

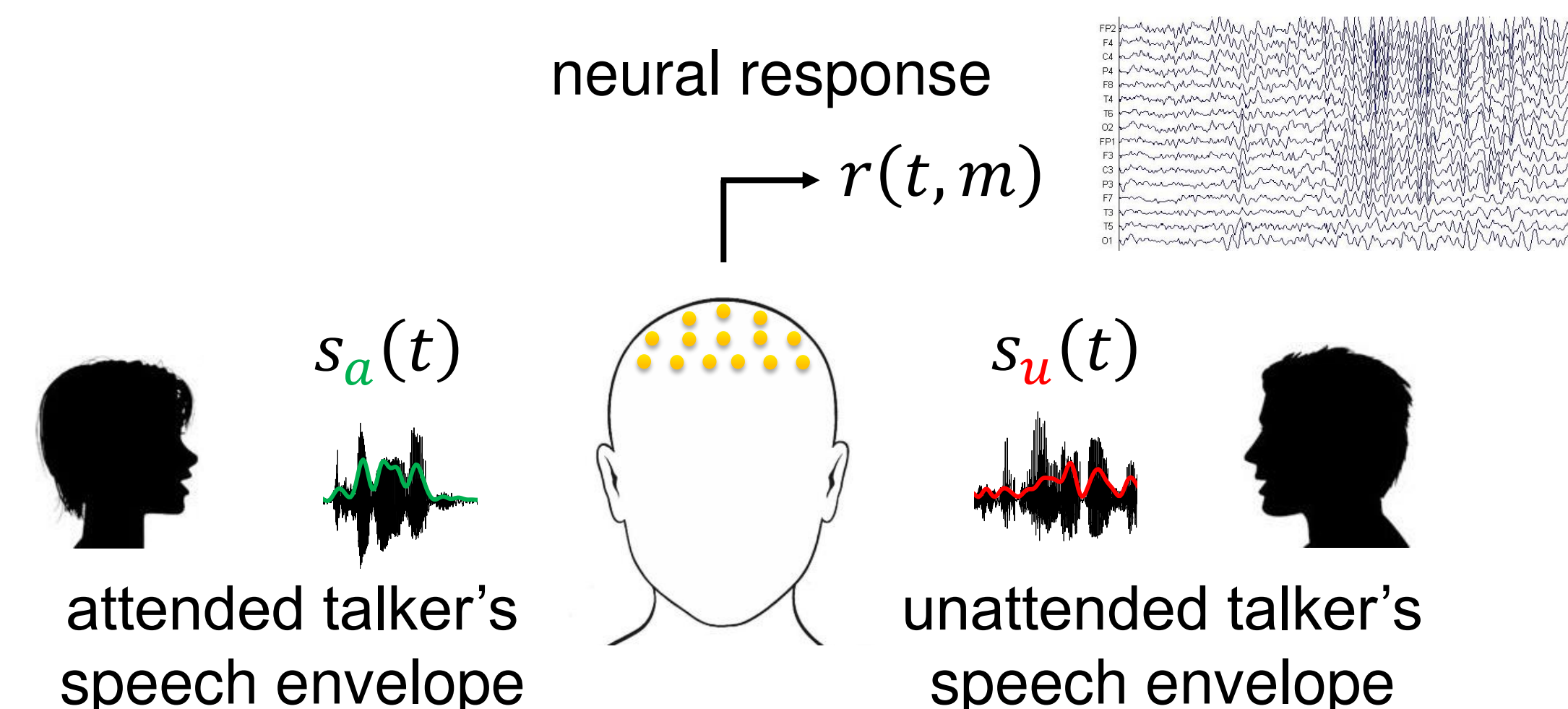
## Overview

- **Goal:** Given two simultaneous speech sources, detect which is being “attended” to and which is being “unattended” to, using the listener’s electroencephalography (EEG) data.
- **Innovation:** Extensions to conventional methods for auditory attention detection are proposed:
  - (1) selective channel deconvolution,
  - (2) maximally correlated multimodal projections,
  - (3) balanced correlation decoders.
- **Importance:** These are tools aimed to improve the understanding of how humans solve “the cocktail party problem.” Applications include, for example, attention driven acoustic beamforming for hearing prostheses.

