

Neural Encoding of Speech in Auditory Cortex

Jonathan Z. Simon

Department of Biology

Department of Electrical & Computer Engineering

Institute for Systems Research

University of Maryland

Acknowledgements

Current (Simon Lab & Affiliates)

Francisco Cervantes
Natalia Lapinskaya
Mahshid Najafi
Alex Presacco
Krishna Puvvada
Lisa Uible
Peng Zan

Past (Simon Lab & Affiliate Labs)

Nayef Ahmar
Sahar Akram
Murat Aytekin
Claudia Bonin
Maria Chait
Marisel Villafane Delgado
Kim Drnec
Nai Ding
Victor Grau-Serrat
Julian Jenkins
David Klein
Ling Ma

Kai Sum Li
Huan Luo
Raul Rodriguez
Ben Walsh
Juanjuan Xiang
Jiachen Zhuo

Collaborators

Pamela Abshire
Samira Anderson
Behtash Babadi
Catherine Carr
Monita Chatterjee
Alain de Cheveigné
Didier Depireux
Mounya Elhilali
Bernhard Englitz
Jonathan Fritz
Cindy Moss
David Poeppel
Shihab Shamma

Past Postdocs & Visitors

Aline Gesualdi Manhães
Dan Hertz
Yadong Wang

Undergraduate Students

Abdulaziz Al-Turki
Nicholas Asendorf
Sonja Bohr
Elizabeth Camenga
Corinne Cameron
Julien Dagenais
Katya Dombrowski
Kevin Hogan
Kevin Kahn
Alexandria Miller
Isidora Ranovadovic
Andrea Shome
Madeleine Varmer
Ben Walsh

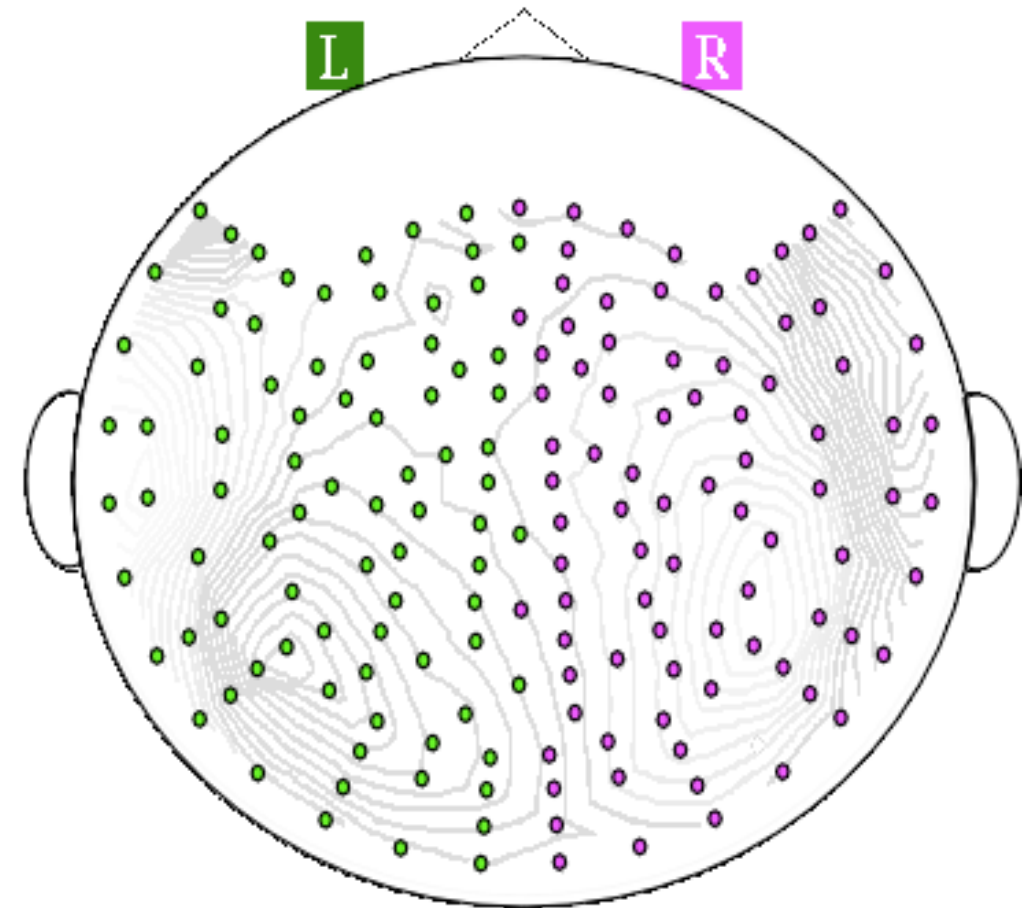
Funding NIH (**NIDCD**, NIA, NIBIB); USDA

Outline

- Magnetoencephalography (MEG)
- Cortical Representations of Speech
 - Encoding vs. Decoding
 - Attended vs. Unattended Speech
- Work in Progress
 - Attentional Dynamics
 - Aging and the Cocktail Party Problem
 - Foreground vs. Background

Magnetoencephalography

- Non-invasive, Passive, Silent Neural Recordings
- Simultaneous Whole-Head Recording (~200 sensors)
- Sensitivity
 - high: ~100 fT (10^{-13} Tesla)
 - low: $\sim 10^4 - \sim 10^6$ neurons
- Temporal Resolution: ~1 ms
- Spatial Resolution
 - coarse: ~1 cm
 - ambiguous



Neural Signals & MEG

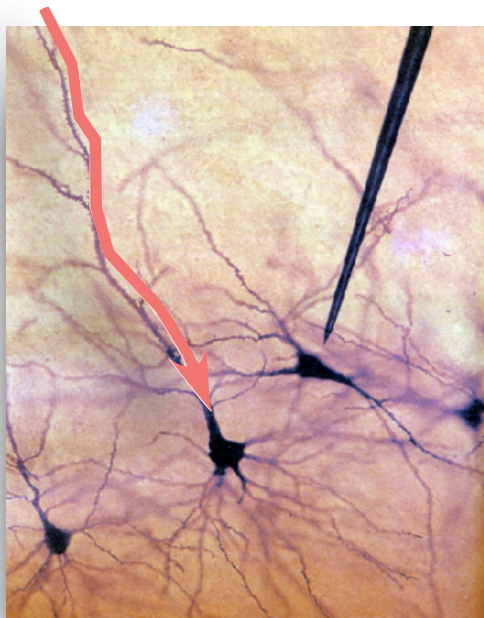
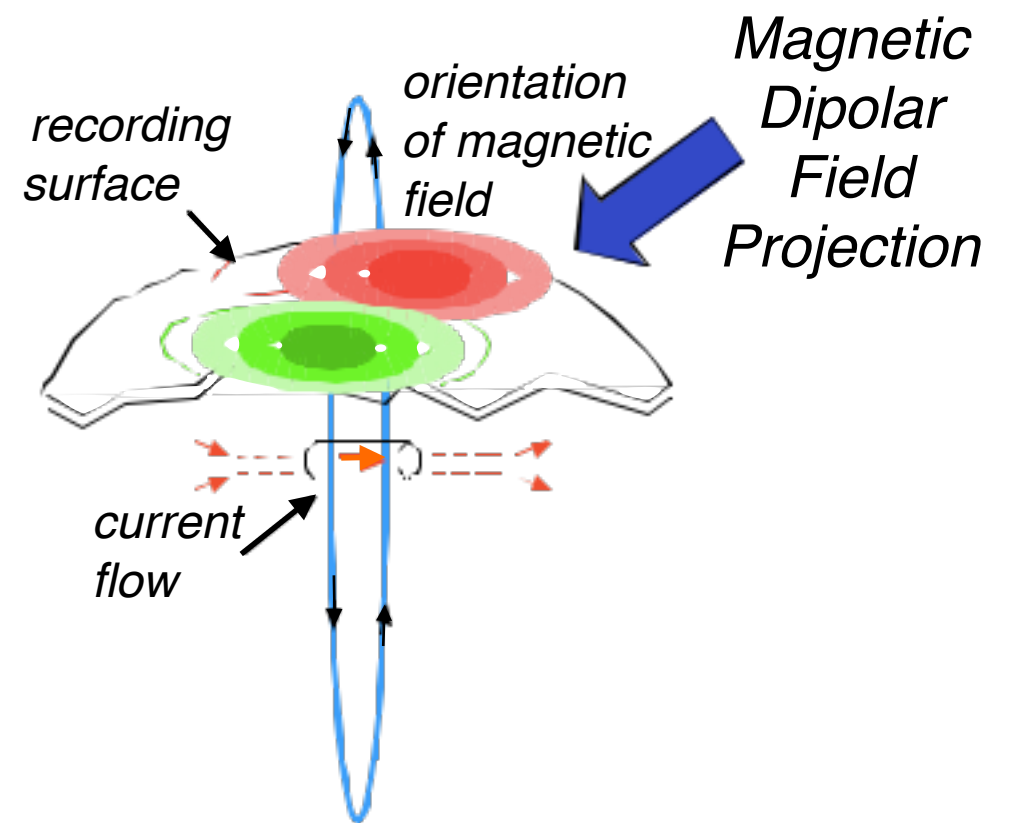
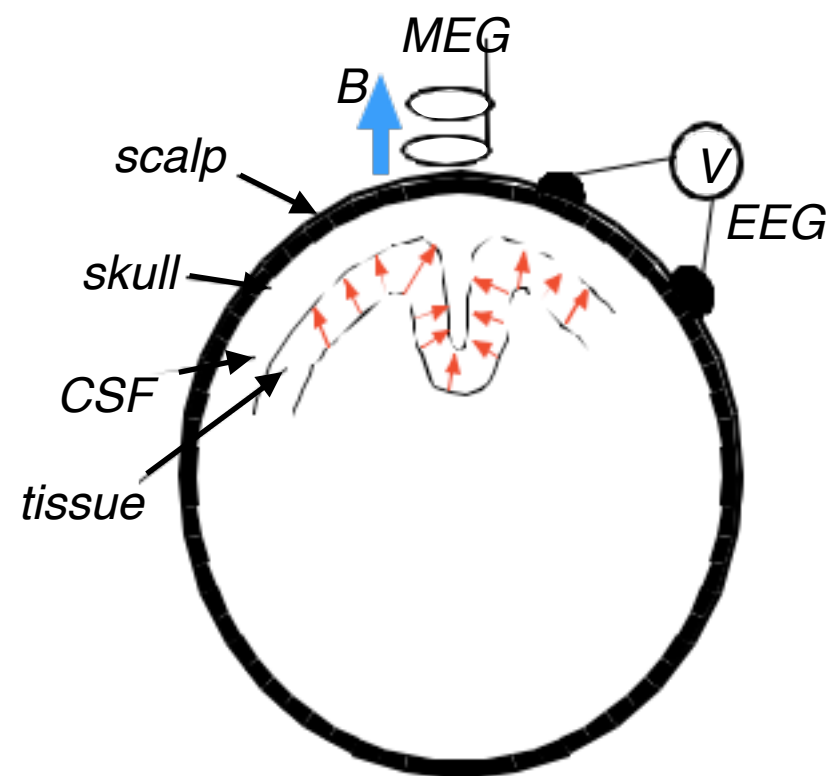


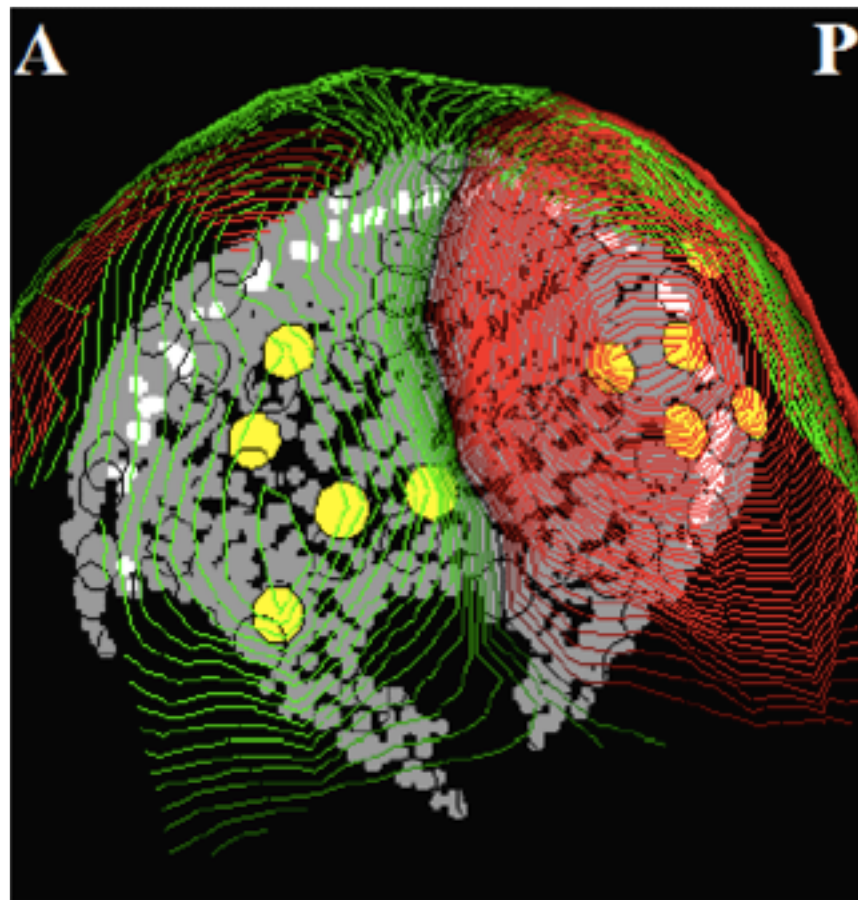
Photo by Fritz Goro



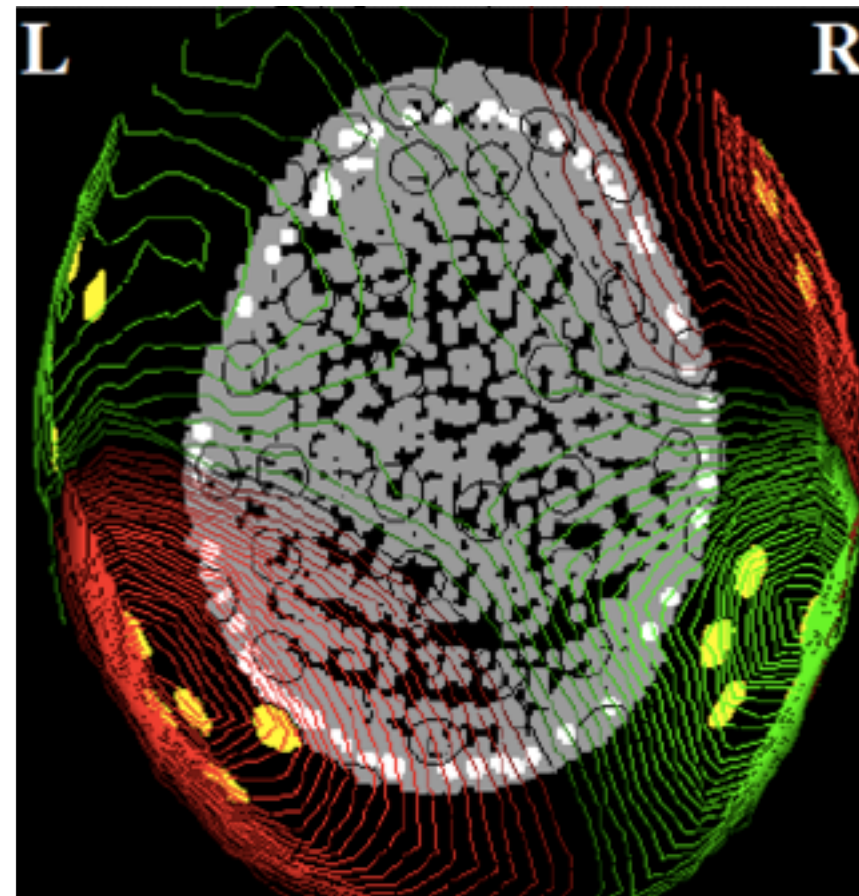
- Direct electrophysiological measurement
 - not hemodynamic
 - real-time
- No unique solution for distributed source

- Measures spatially synchronized cortical activity
- Fine temporal resolution (~ 1 ms)
- Moderate spatial resolution (~ 1 cm)

