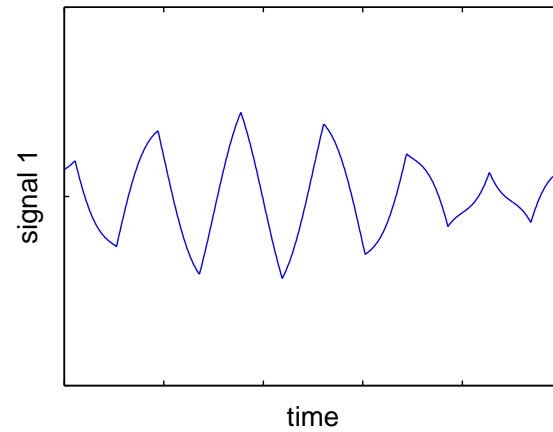


Whitening and rotation





Whitening and rotation



Whitening and rotation