## Paradigm

### Dichotic Listening Scenario

EEG electrodes  $\{m\}_1^M$



**Experiment:** An EEG subject is presented with 36 one-minute duration segments of competing spoken stories, one in each ear, and attends to only one.

## Auditory Attention Decoding Goal

- Using training data (35 of 36 segments in a leave-one-out cross validation approach), learn linear filters to reconstruct  $s_a(t)$  or  $s_u(t)$  from  $r(t, m)$ .
- Using test data (the remaining segment), predict  $\hat{s}_a(t)$  and  $\hat{s}_u(t)$  from  $r(t, m)$  and compare to ground truth  $s_a(t)$  and  $s_u(t)$ .

## Assumed Linear System

neural response  $r(t, m)$  additive noise  $\eta(t, m)$

$$r(t, m) = \int_{\tau=0}^{\tau_{max}} (a(\tau, m)s_a(t - \tau) + u(\tau, m)s_u(t - \tau)) + \eta(t, m)$$

attended and unattended channels

## Discrete Time Definitions

$$\mathbf{r} = \begin{bmatrix} r(t + \tau_{max}, 1) \\ \vdots \\ r(t, 1) \\ r(t + \tau_{max}, 2) \\ \vdots \\ r(t, M) \end{bmatrix} \quad \mathbf{s}_a = \begin{bmatrix} s_a(t + 2\tau_{max}) \\ \vdots \\ s_a(t + 1) \\ s_a(t) \end{bmatrix}$$

$$\mathbf{A} = \begin{bmatrix} a(0,1) & \dots & a(\tau_{max}, 1) & 0 & \dots & 0 \\ \vdots & & \vdots & & & \\ 0 & \dots & 0 & a(0,1) & \dots & a(\tau_{max}, 1) \\ a(0,2) & \dots & a(\tau_{max}, 2) & 0 & \dots & 0 \\ \vdots & & \vdots & & & \\ 0 & \dots & 0 & a(0, M) & \dots & a(\tau_{max}, M) \end{bmatrix}$$

Discrete time is assumed for both  $t$  and  $\tau$ .  $\eta$ ,  $s_u$ , and  $\mathbf{U}$  are analogous to the equations above, respectively.

**Compact model:**  $\mathbf{r} = \mathbf{A}\mathbf{s}_a + \mathbf{U}\mathbf{s}_u + \eta$

## Baseline Method: Stimulus Reconstruction

**MMSE:** *Minimum-Mean Square Error*

Learn a reconstruction filter,  $g(\tau, m)$ , to linearly combine spatiotemporal EEG observations using MMSE criteria.

$$\mathbf{g}_{a-MMSE} = \operatorname{argmin}_g \mathbb{E}\{|\mathbf{g}^T \mathbf{r} - s_a(t)|^2\}$$

## (1) Selective Channel Deconvolution

**MVDR:** *Minimum Variance Distortionless Response*

Reconstruct “attended” stimulus while minimizing any presence of deconvolved “unattended” stimulus & noise.

$$\mathbf{g}_{a-MVDR} = \operatorname{argmin}_g \mathbb{E}\{|\mathbf{g}^T (\mathbf{U}\mathbf{s}_u + \eta)|^2\}$$

s. t. :  $\mathbf{g}^T \mathbf{A} = [0, \dots, 0, 1]$

## (2) Multimodal Projections

**Multimodal (EEG & Speech) Data Structure**

$$\mathbf{r}_c = \begin{bmatrix} r(t, 1) \\ \vdots \\ r(t, M) \end{bmatrix} \quad \mathbf{s}_{a,c} = \begin{bmatrix} s_a(t + \tau_{max}) \\ \vdots \\ s_a(t) \end{bmatrix} \quad \mathbf{s}_{u,c} = \begin{bmatrix} s_u(t + \tau_{max}) \\ \vdots \\ s_u(t) \end{bmatrix}$$

**CCA:** *Canonical Correlation Analysis*

Bypass direct estimation of  $\hat{s}(t)$  to produce maximally correlated projections for both modalities, reducing overall number of feature dimensions by using spatial structure of EEG and temporal structure of stimuli.