MEG Auditory Field



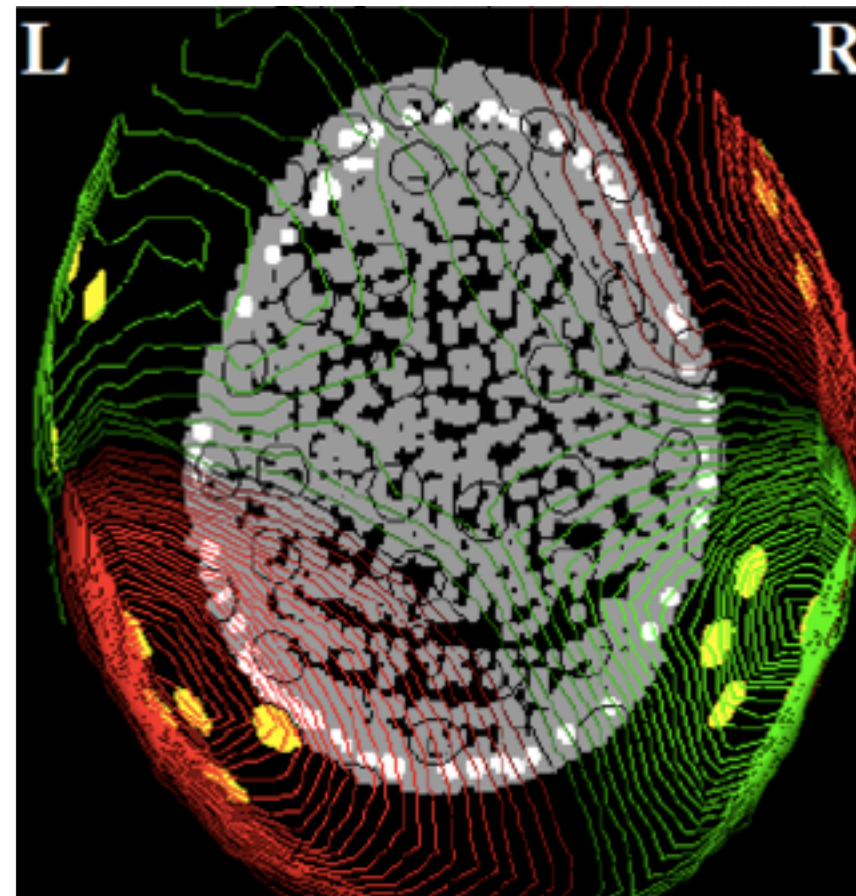
Sagittal View



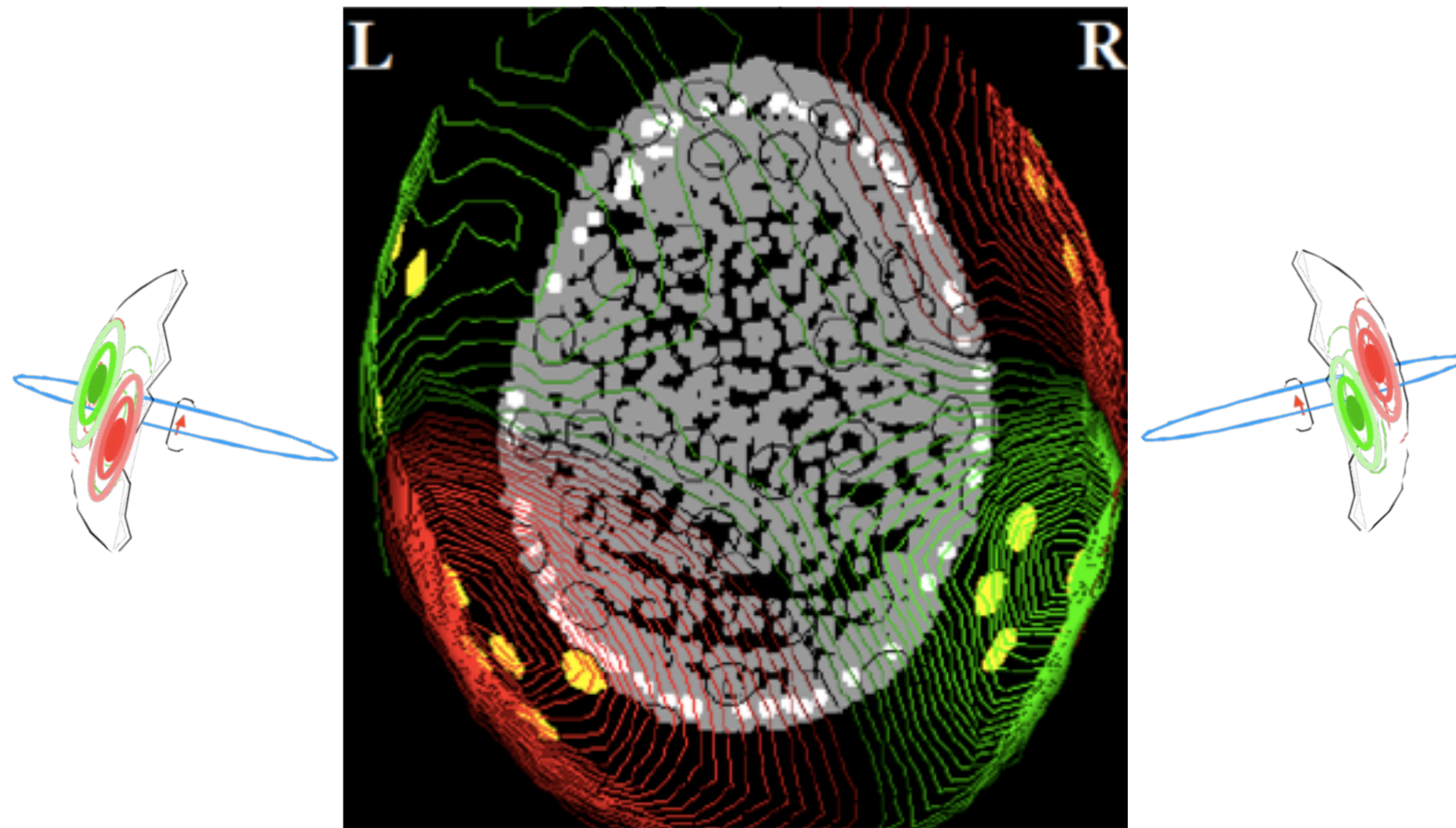
Axial View

Strongly
Lateralized

MEG Auditory Field



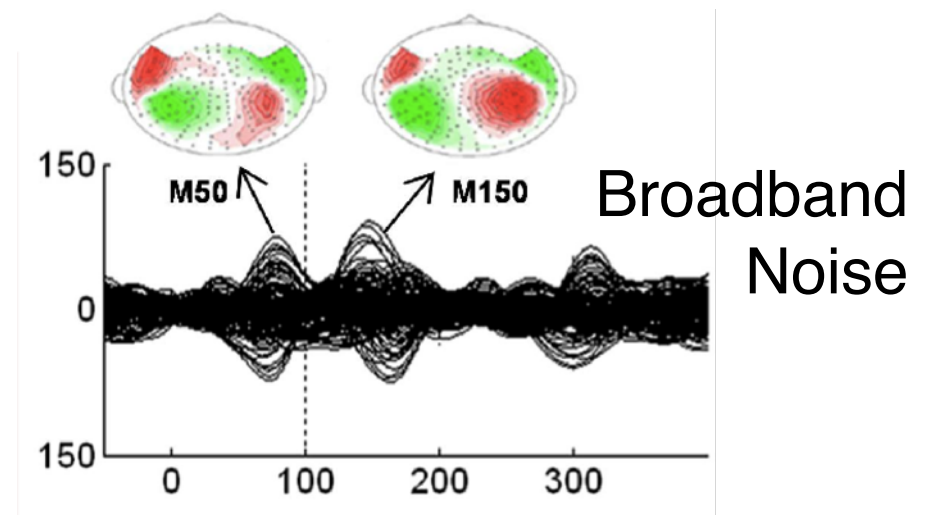
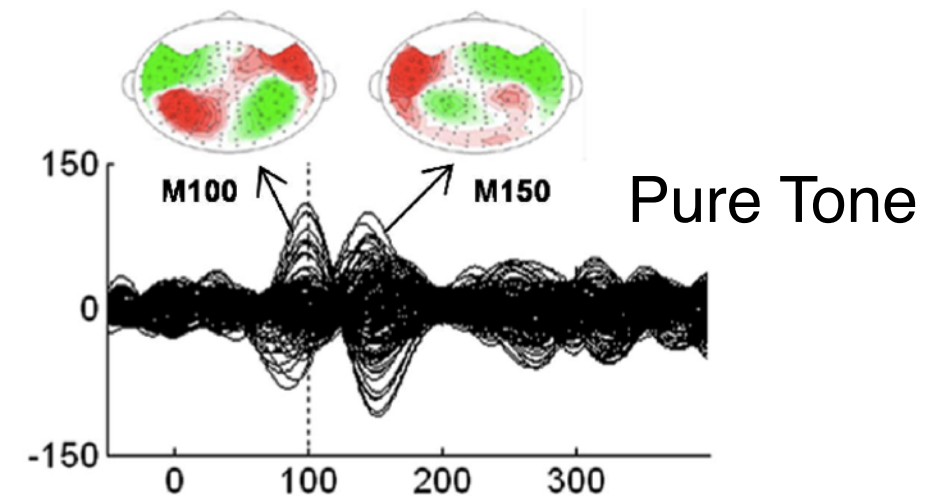
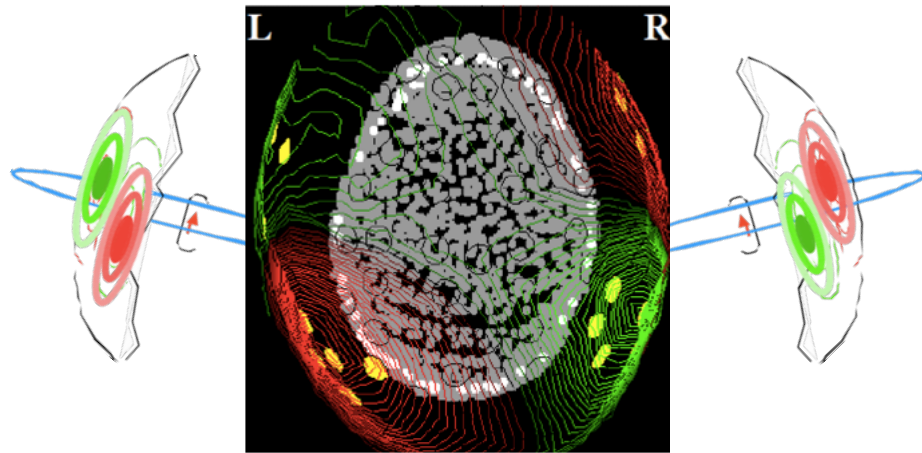
MEG Auditory Field



Time Course of MEG Responses

Auditory Evoked Responses

- MEG Response Patterns Time-Locked to Stimulus Events
- Robust
- Strongly Lateralized

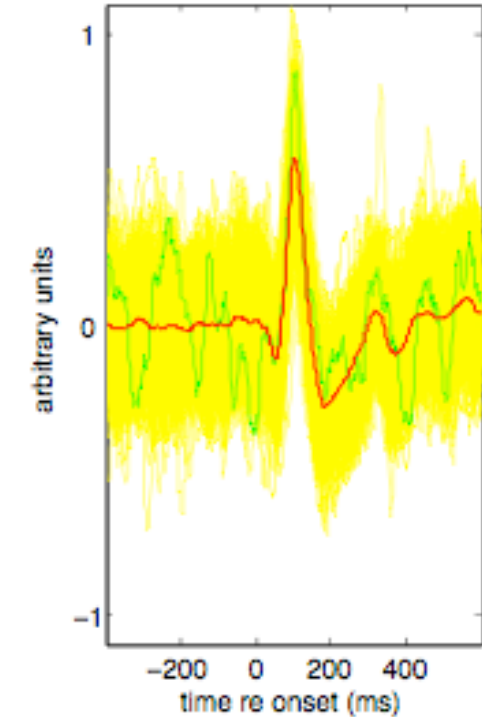
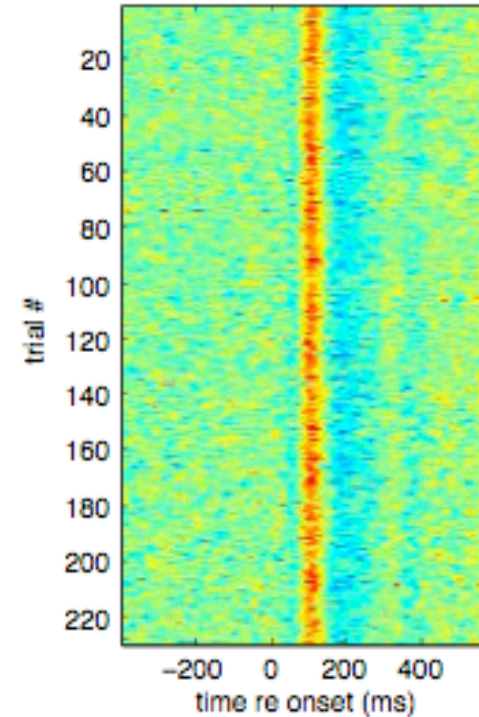
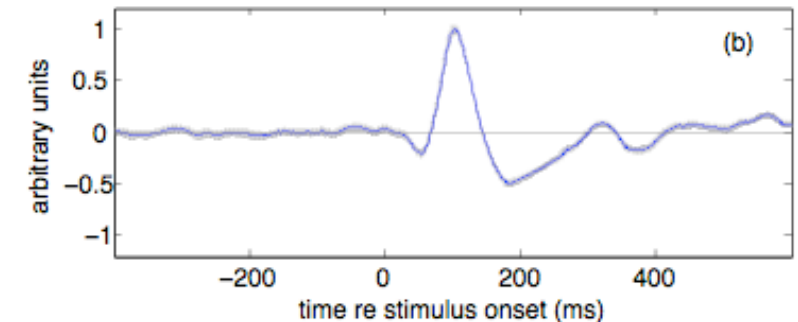
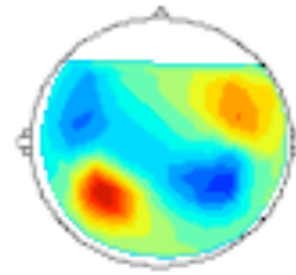


MEG Component Analysis

- Data driven spatial filtering:
many available methods—IICA, PCA, DSS
- Generate spatial filters & their outputs (“components”)
- DSS: Denoising Source Separation:
Särelä & Valpola (2005)
- DSS components ordered by reproducibility
 - 1st component “maximally reproducible” = most stimulus driven

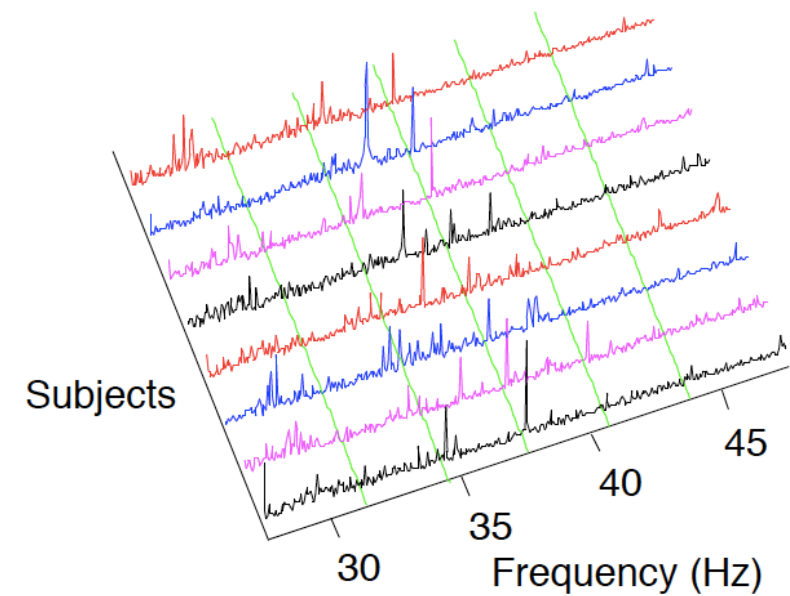
Component Analysis

- Each component has both spatial and temporal profile
- Data driven, e.g., PCA, ICA, DSS
- DSS: ordered by trial-to-trial reproducibility
- → Spatial Filter, e.g. for single trials
- Can analyze temporal processing separately from anatomical origin

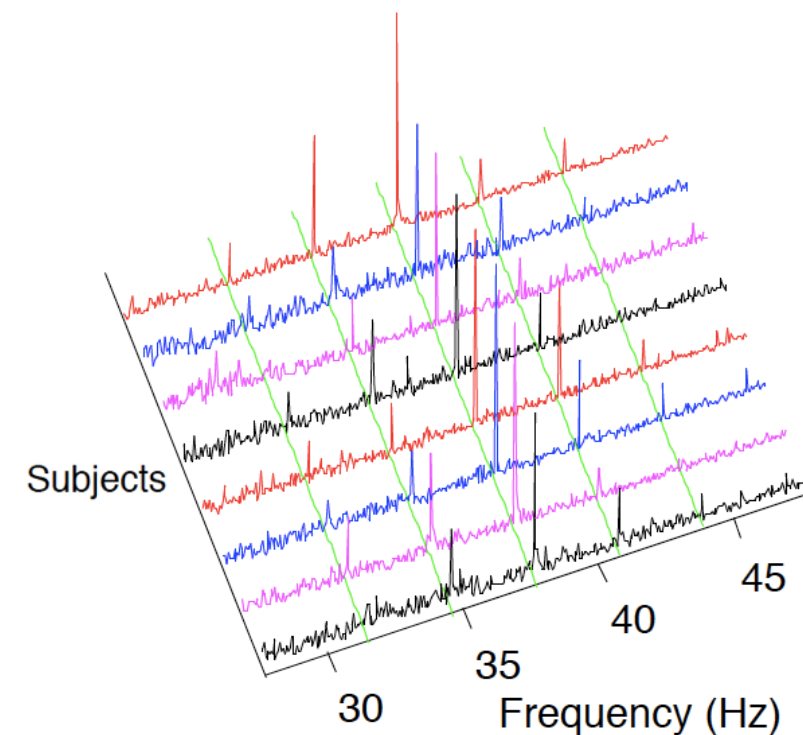


DSS Example: Spectral

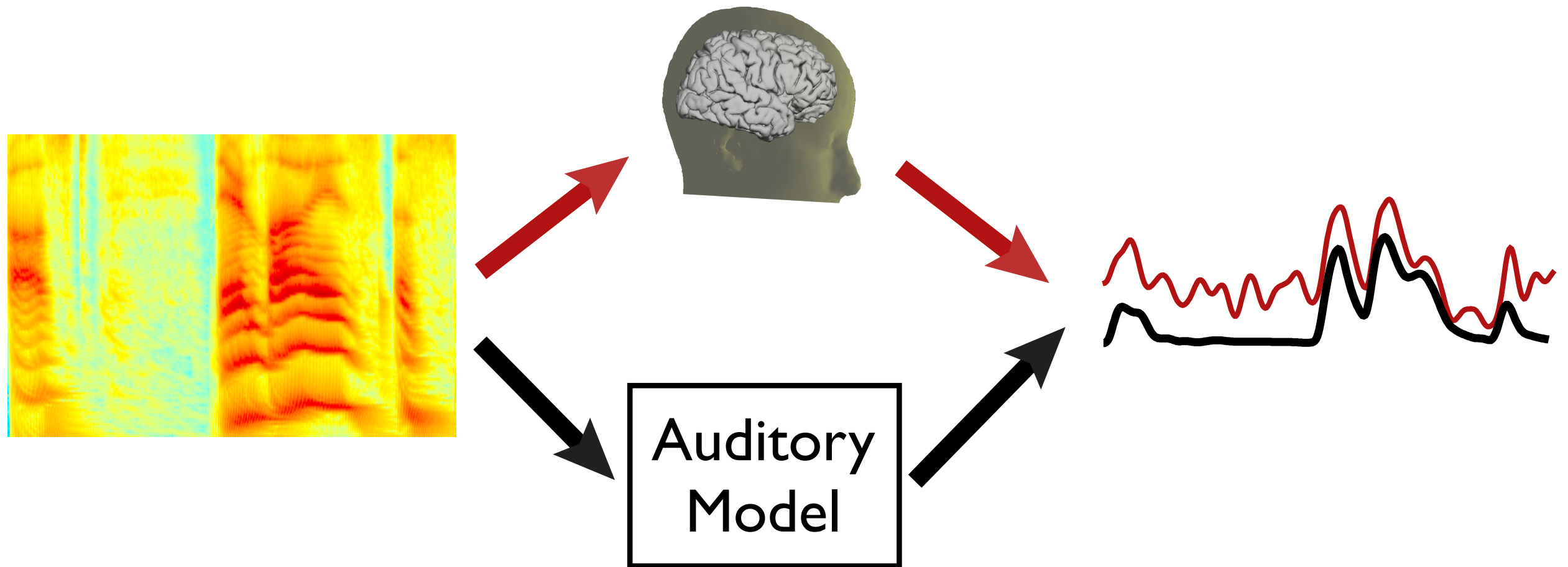
Frequency Spectrum before DSS



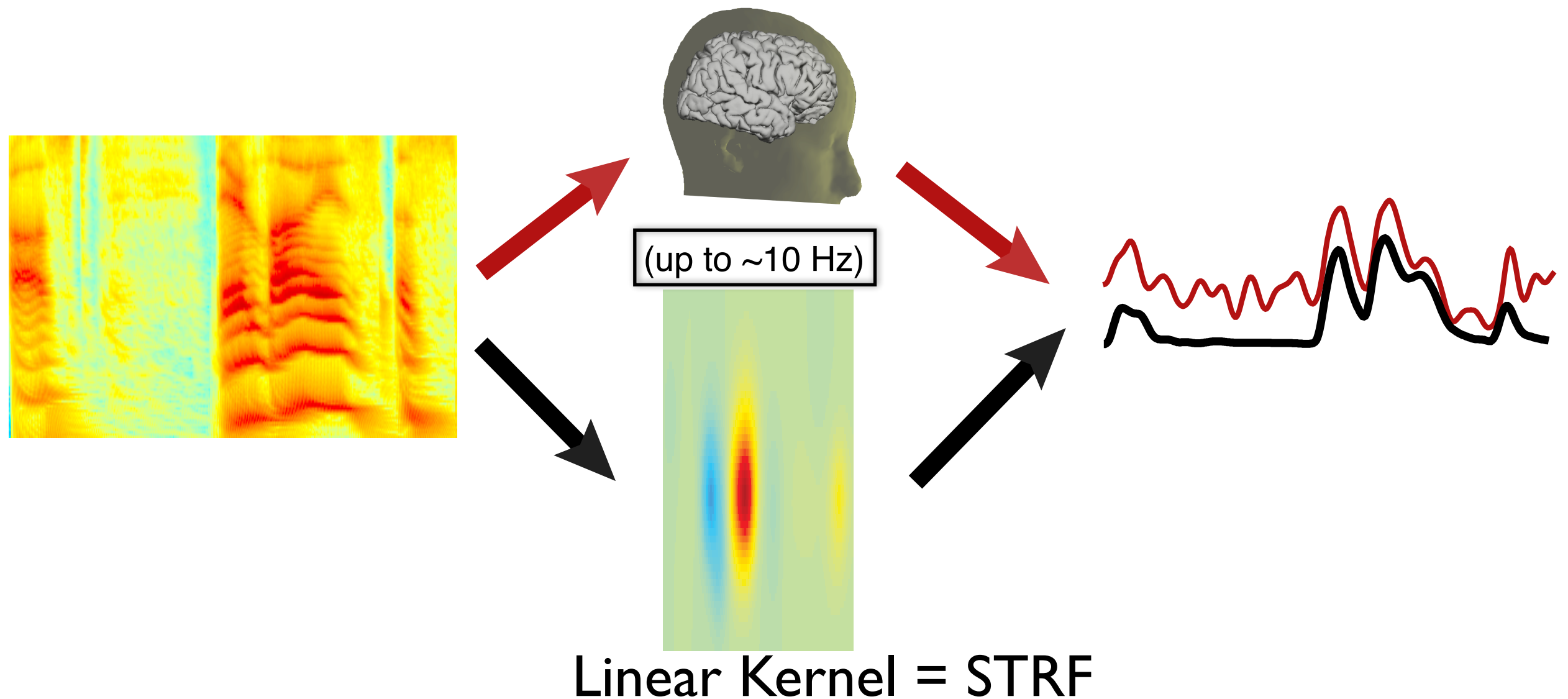
Frequency Spectrum after DSS



MEG Responses to Speech Modulations

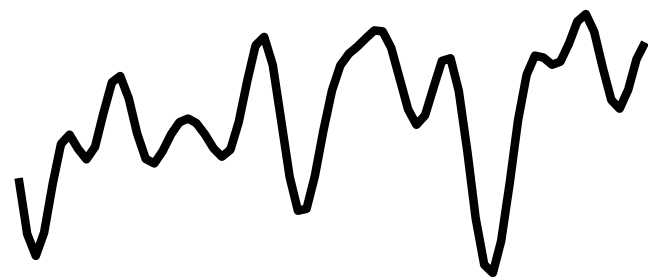


MEG Responses Predicted by STRF Model

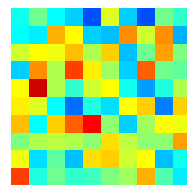


Neural Reconstruction of Speech Envelope

Speech Envelope

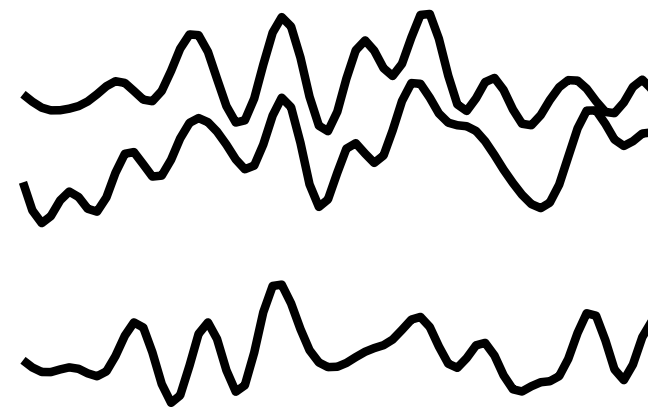


Decoder

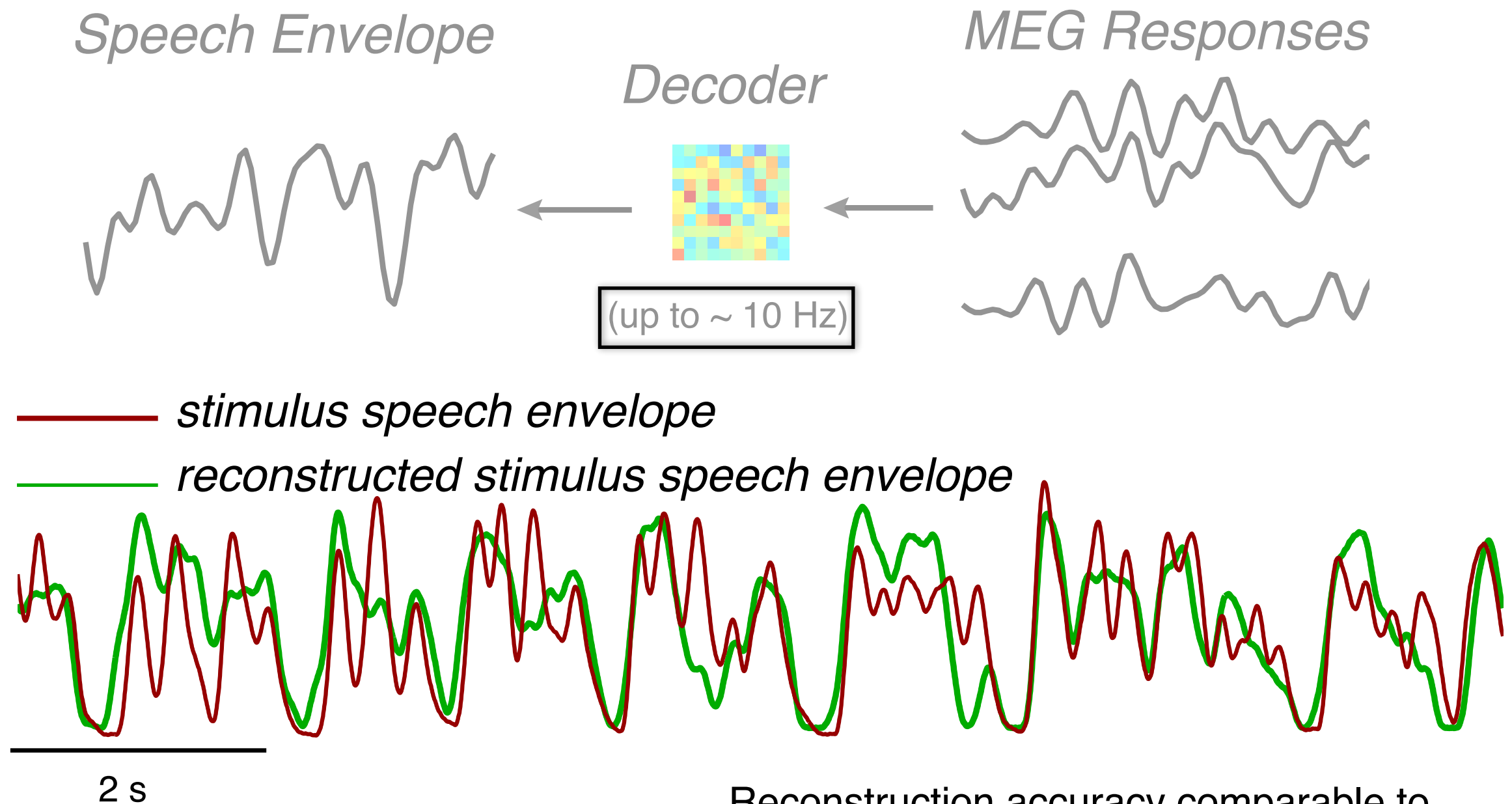


(up to ~ 10 Hz)