$$\mathbf{g}_{r_a-CCA}, \mathbf{g}_{s_a-CCA} = \operatorname{argmin}_{\mathbf{g}_r, \mathbf{g}_s} \mathbb{E}\{|\mathbf{g}_r^T \mathbf{r}_c - \mathbf{g}_s^T \mathbf{s}_{a,c}|^2\}$$

s. t. :  $\mathbf{g}_r^T \mathbb{E}\{\mathbf{r}_c \mathbf{r}_c^T\} \mathbf{g}_r = \mathbf{g}_s^T \mathbb{E}\{\mathbf{s}_{a,c} \mathbf{s}_{a,c}^T\} \mathbf{g}_s = 1$

## (3) Balanced Correlation Decoders

**BMMSE:** *Balanced MMSE*

Optimize the reconstruction filter to the detection statistic by jointly maximizing the correlation & anticorrelation of the reconstruction with the “attended” & “unattended” stimuli, respectively.

$$\mathbf{g}_{a-BMMSE} = \operatorname{argmin}_g \mathbb{E}\{|\mathbf{g}^T \mathbf{r} - (s_a(t) - s_u(t))|^2\}$$

**BCCA:** *Balanced CCA*

Similar to BMMSE, but maximizing canonical correlation and canonical anticorrelation of multimodal projections.

$$\mathbf{g}_{r_a-BCCA}, \mathbf{g}_{s_a-BCCA} = \operatorname{argmin}_{\mathbf{g}_r, \mathbf{g}_s} \mathbb{E}\{|\mathbf{g}_r^T \mathbf{r}_c - \mathbf{g}_s^T (\mathbf{s}_{a,c} - \mathbf{s}_{u,c})|^2\}$$

s. t. :  $\mathbf{g}_s^T \mathbb{E}\{(\mathbf{s}_{a,c} - \mathbf{s}_{u,c})(\mathbf{s}_{a,c} - \mathbf{s}_{u,c})^T\} \mathbf{g}_s = 1$   
 $\mathbf{g}_r^T \mathbb{E}\{\mathbf{r}_c \mathbf{r}_c^T\} \mathbf{g}_r = 1$