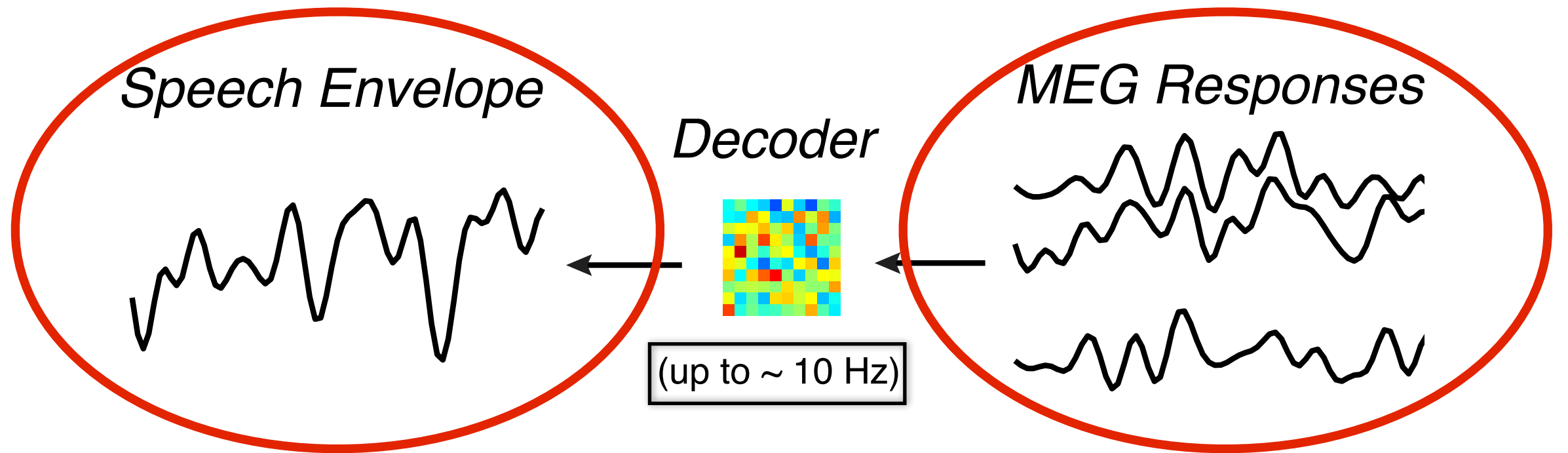
MEG Responses



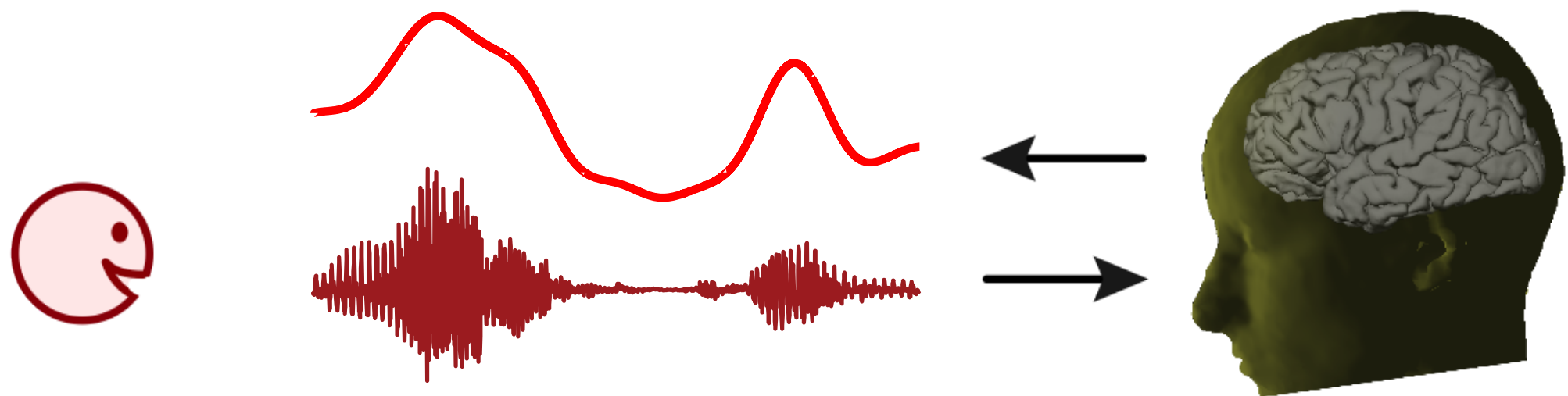
Neural Reconstruction of Speech Envelope



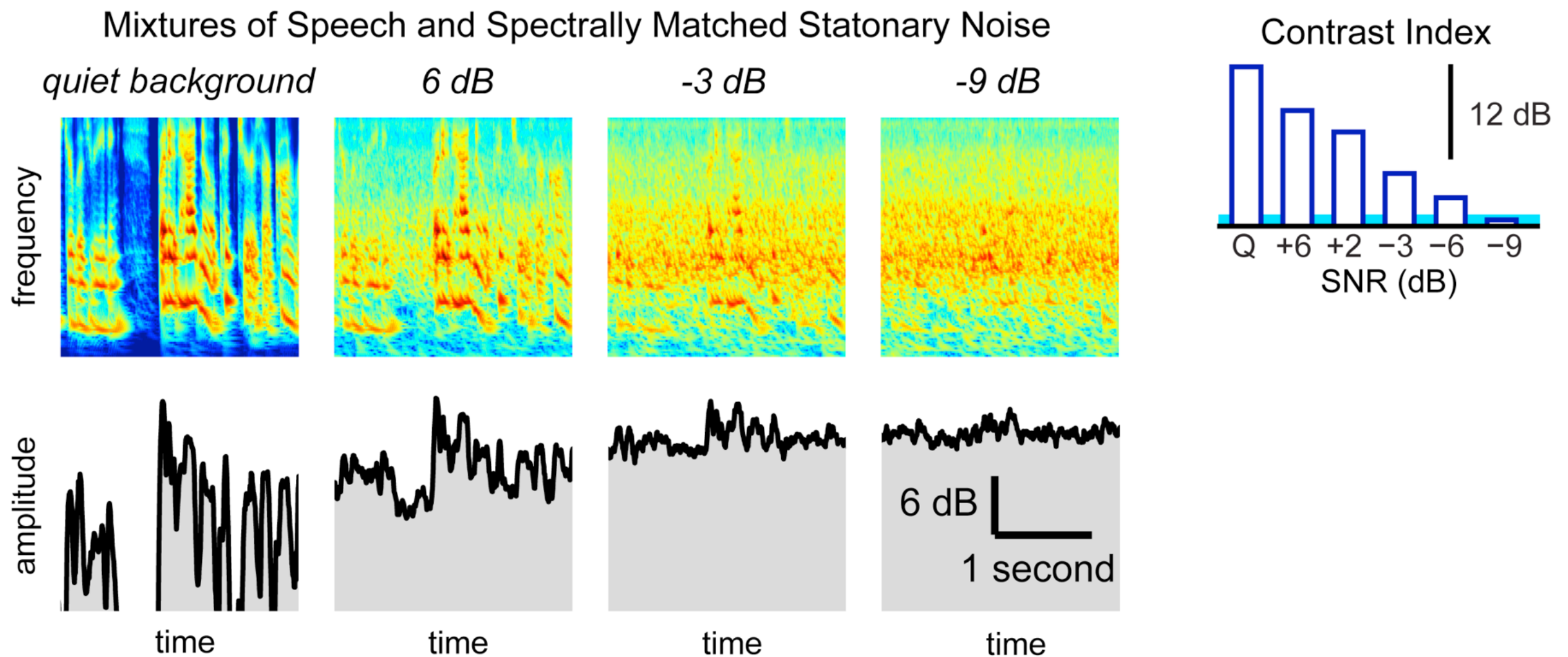
Reconstruction accuracy comparable to
single unit & ECoG recordings



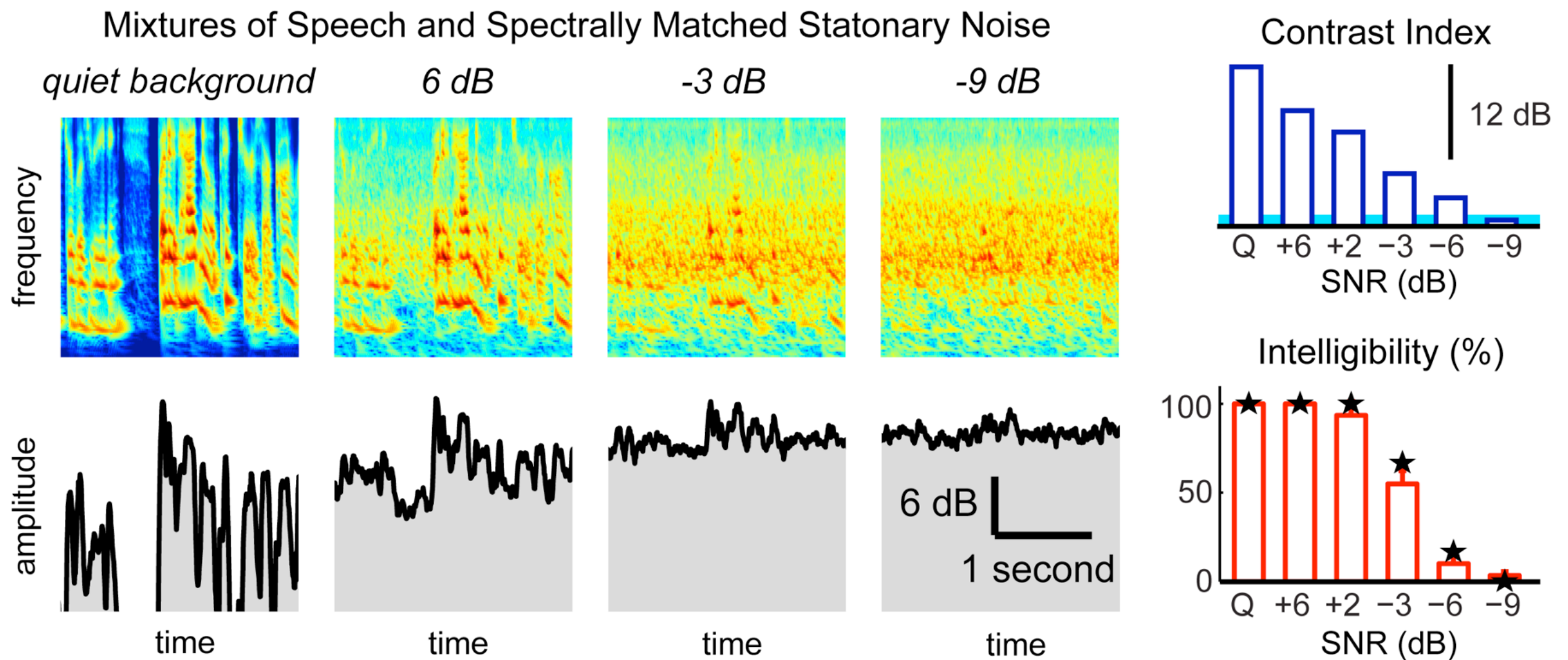
Neural Representation of Speech: Temporal



Speech in Noise

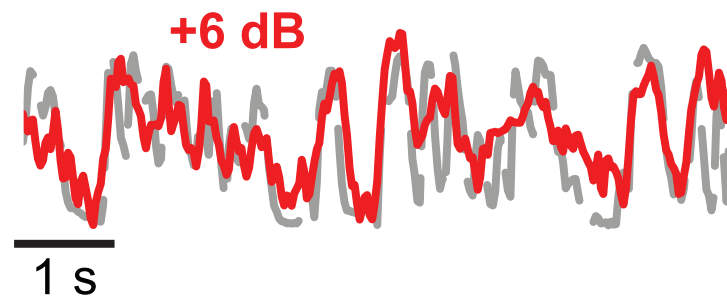


Speech in Noise



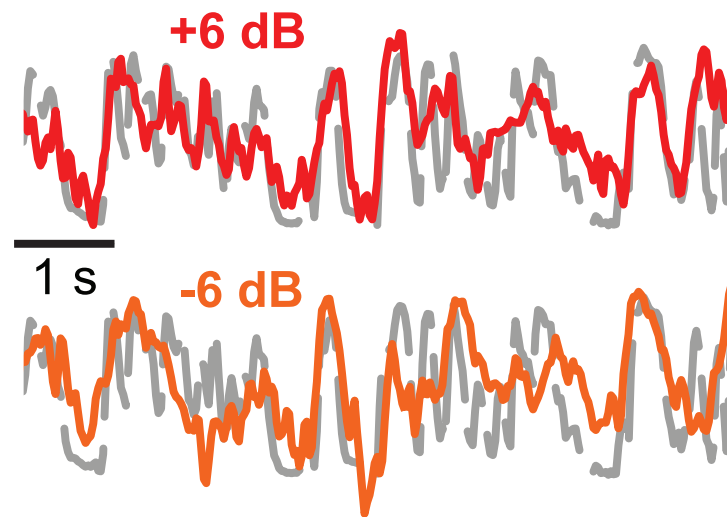
Speech in Noise: Results

Neural Reconstruction of
Underlying Speech Envelope



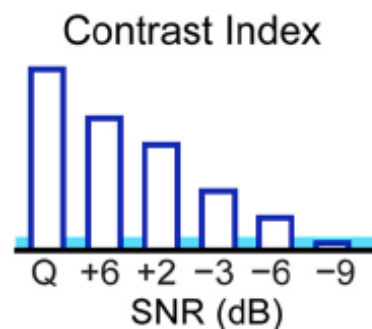
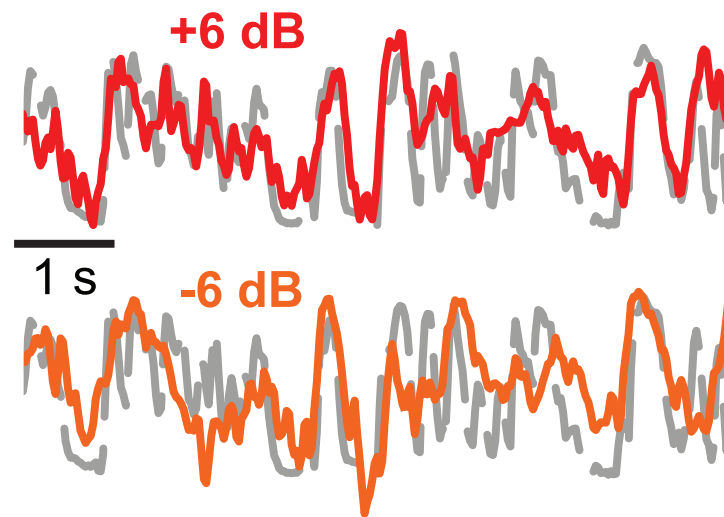
Speech in Noise: Results

Neural Reconstruction of
Underlying Speech Envelope



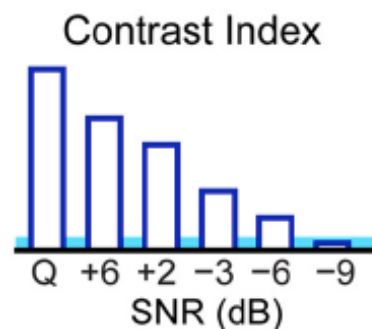
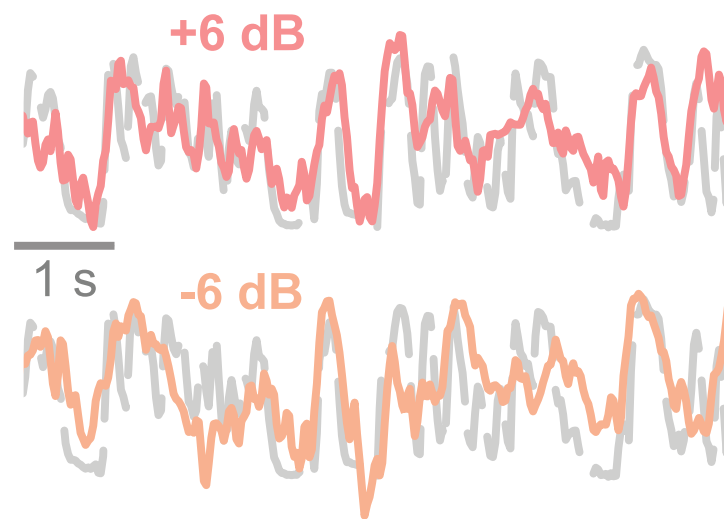
Speech in Noise: Results

Neural Reconstruction of
Underlying Speech Envelope

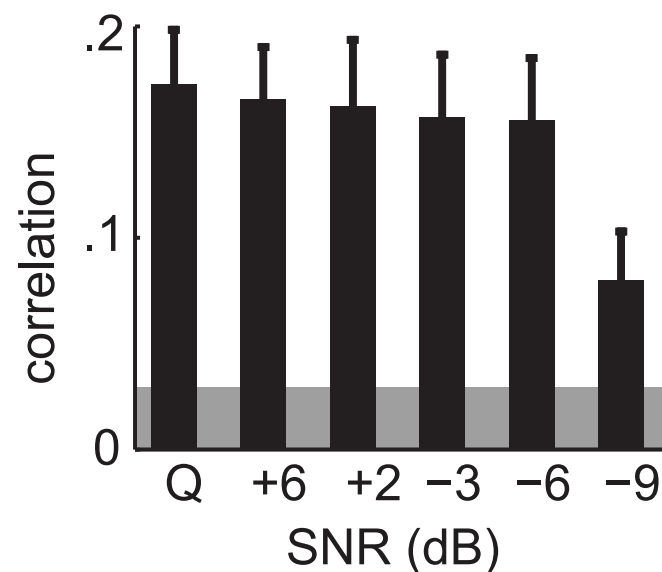


Speech in Noise: Results

Neural Reconstruction of
Underlying Speech Envelope

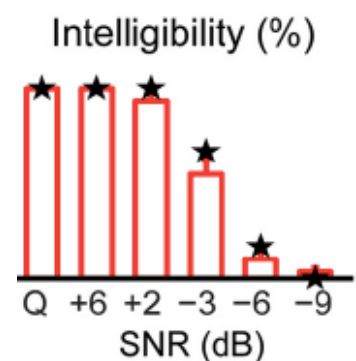
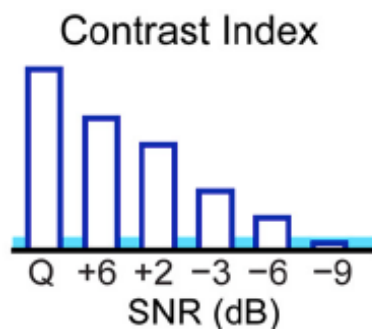
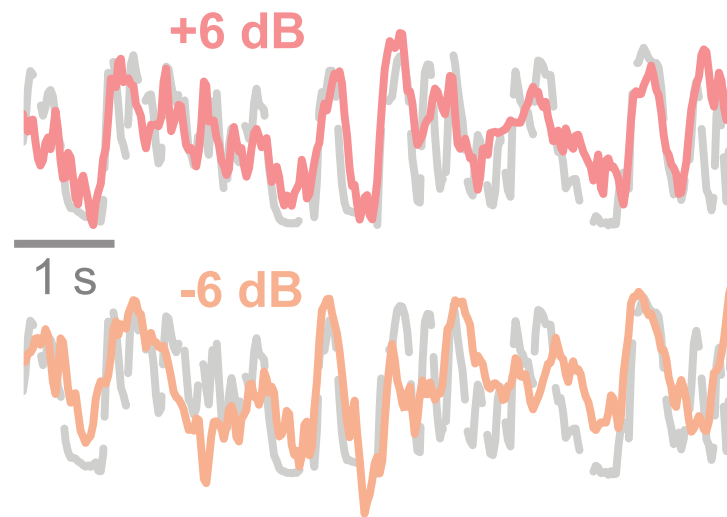


Reconstruction Accuracy

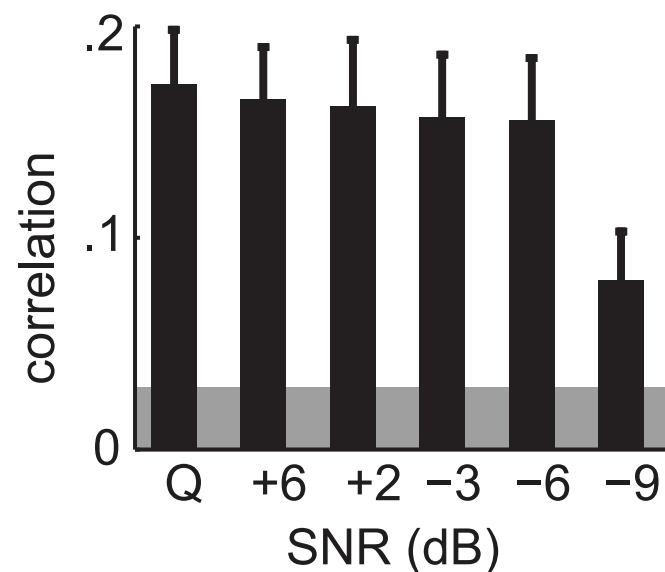


Speech in Noise: Results

Neural Reconstruction of Underlying Speech Envelope

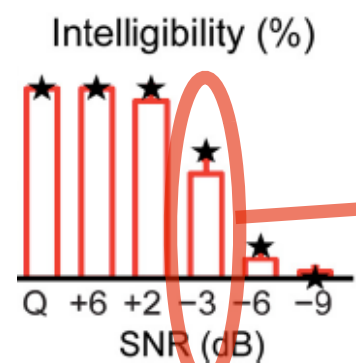
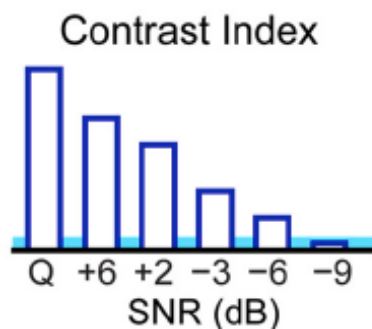
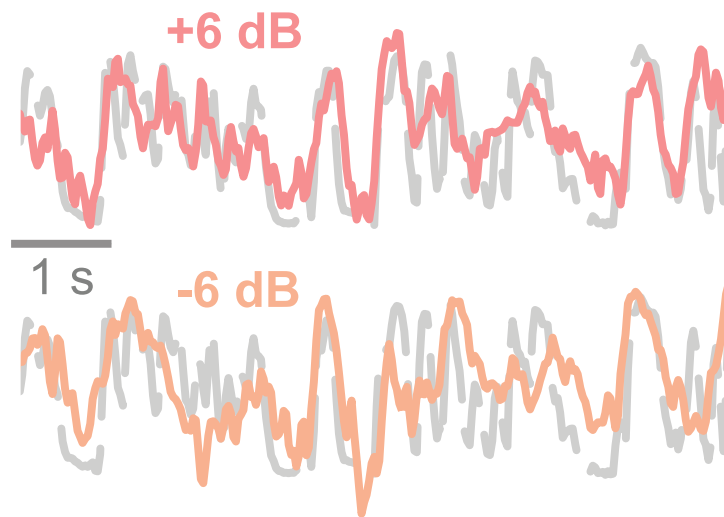


Reconstruction Accuracy

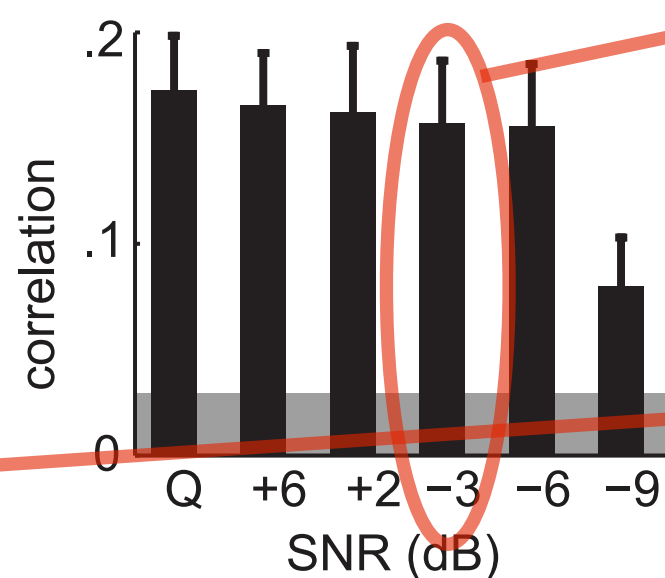


Speech in Noise: Results

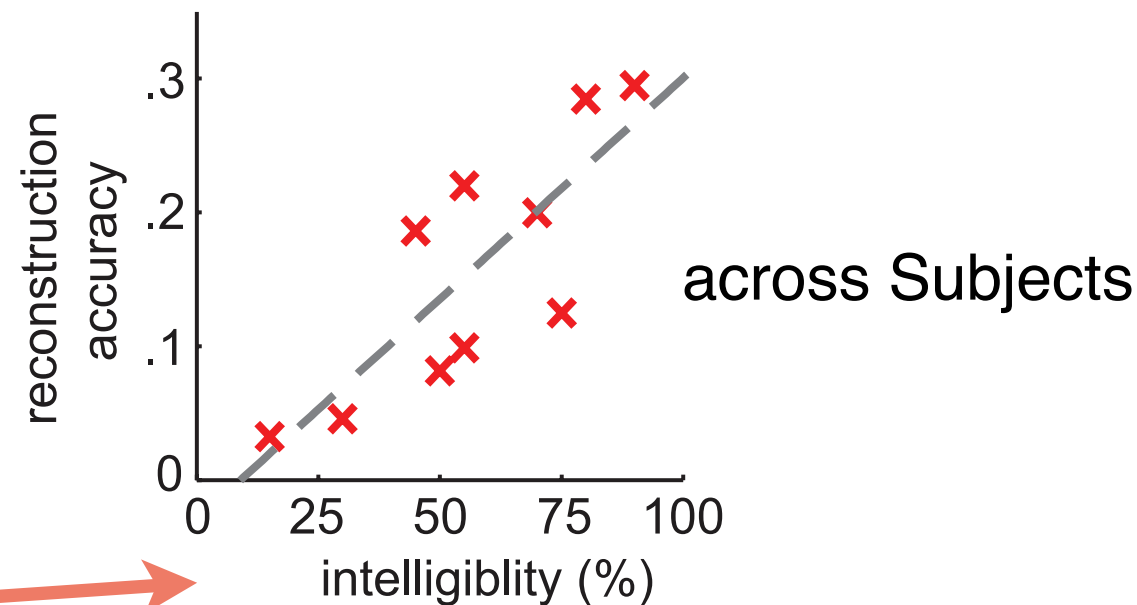
Neural Reconstruction of Underlying Speech Envelope



Reconstruction Accuracy



Correlation with Intelligibility



Ding & Simon, J Neuroscience (2013)

Cortical Speech Representations

- Neural Representations: Encoding & Decoding
- Linear models: Useful & Robust
- Speech **Envelope** only (as seen by MEG)
- Envelope Rates: $\sim 1 - 10$ Hz

The Cocktail Party



Alex Katz,
The Cocktail Party

The Cocktail Party



Alex Katz,
The Cocktail Party

The Cocktail Party



Alex Katz,
The Cocktail Party

The Cocktail Party



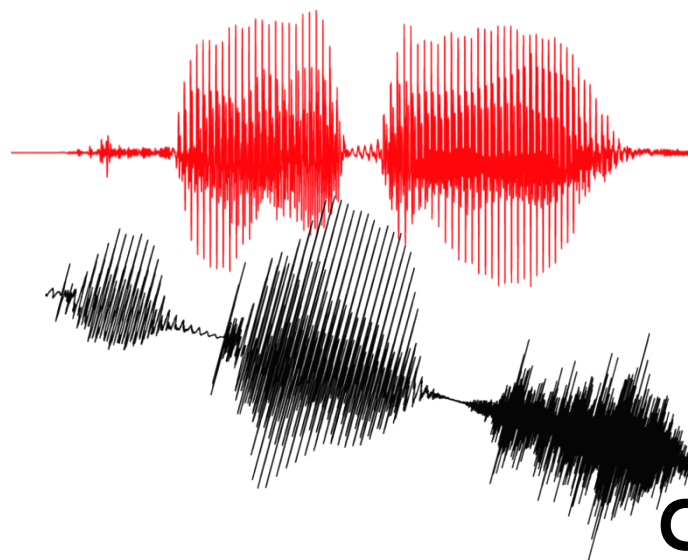
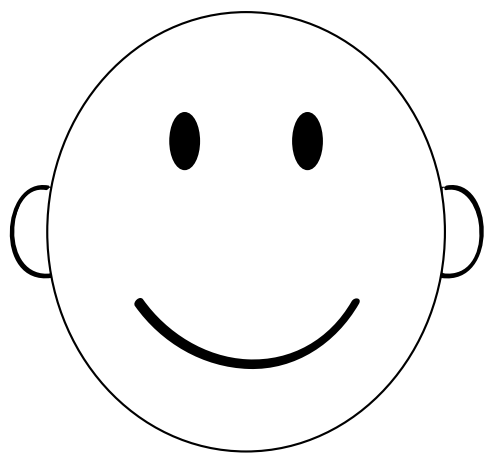
Alex Katz,
The Cocktail Party

The Cocktail Party



Alex Katz,
The Cocktail Party

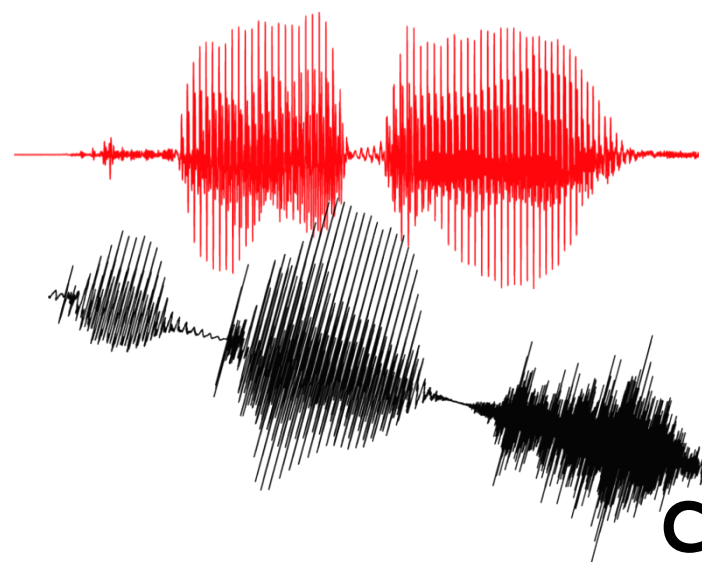
Experiments



speech

competing speech

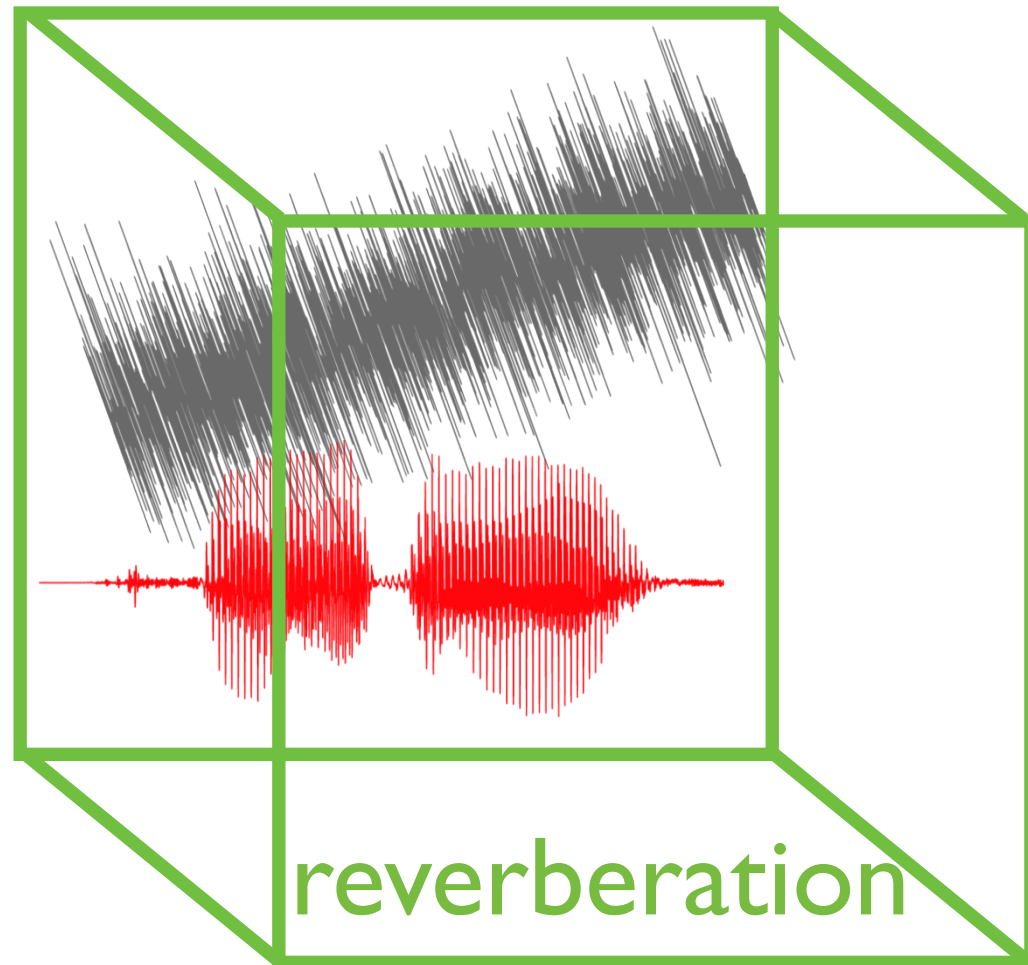
Experiments



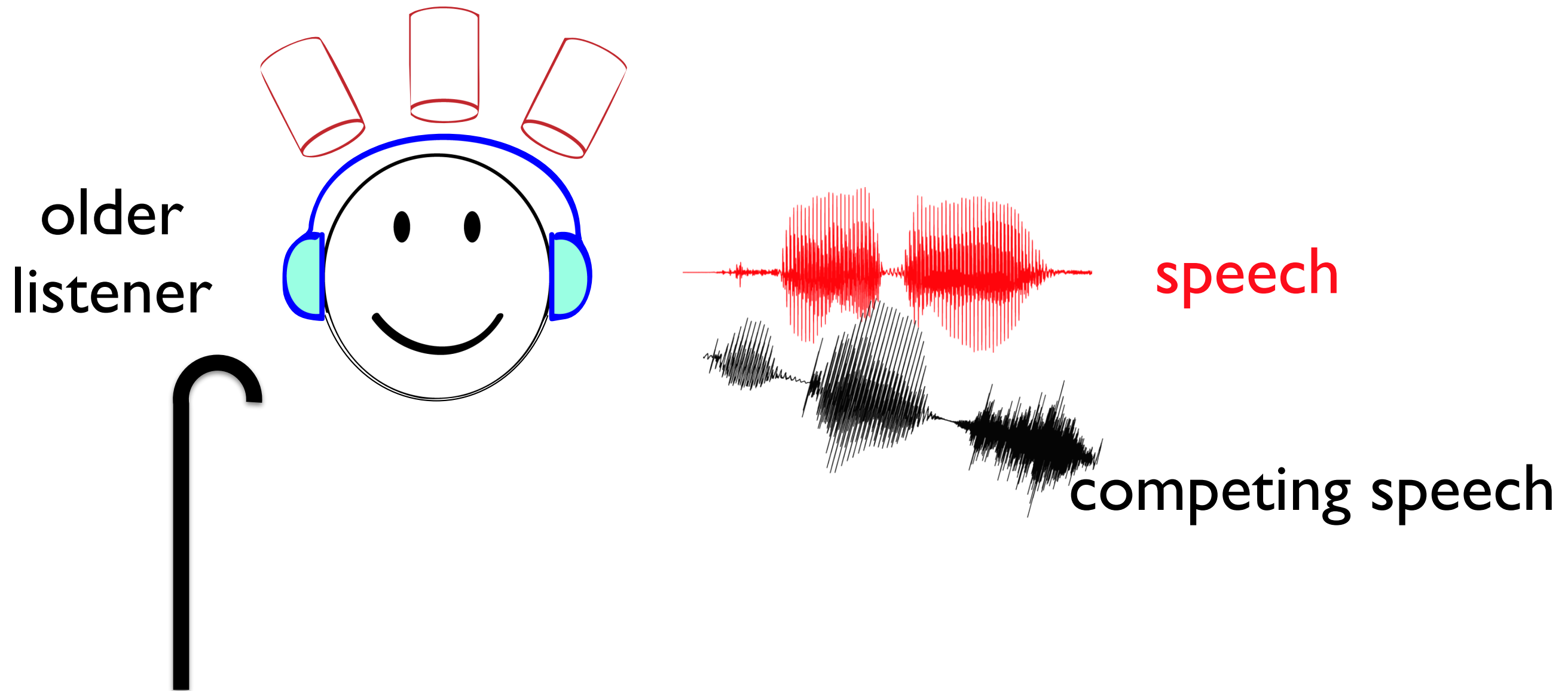
speech

competing speech

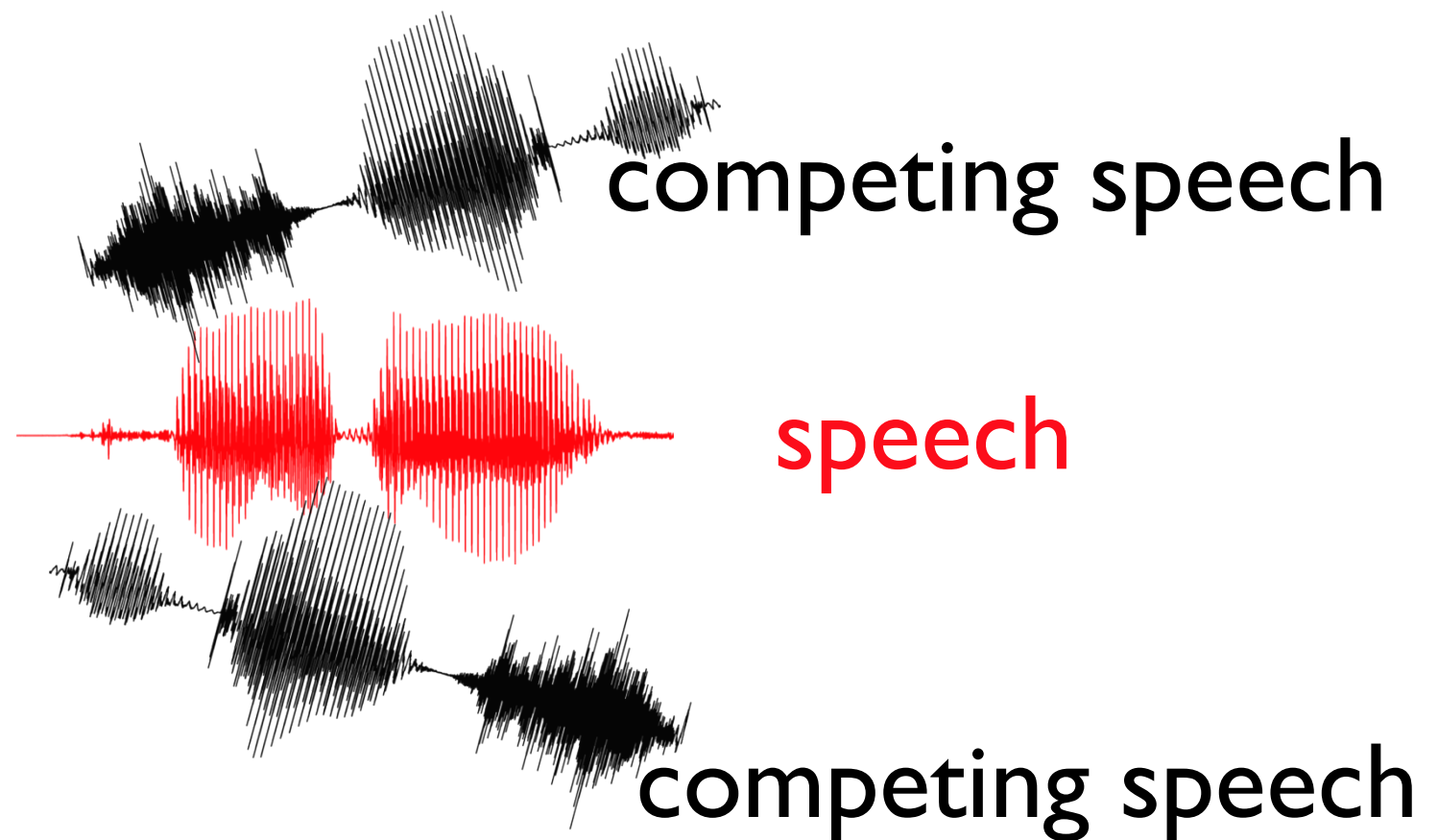
Experiments in Progress



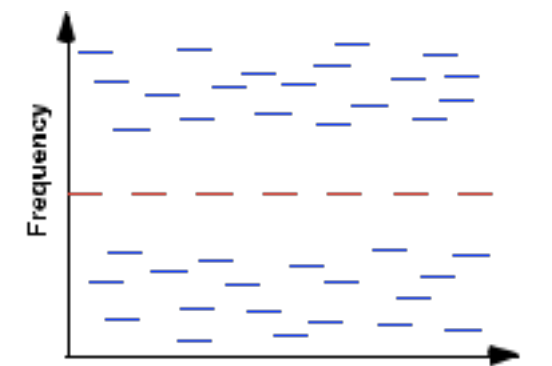
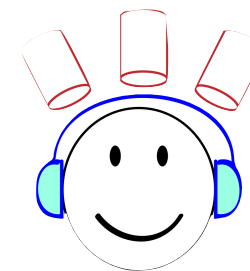
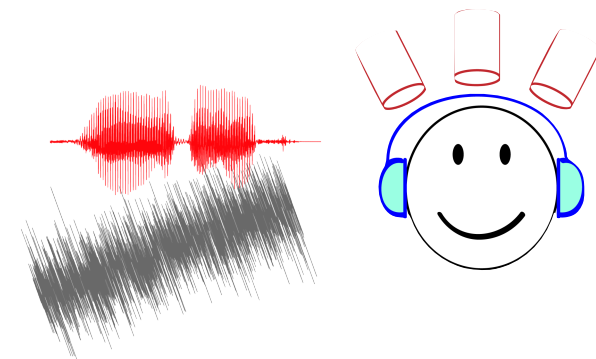
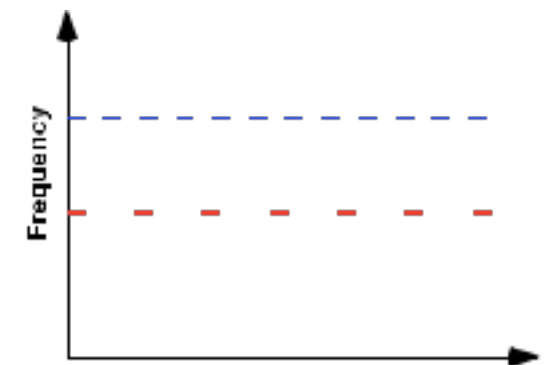
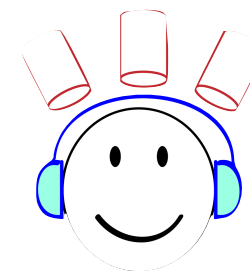
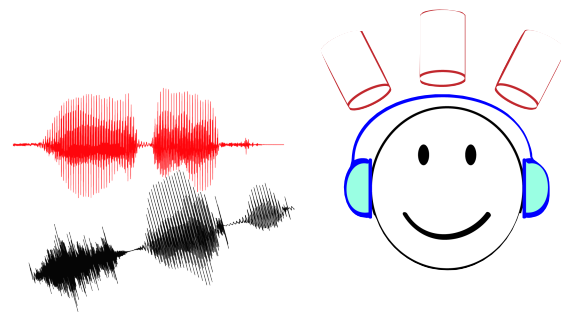
Experiments in Progress