## Auditory Attention Detection Statistic

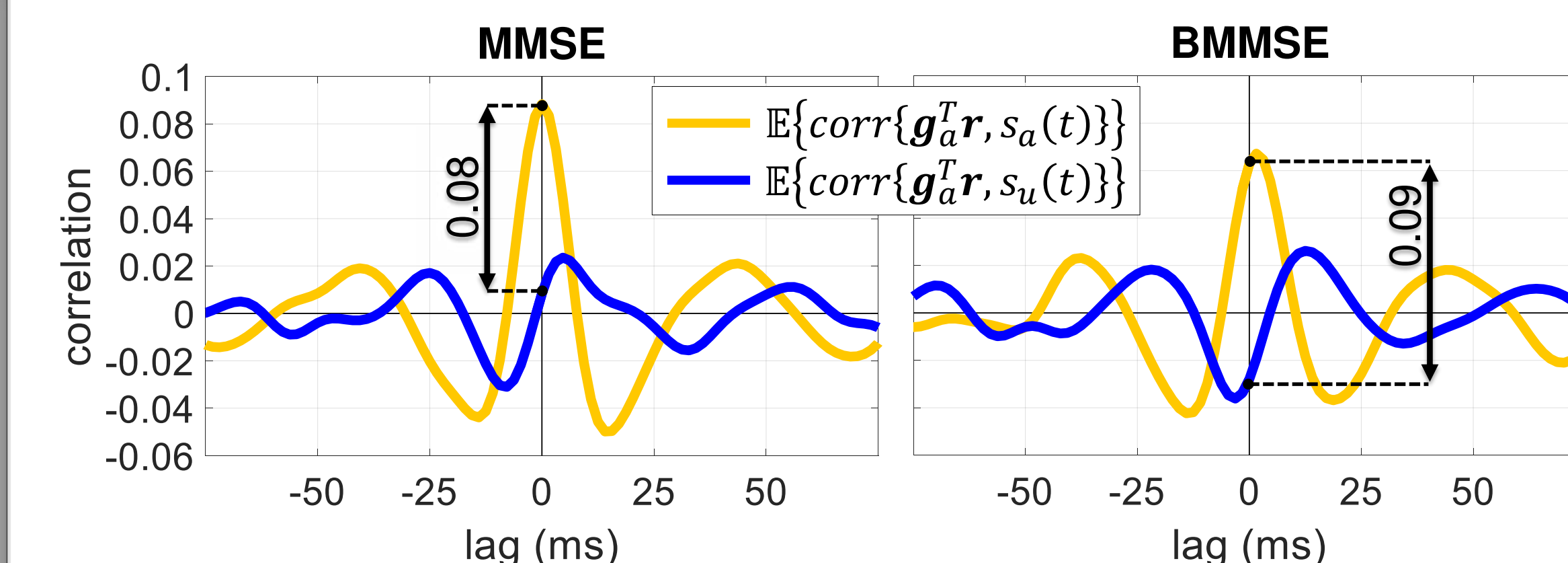
**MMSE, MVDR, and BMMSE Decoders**

Attended detected if:  $\operatorname{corr}\{\mathbf{g}_a^T \mathbf{r}, s_a(t)\} > \operatorname{corr}\{\mathbf{g}_a^T \mathbf{r}, s_u(t)\}$   
Unattended detected if:  $\operatorname{corr}\{\mathbf{g}_u^T \mathbf{r}, s_u(t)\} > \operatorname{corr}\{\mathbf{g}_u^T \mathbf{r}, s_a(t)\}$

**CCA and BCCA Decoders**

Attended detected if:  $\operatorname{corr}\{\mathbf{g}_{r_a}^T \mathbf{r}_c, \mathbf{g}_{s_a}^T \mathbf{s}_{a,c}\} > \operatorname{corr}\{\mathbf{g}_{r_a}^T \mathbf{r}_c, \mathbf{g}_{s_a}^T \mathbf{s}_{u,c}\}$   
Unattended detected if:  $\operatorname{corr}\{\mathbf{g}_{r_u}^T \mathbf{r}_c, \mathbf{g}_{s_u}^T \mathbf{s}_{u,c}\} > \operatorname{corr}\{\mathbf{g}_{r_u}^T \mathbf{r}_c, \mathbf{g}_{s_u}^T \mathbf{s}_{a,c}\}$