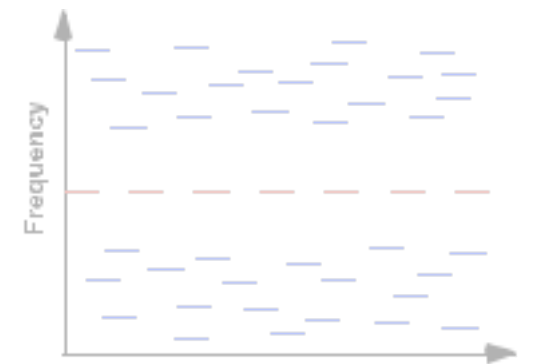
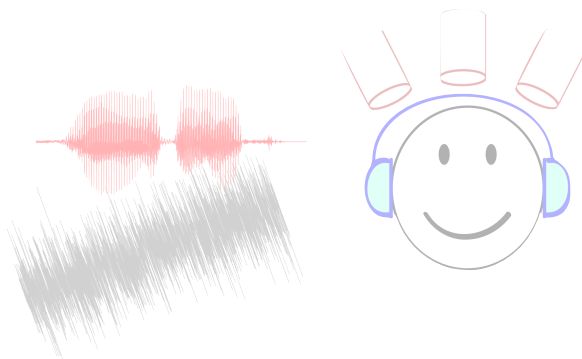
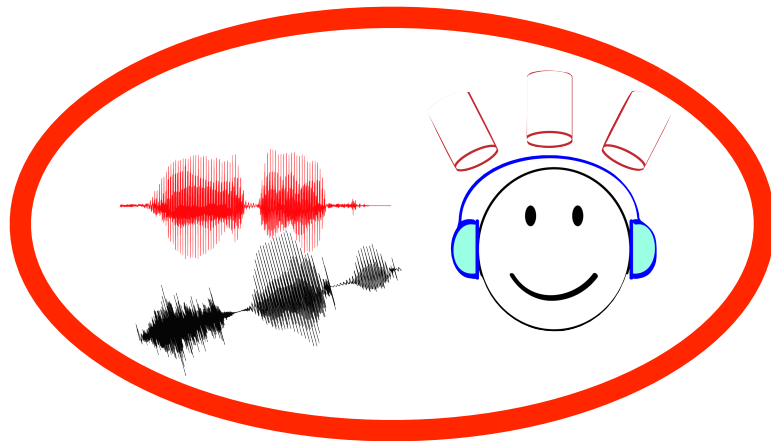
Experiments in Progress



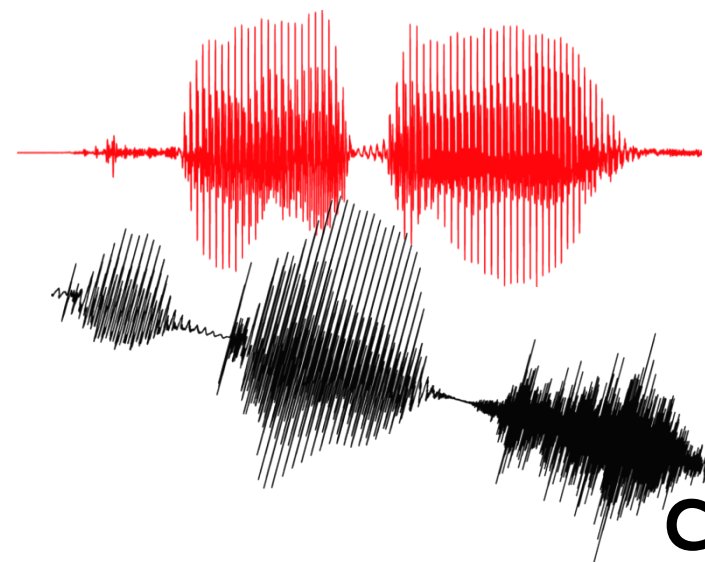
Experiments



Experiments



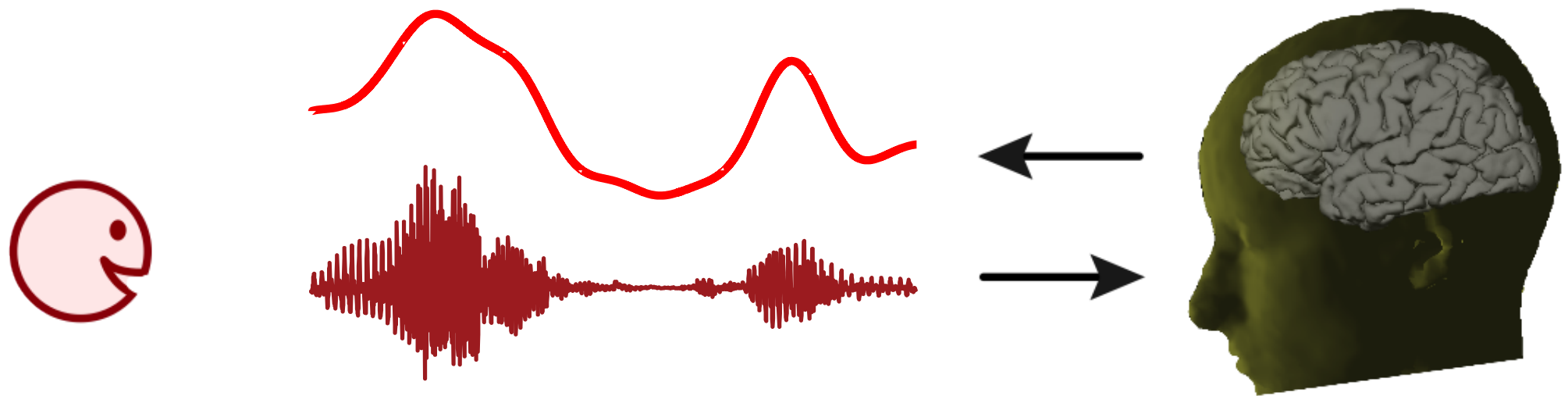
Two Competing Speakers



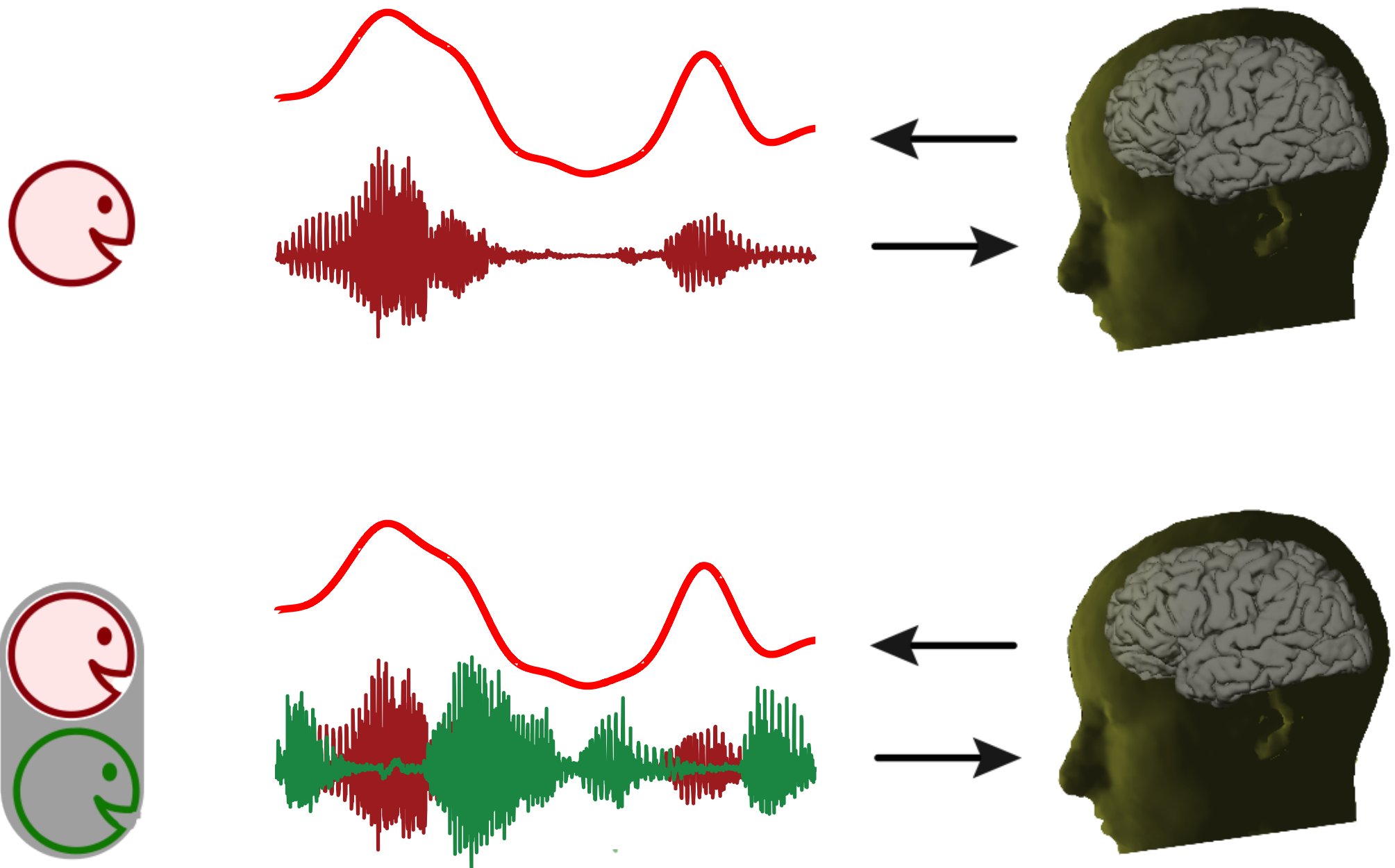
speech

competing speech

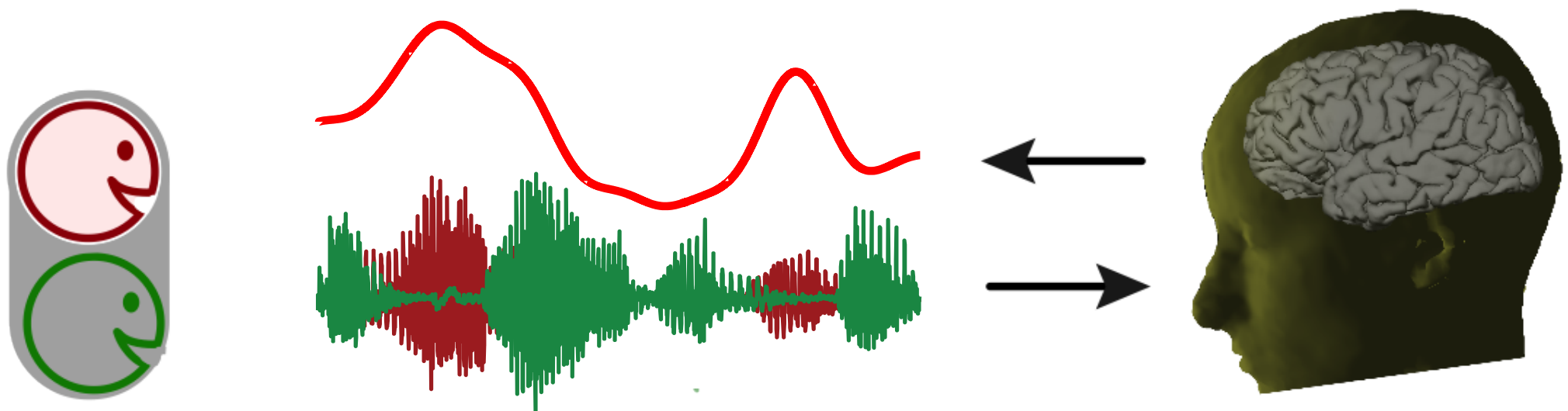
Selective Neural Encoding



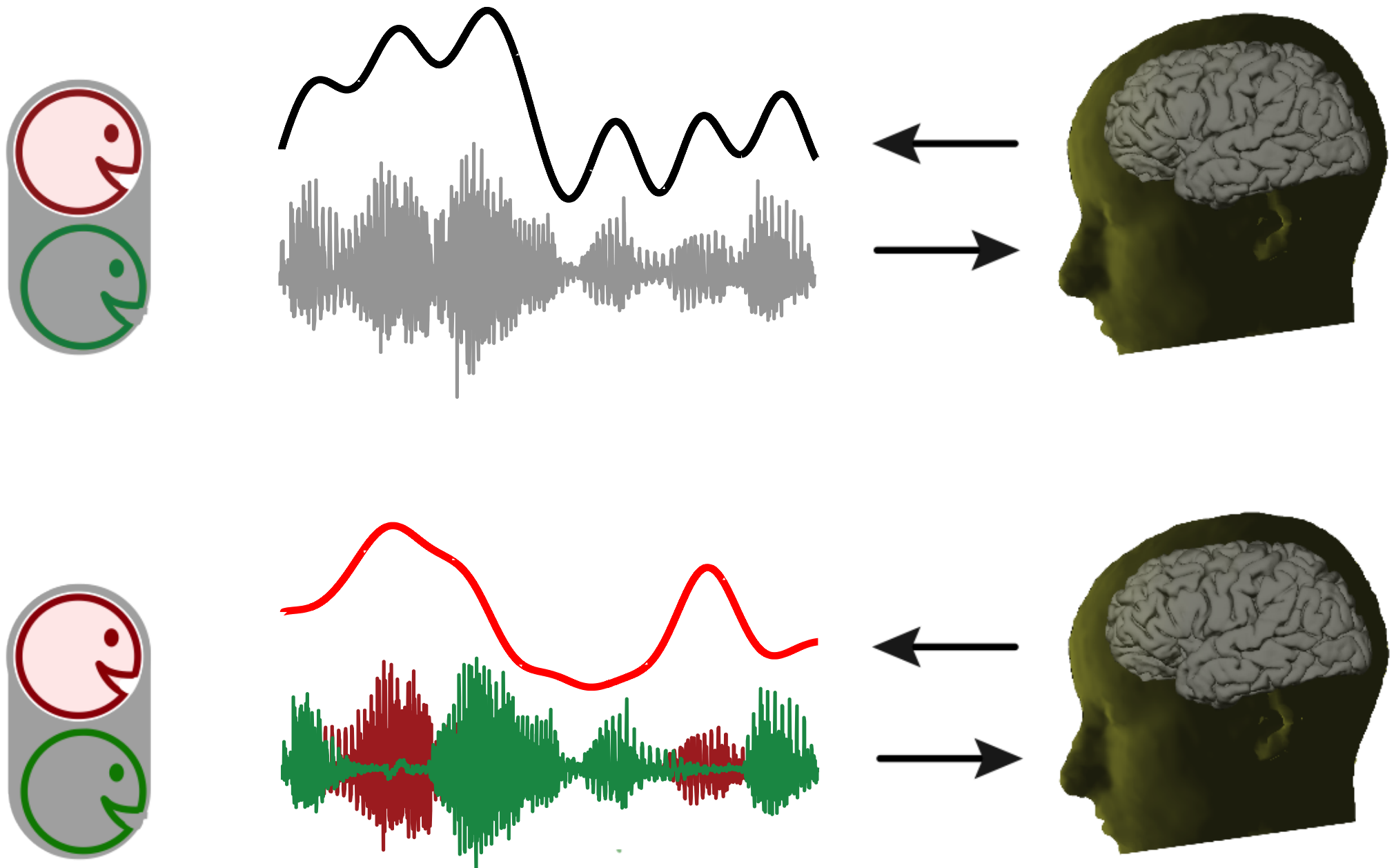
Selective Neural Encoding



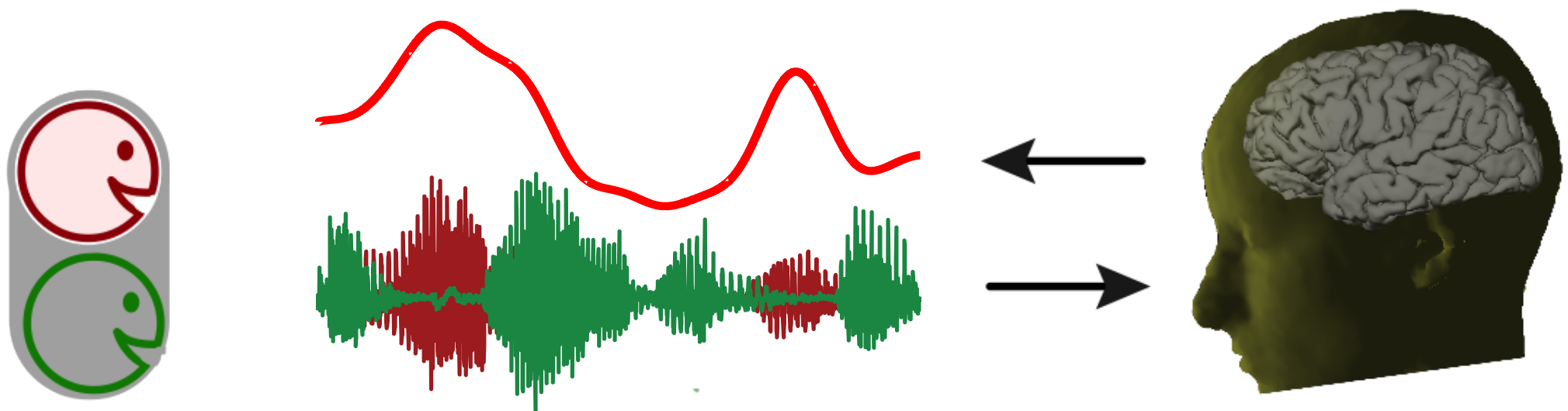
Selective Neural Encoding



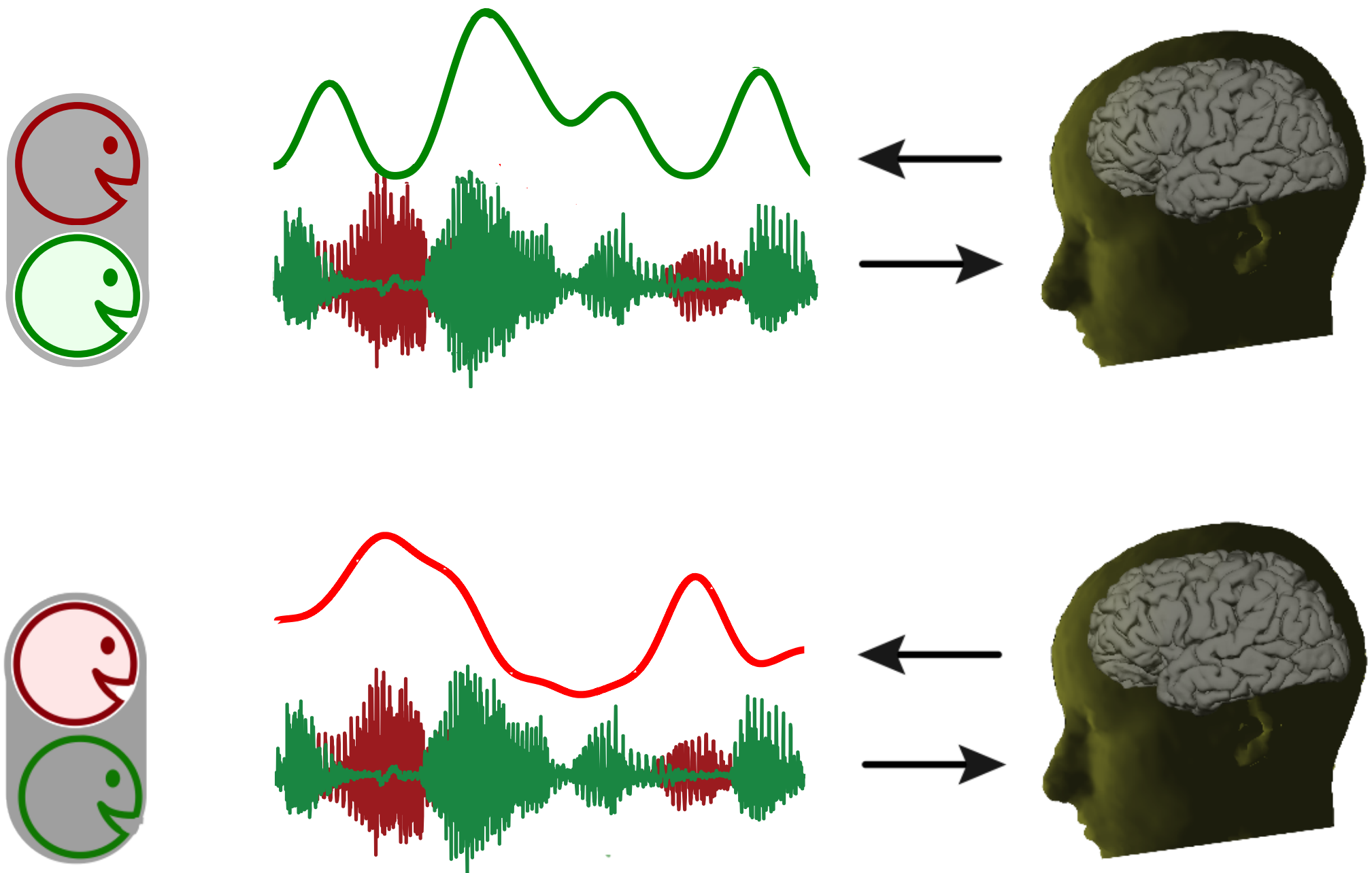
Unselective vs. Selective Neural Encoding



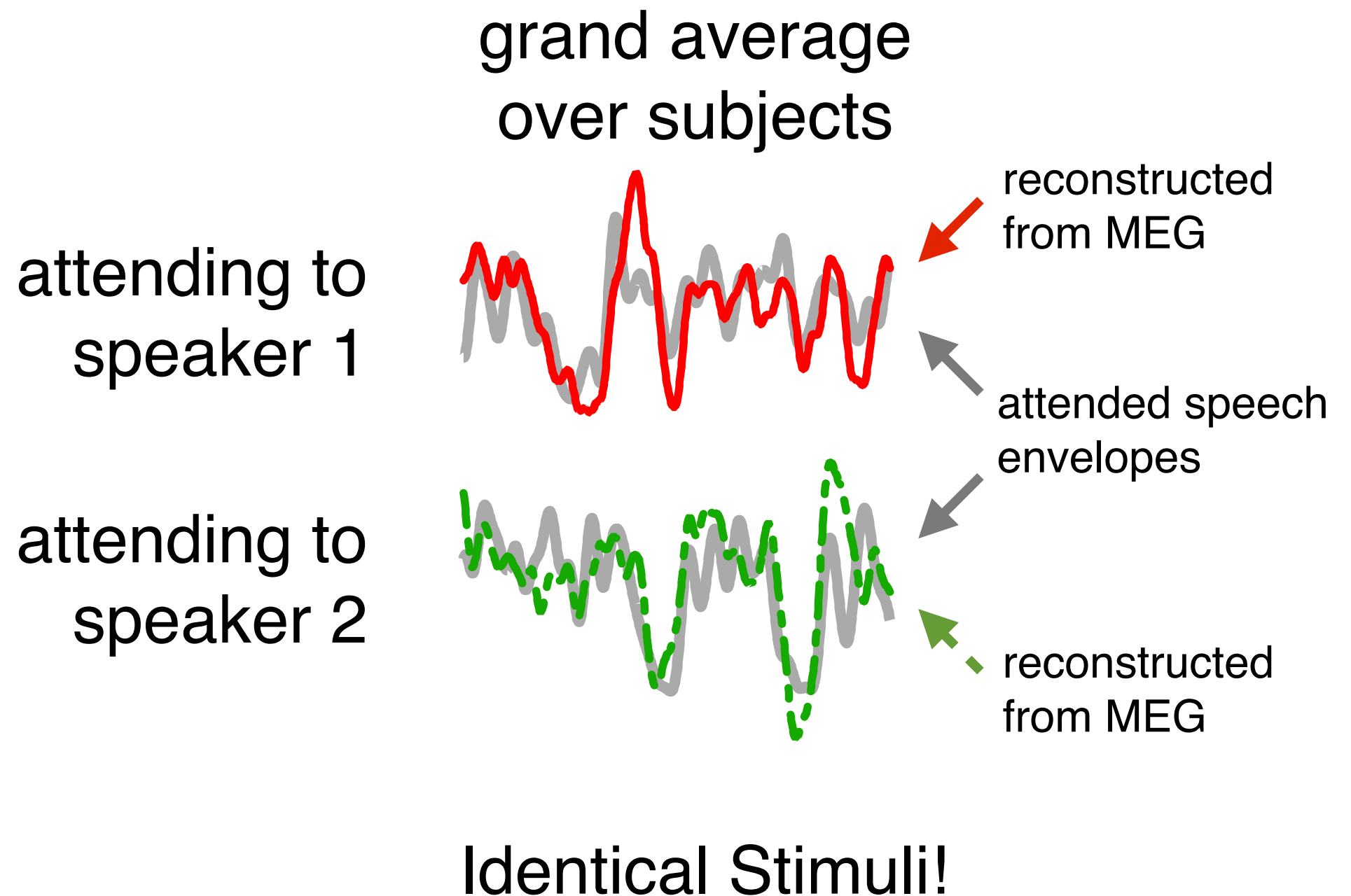
Unselective vs. Selective Neural Encoding