## Estimation & Detection Accuracy

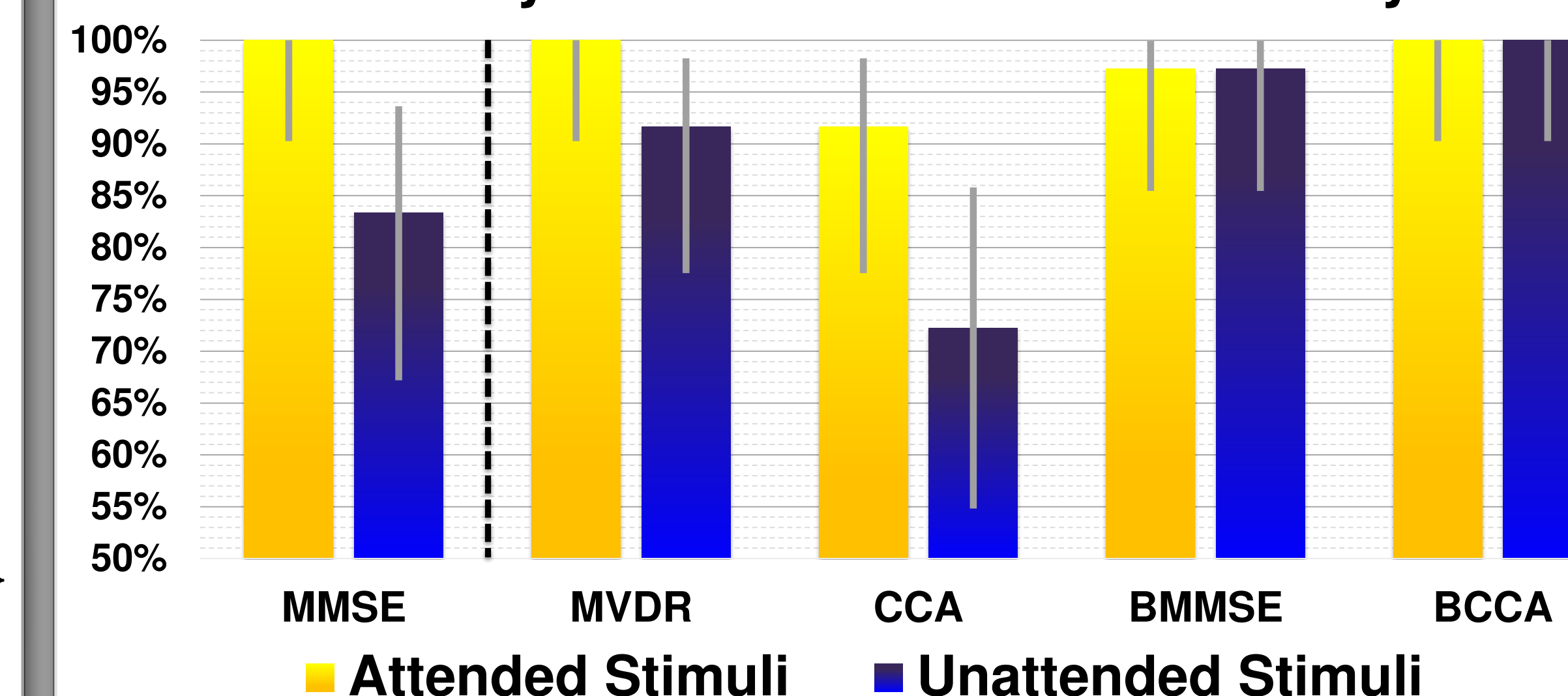


Average cross-correlation of reconstructed stimuli with true stimuli shows how BMMSE takes advantage of selective-attention structures embedded within EEG to maximize correlation separation used in the detection statistic.

	MMSE	MVDR	BMMSE
$\mathbb{E}\{\operatorname{corr}\{\mathbf{g}_a^T \mathbf{r}, s_a(t)\} - \operatorname{corr}\{\mathbf{g}_a^T \mathbf{r}, s_u(t)\}\}$	0.0800	0.0740	0.0908
$\mathbb{E}\{\operatorname{corr}\{\mathbf{g}_u^T \mathbf{r}, s_u(t)\} - \operatorname{corr}\{\mathbf{g}_u^T \mathbf{r}, s_a(t)\}\}$	0.0405	0.0367	0.0908
	CCA	BCCA	
$\mathbb{E}\{\operatorname{corr}\{\mathbf{g}_{r_a}^T \mathbf{r}_c, \mathbf{g}_{s_a}^T \mathbf{s}_{a,c}\} - \operatorname{corr}\{\mathbf{g}_{r_a}^T \mathbf{r}_c, \mathbf{g}_{s_a}^T \mathbf{s}_{u,c}\}\}$	0.0643	<b>0.0967</b>	
$\mathbb{E}\{\operatorname{corr}\{\mathbf{g}_{r_u}^T \mathbf{r}_c, \mathbf{g}_{s_u}^T \mathbf{s}_{u,c}\} - \operatorname{corr}\{\mathbf{g}_{r_u}^T \mathbf{r}_c, \mathbf{g}_{s_u}^T \mathbf{s}_{a,c}\}\}$	0.0336	<b>0.0967</b>	

Mean difference in detection statistics for decoder models.

## Auditory Attention Detection Accuracy



Mean detection accuracies for each decoding method. Error bars based on a binomial distribution ( $p = 0.5$ ,  $n = 36$ ,  $\alpha = 0.05$ ). Methods evaluated on one subject.

## Summary

- We developed 4 auditory attention decoders, each an extension to the traditional MMSE optimization criteria.
- See paper for details on how utilizing channel estimations via **MVDR** improves decoding accuracy.
- Best accuracy for attention detection occurs using **BCCA** by balancing canonical projections using both stimuli.