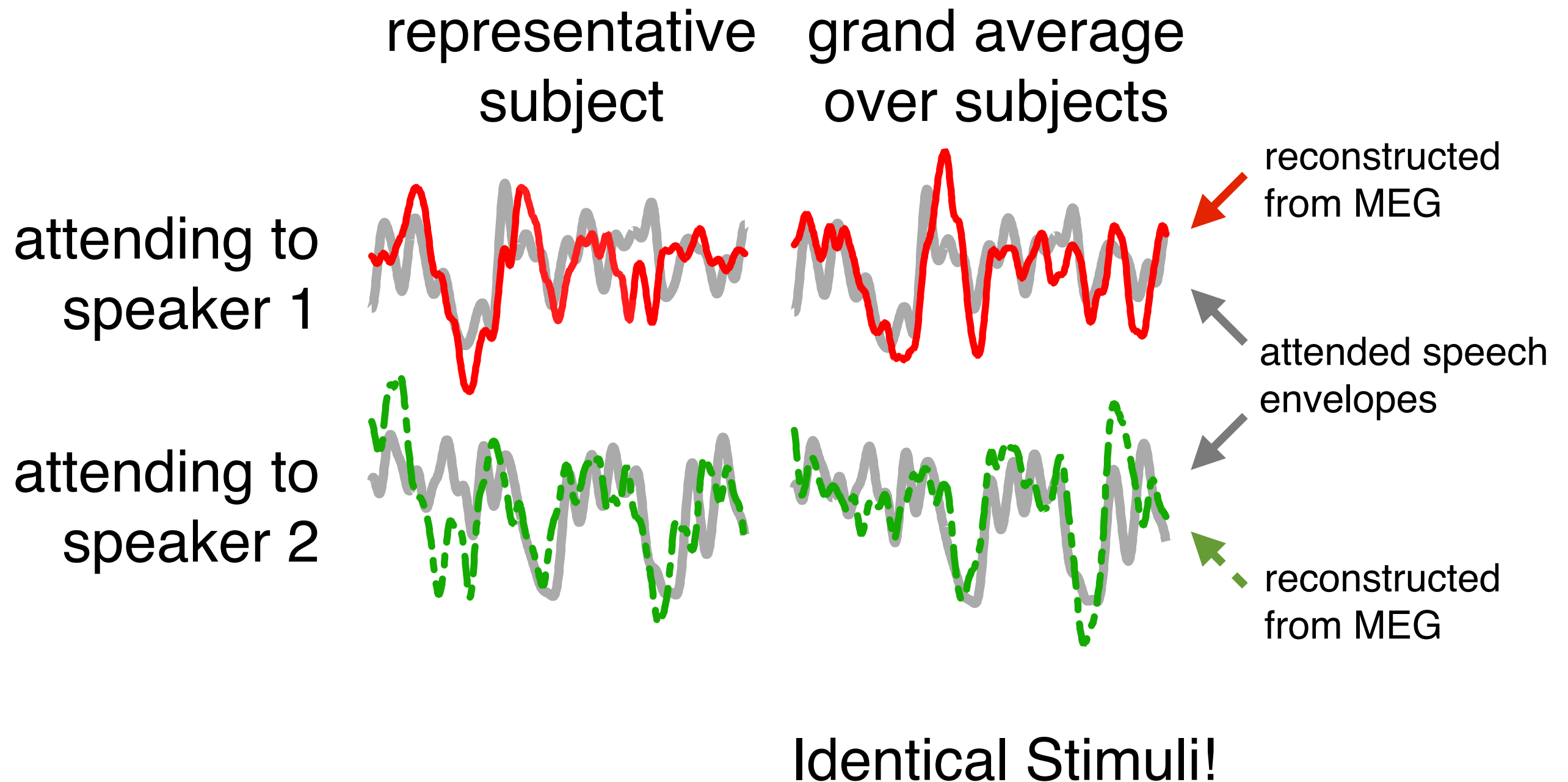
Selective Neural Encoding



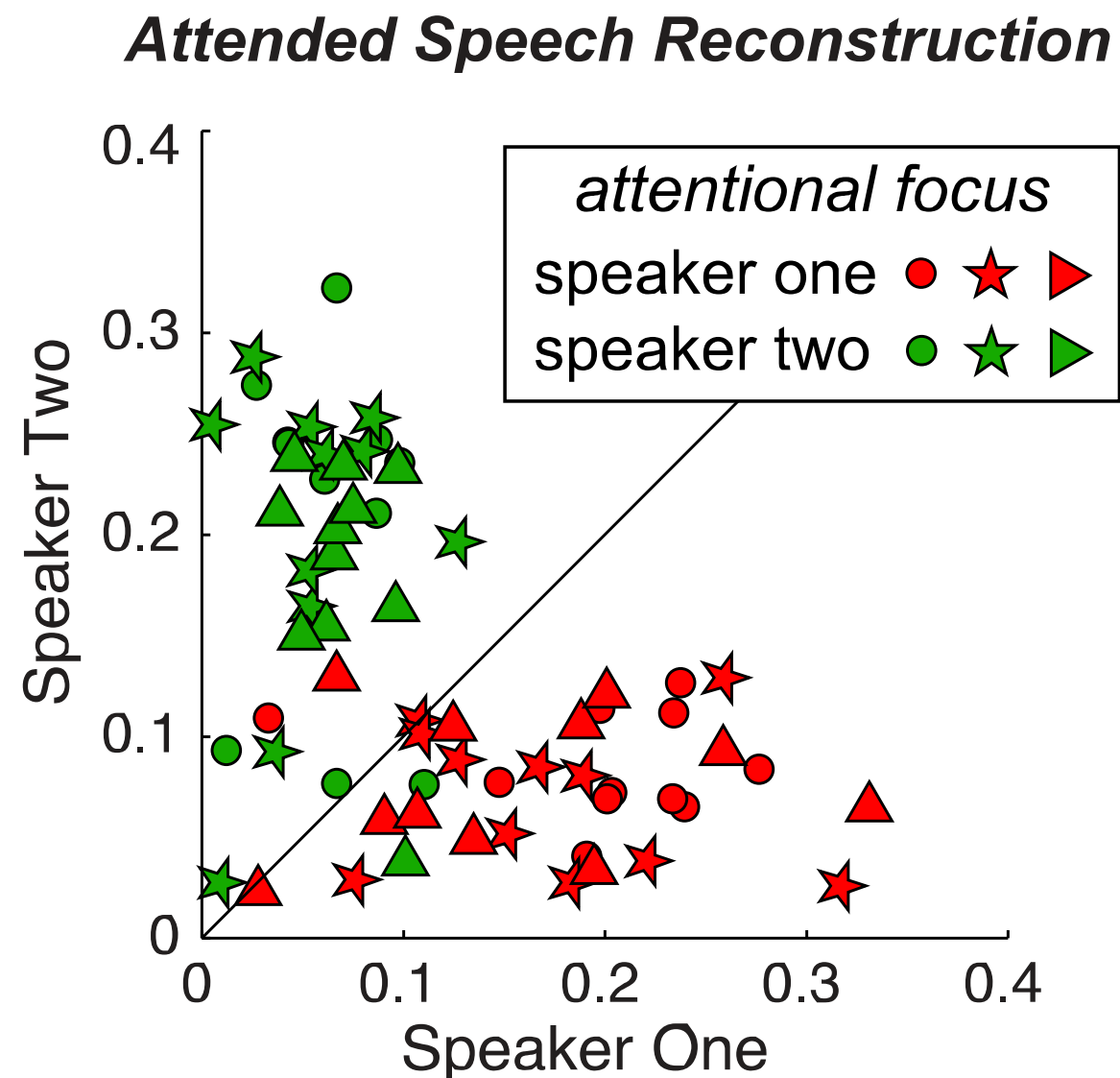
Stream-Specific Representation



Stream-Specific Representation

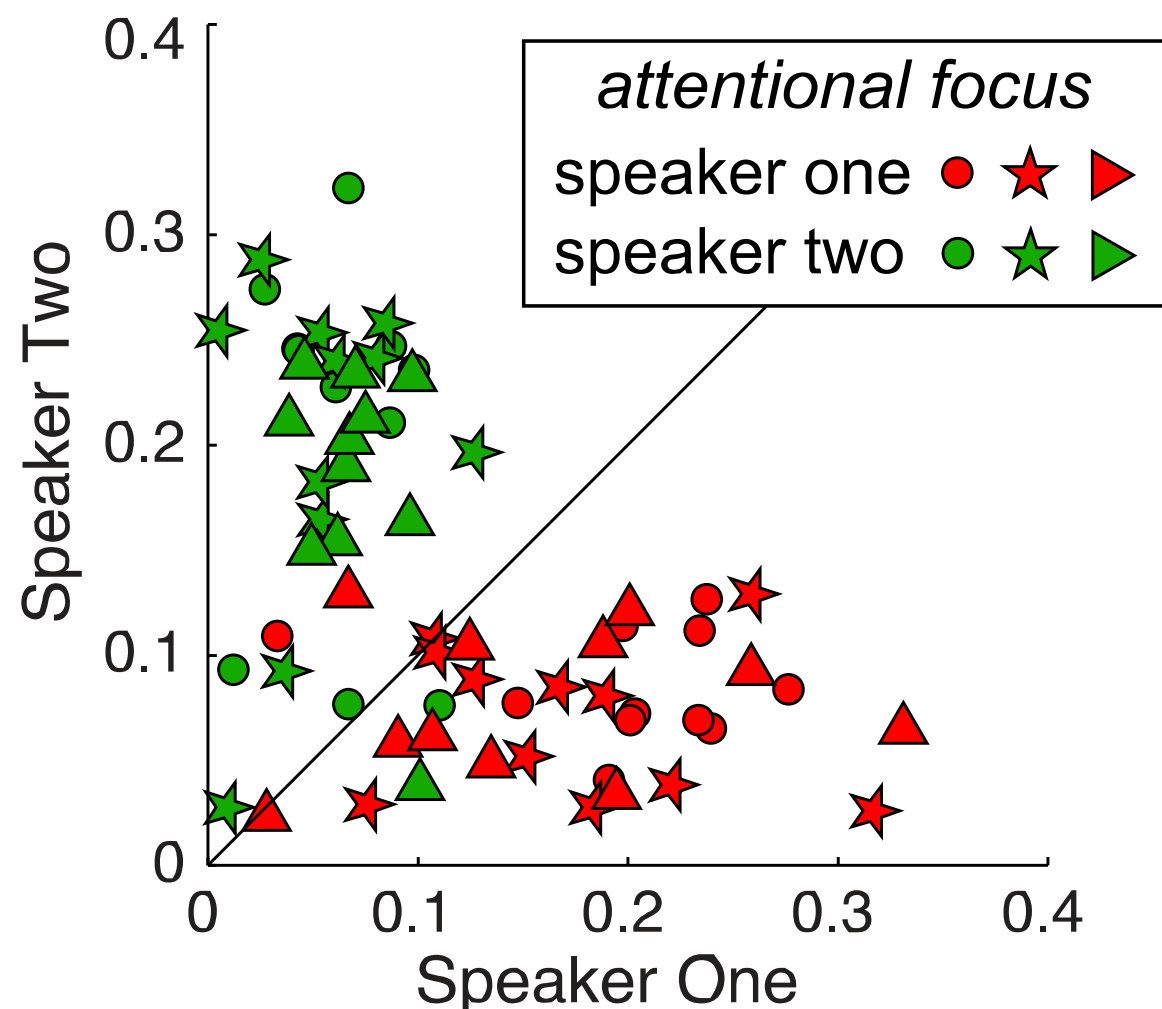


Single Trial Speech Reconstruction

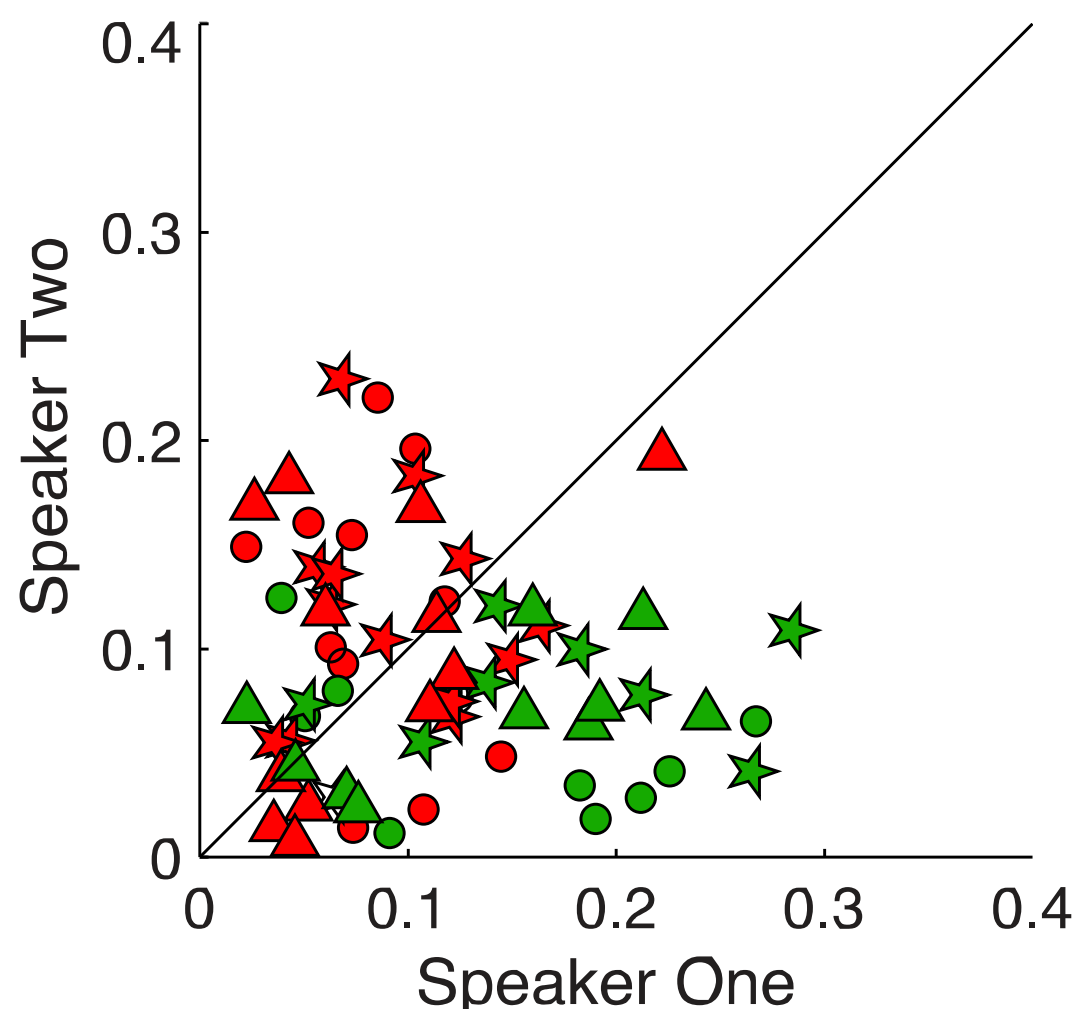


Single Trial Speech Reconstruction

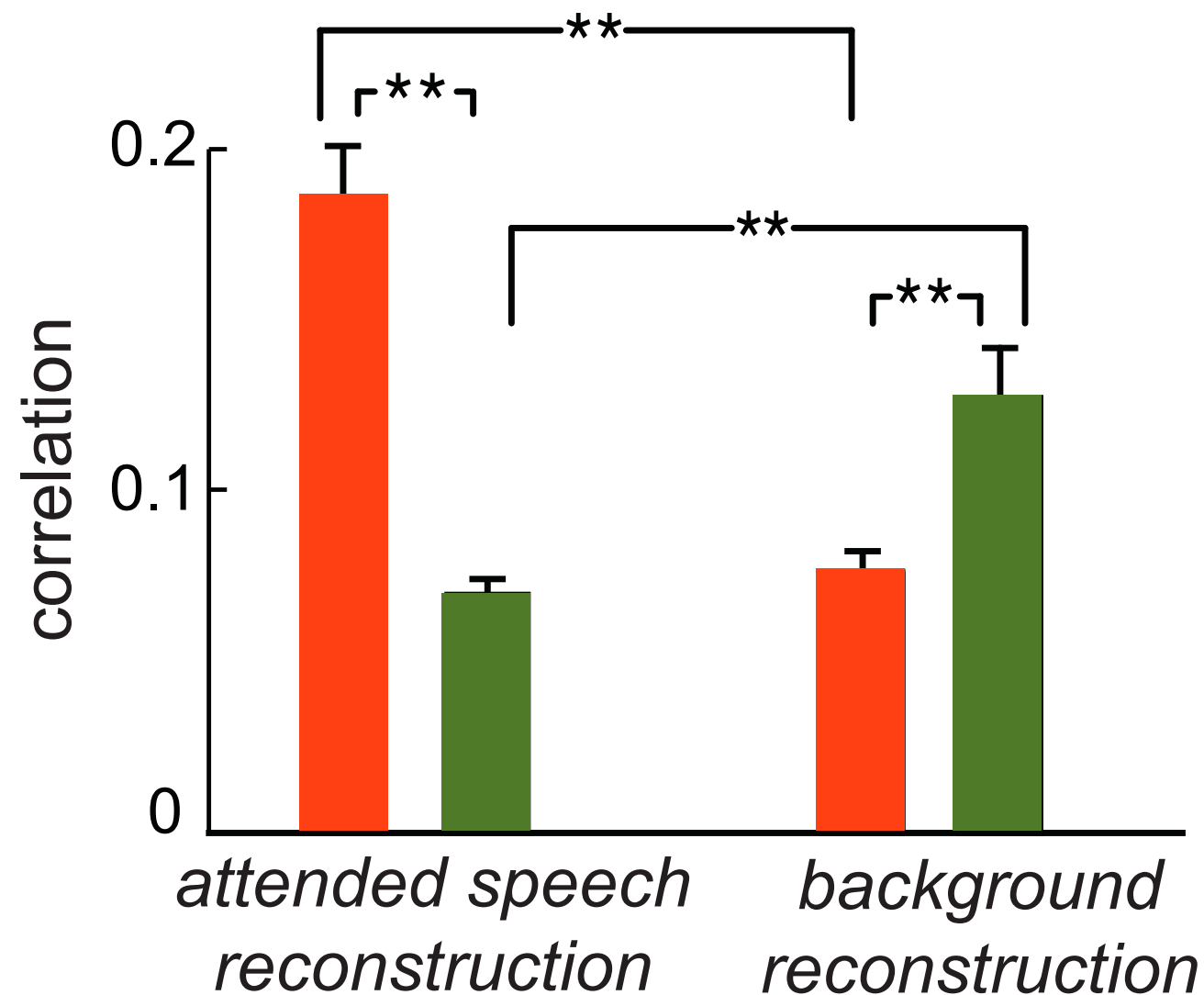
Attended Speech Reconstruction



Background Speech Reconstruction



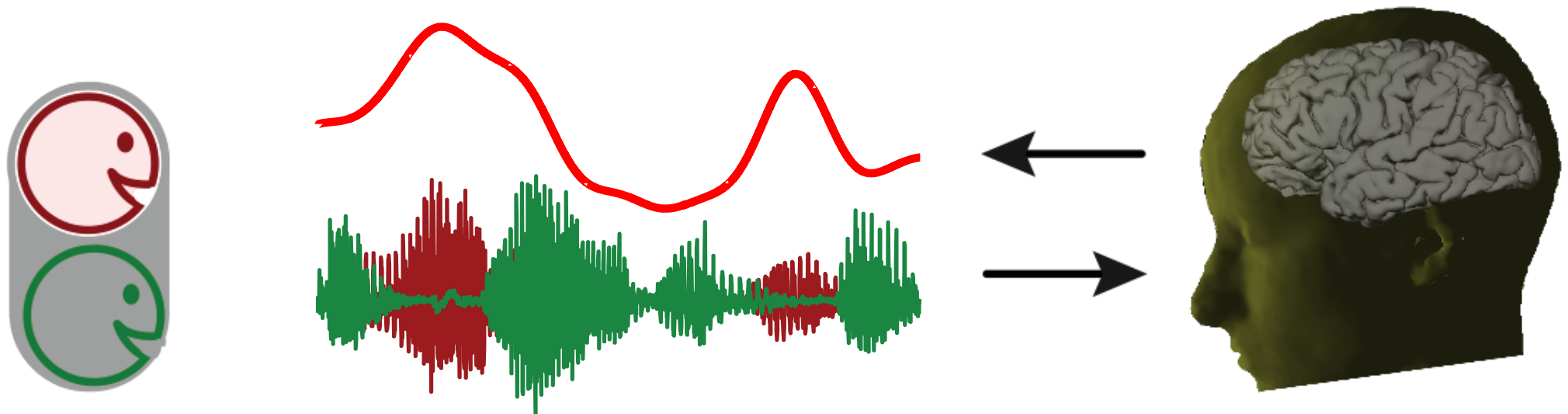
Overall Speech Reconstruction



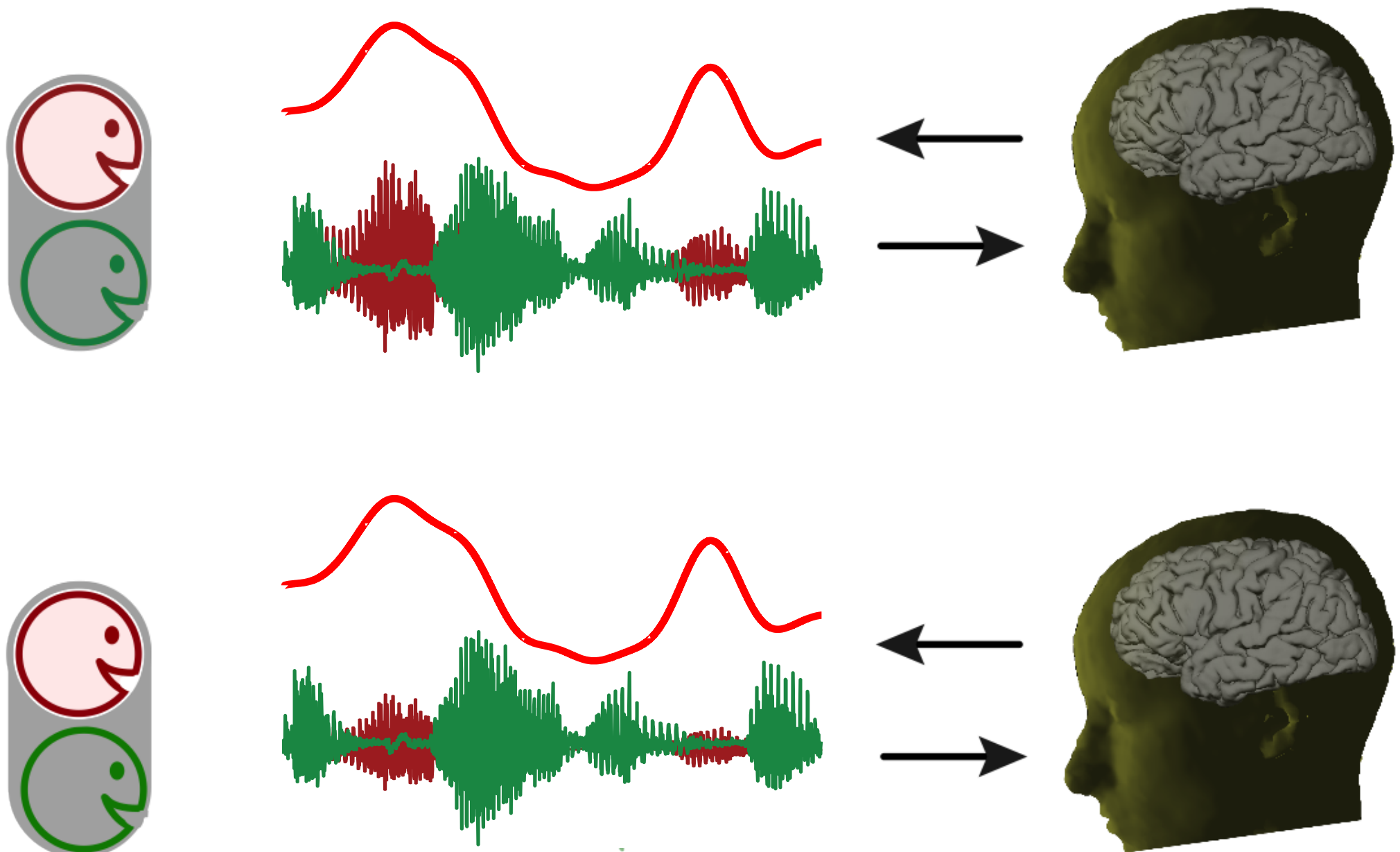
Distinct neural representations for different speech streams

attended speech ■ background ■

Invariance Under Relative Loudness Change?

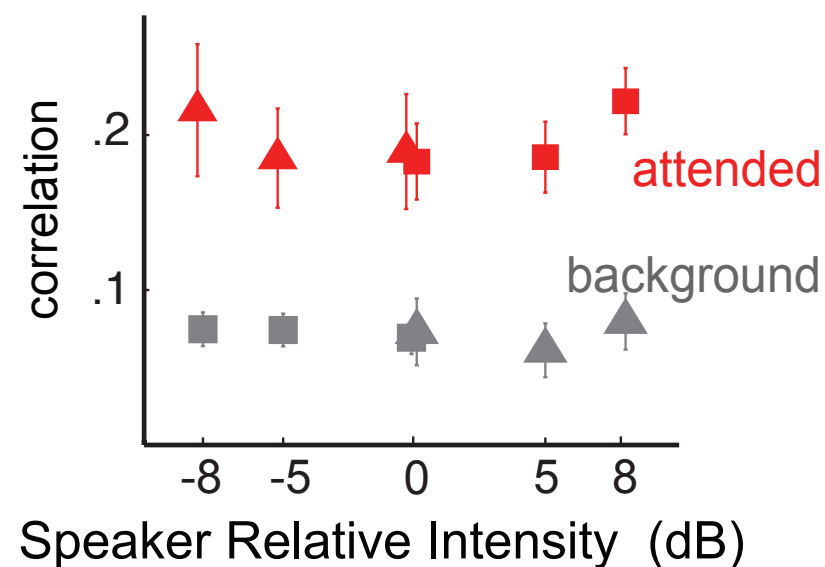


Invariance Under Relative Loudness Change?



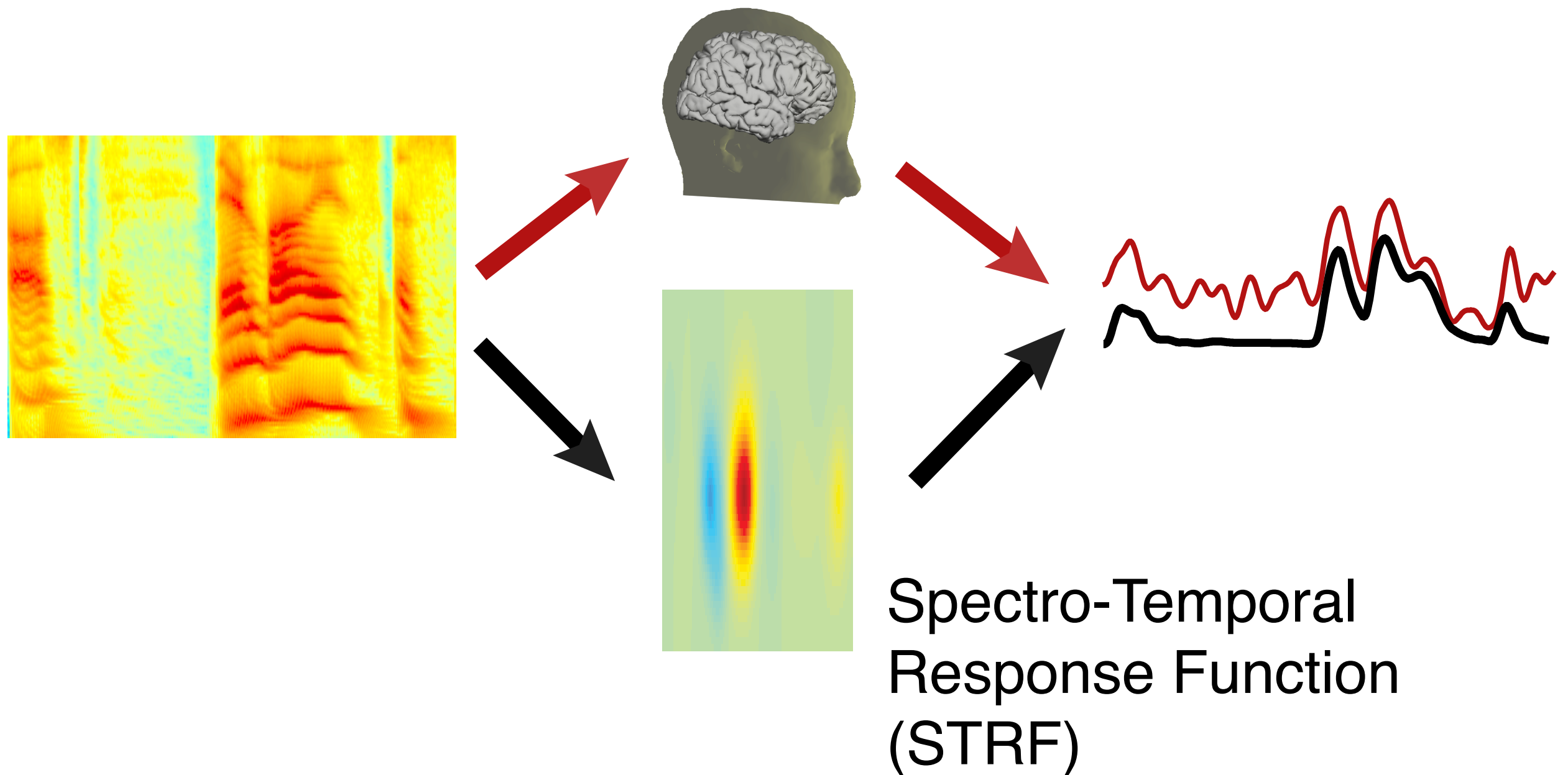
Invariance under Relative Loudness Change

Neural Results

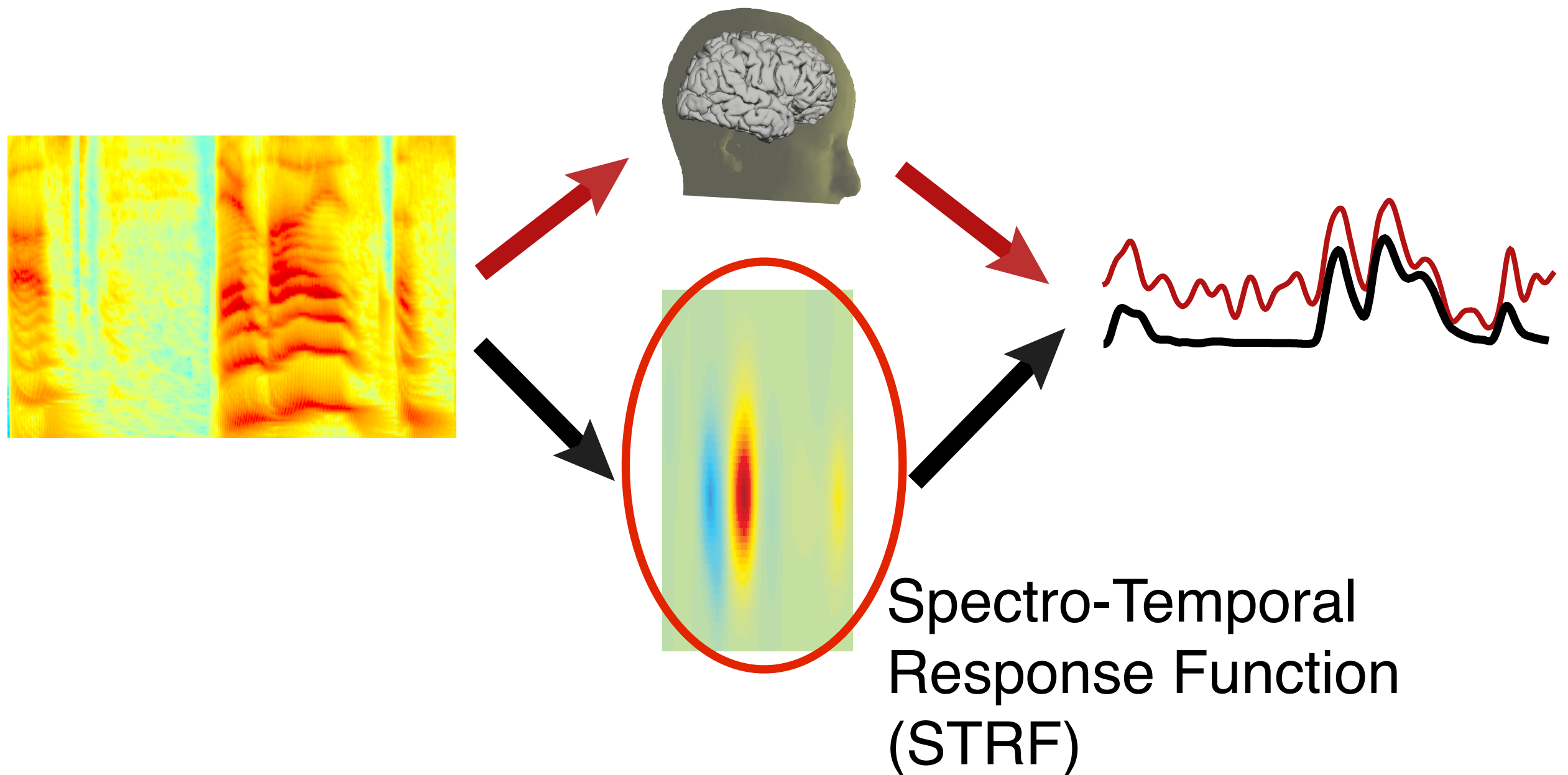


- Neural representation invariant to relative loudness change
- Stream-based Gain Control, not stimulus-based

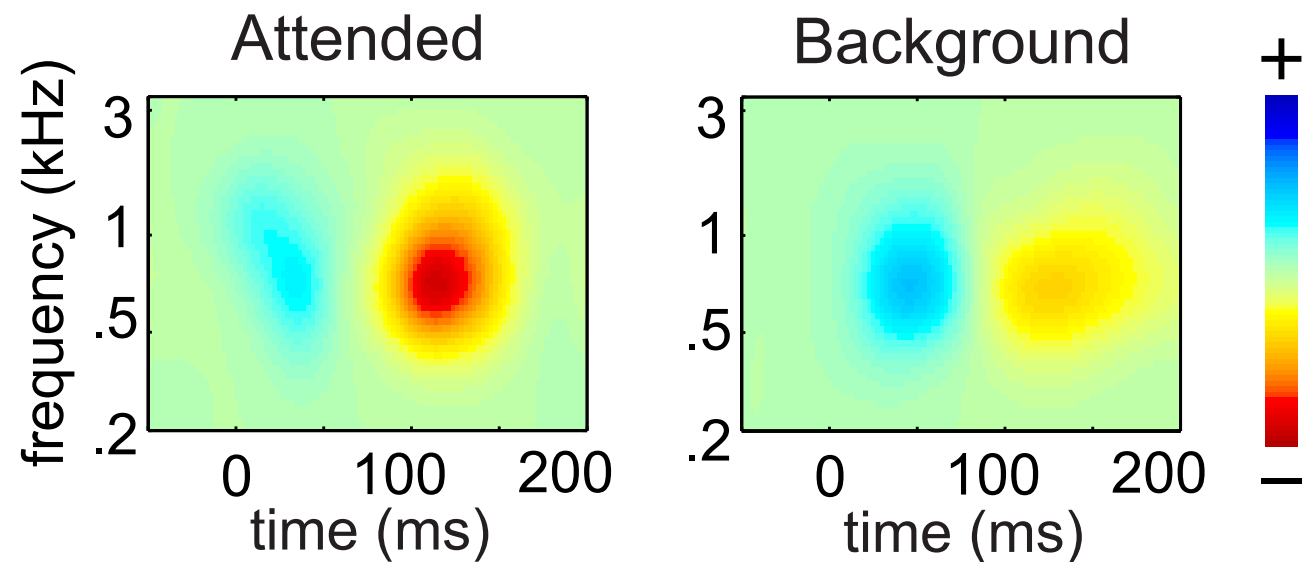
Forward STRF Model



Forward STRF Model

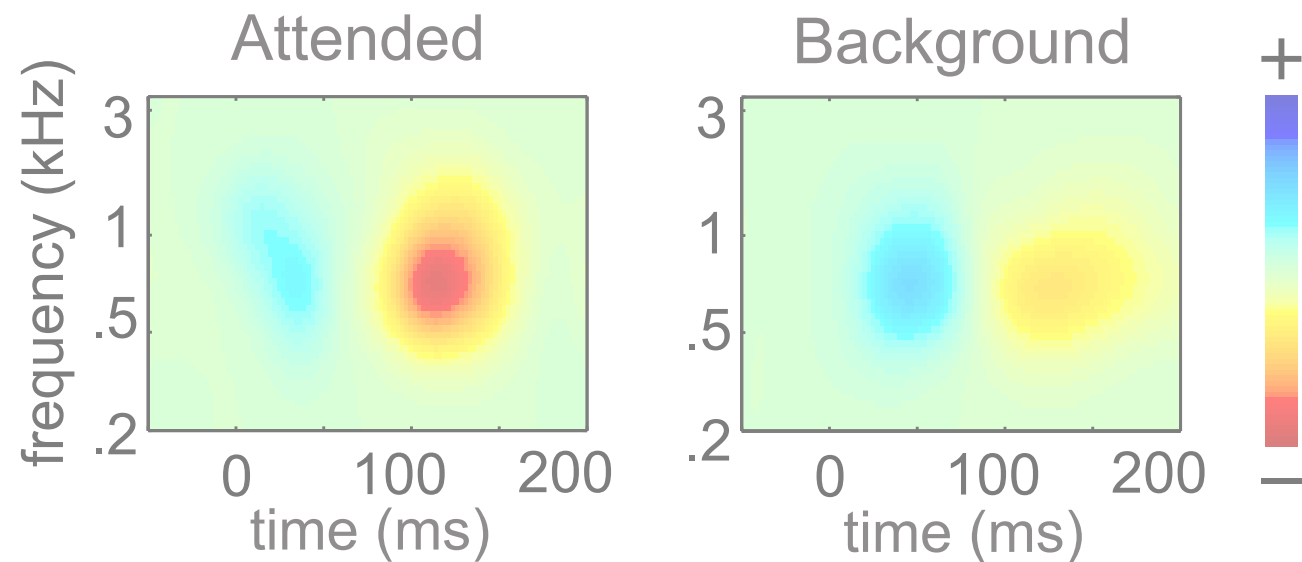


STRF Results

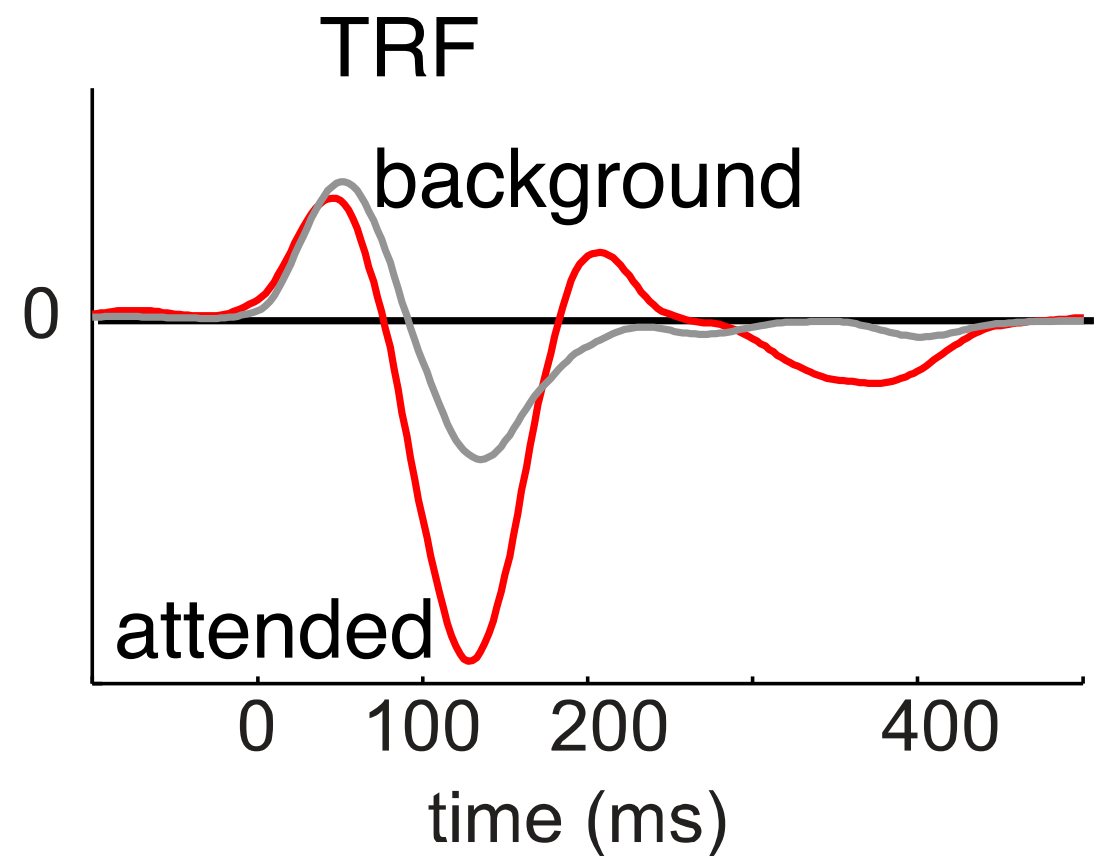


- STRF separable (time, frequency)
- 300 Hz - 2 kHz dominant carriers
- M50_{STRF} positive peak
- M100_{STRF} negative peak

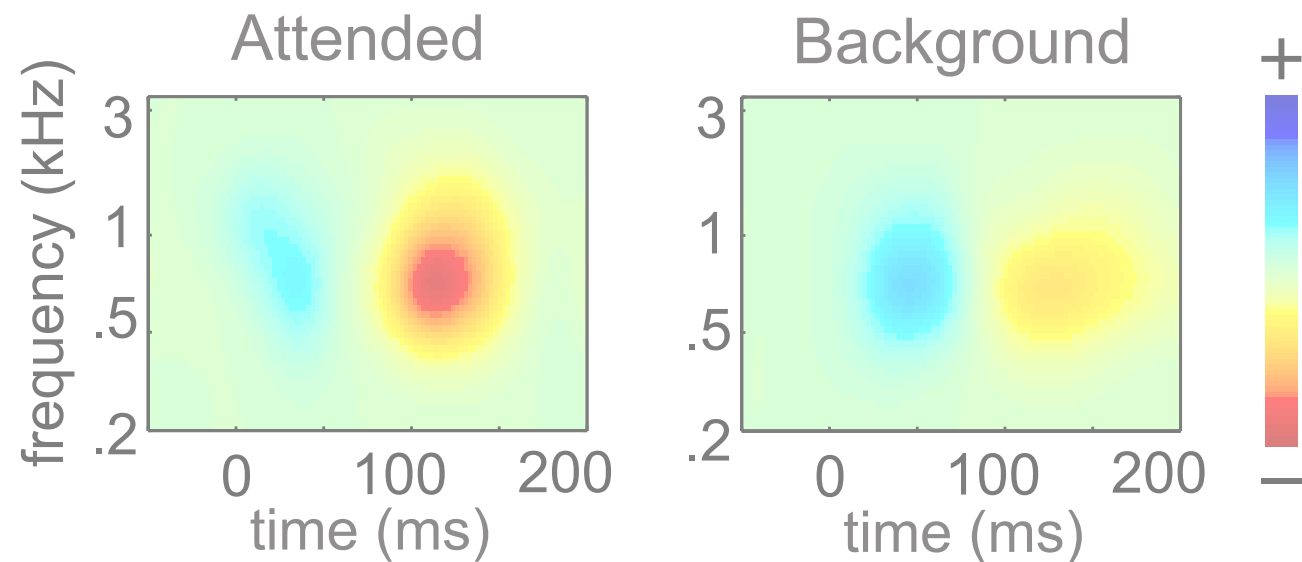
STRF Results



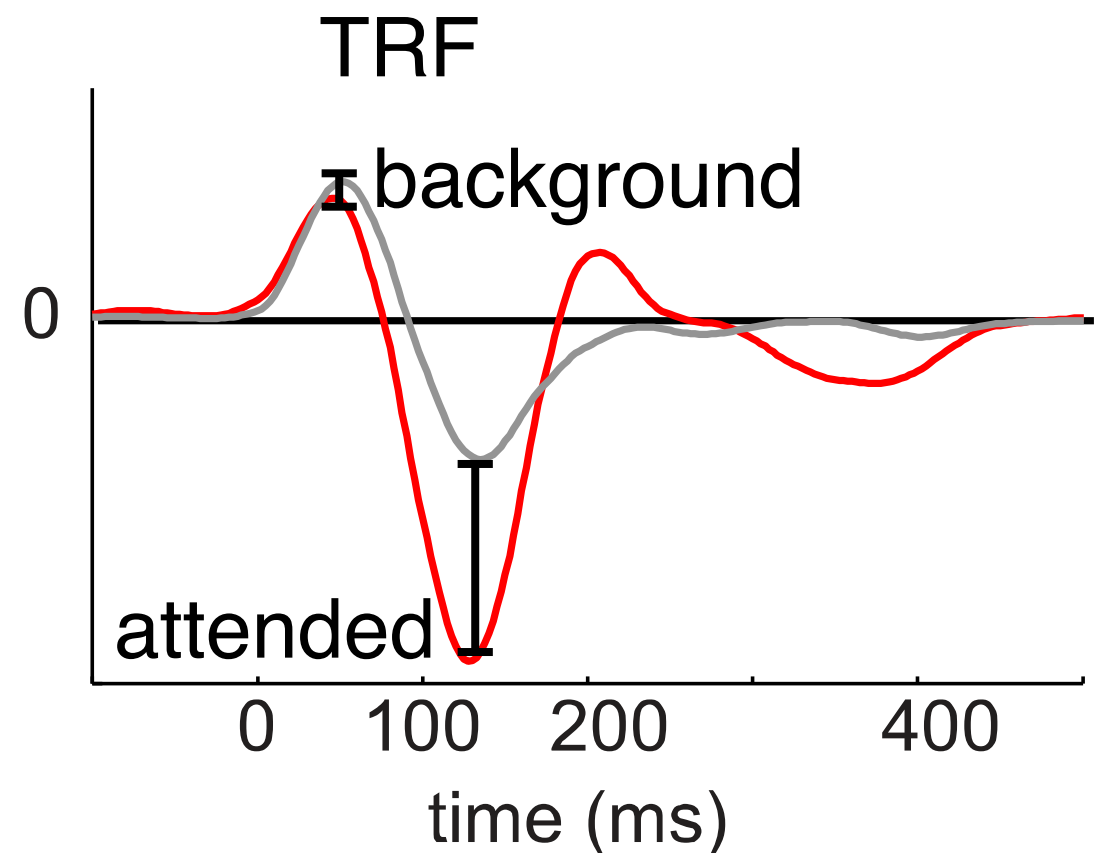
- STRF separable (time, frequency)
- 300 Hz - 2 kHz dominant carriers
- M50_{STRF} positive peak
- M100_{STRF} negative peak



STRF Results

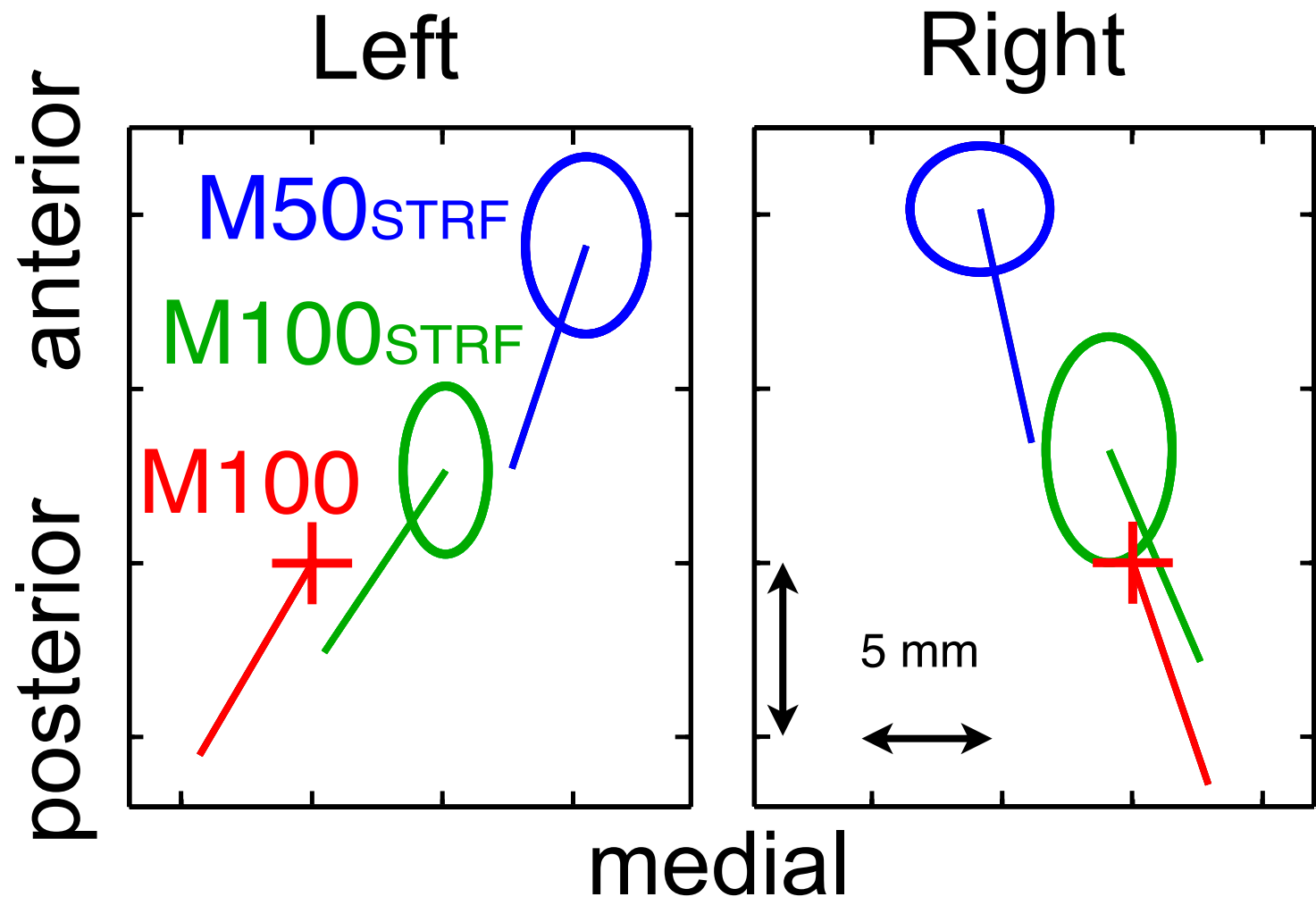


- STRF separable (time, frequency)
- 300 Hz - 2 kHz dominant carriers
- M50_{STRF} positive peak
- M100_{STRF} negative peak
- **M100_{STRF} strongly modulated by attention, *but not* M50_{STRF}**

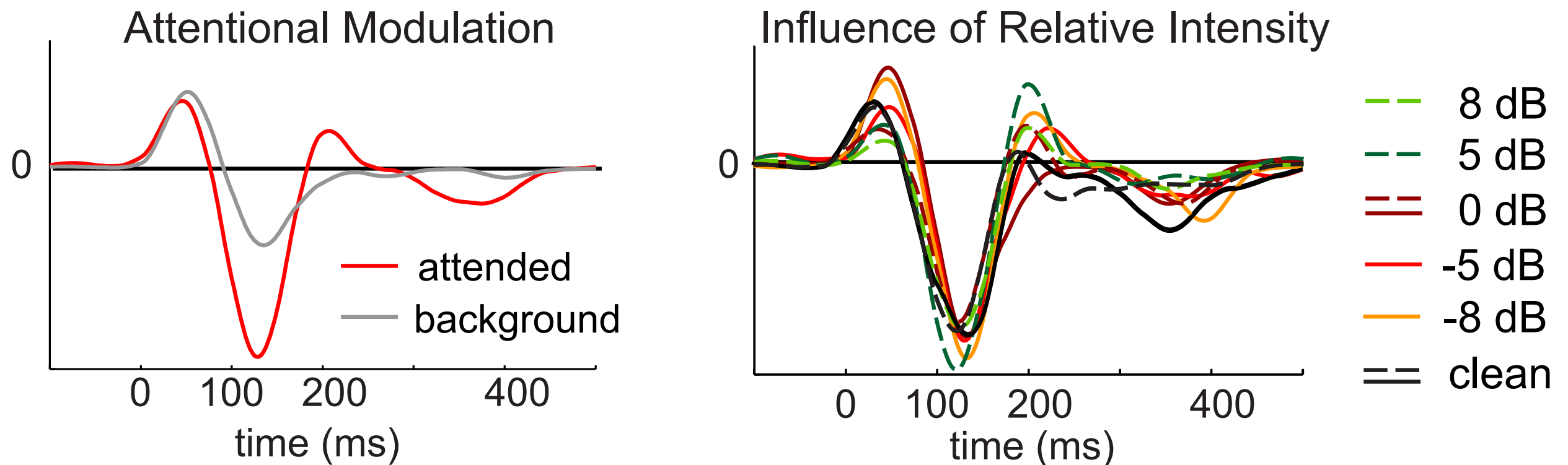


Neural Sources

- M100_{STRF} source near (same as?) M100 source:
Planum Temporale
- M50_{STRF} source is anterior and medial to M100 (same as M50?):
Heschl's Gyrus
- **PT strongly modulated by attention, *but not HG***



Cortical Object-Processing Hierarchy



- M100_{STRF} strongly modulated by attention, but not M50_{STRF}.
- M100_{STRF} invariant against acoustic changes.
- Objects well-neurally represented at 100 ms, but not 50 ms.

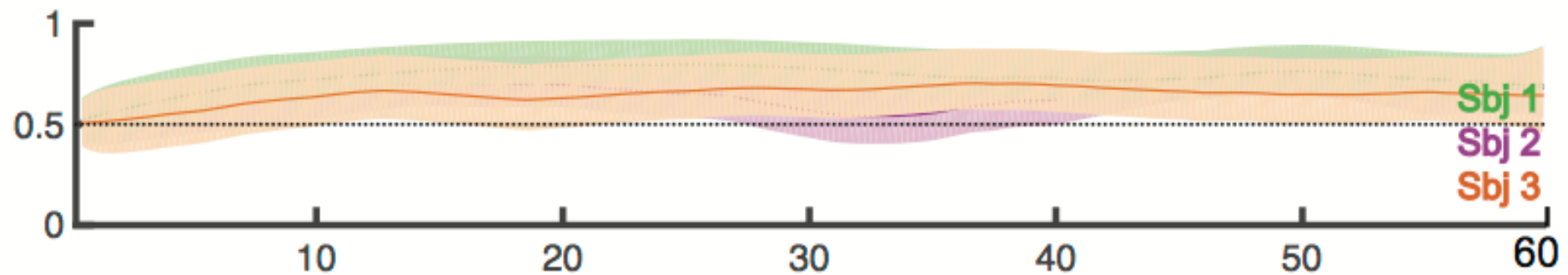
Studies In Progress

- Attentional Dynamics
- Aging & Neural Representations of Speech
- Neural Representations of the Background

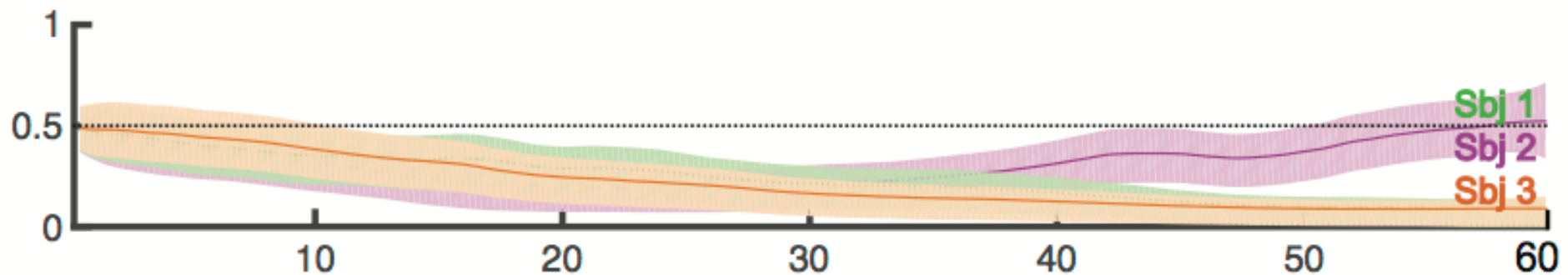
Attentional Dynamics

Attend to Speaker 1

Probability
of attending
Speaker 1



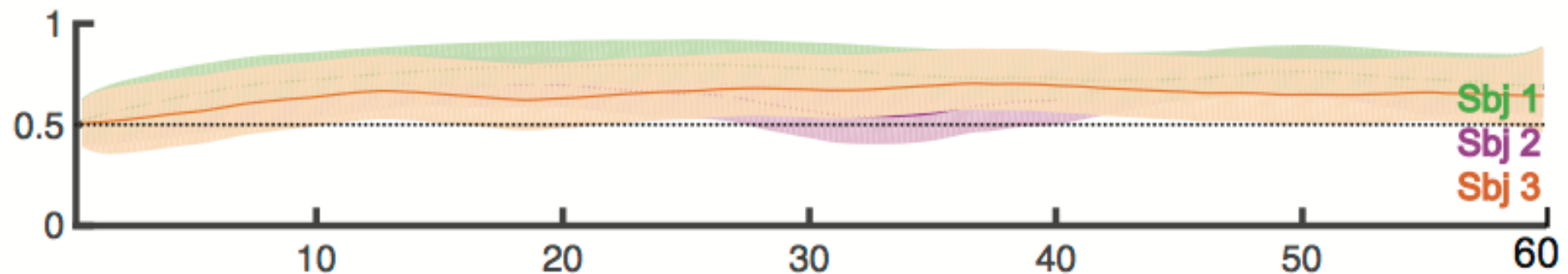
Attend to Speaker 2



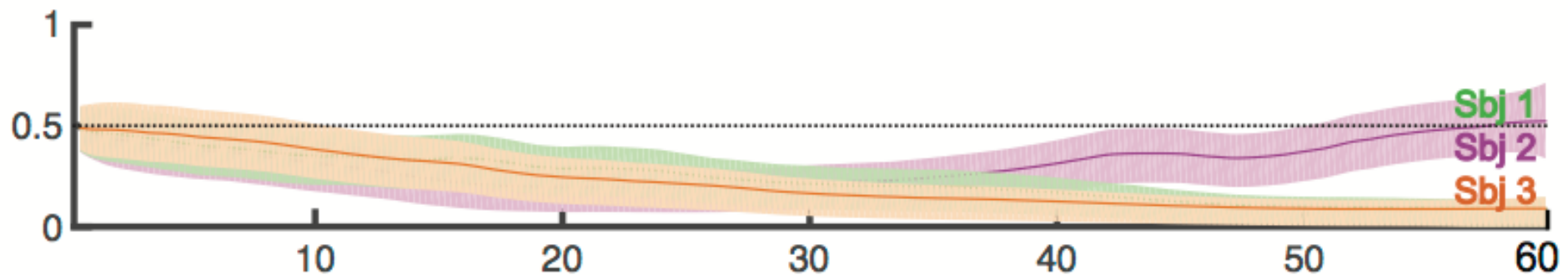
Attentional Dynamics

Attend to Speaker 1

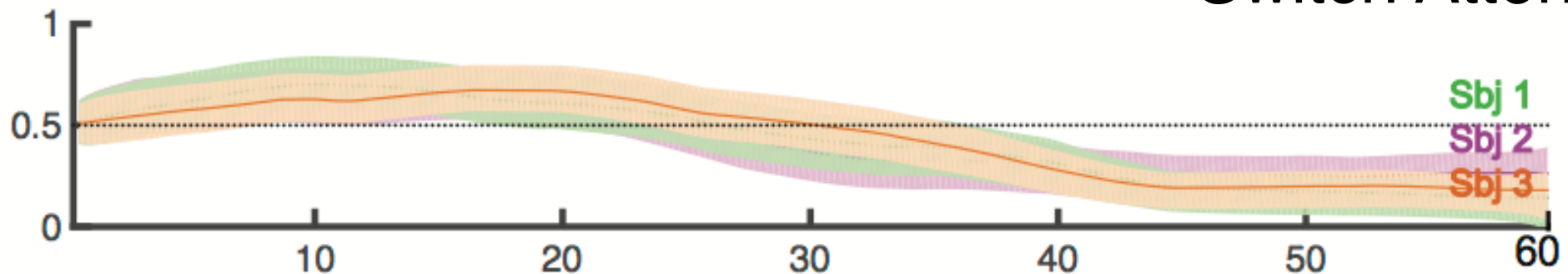
Probability
of attending
Speaker 1



Attend to Speaker 2

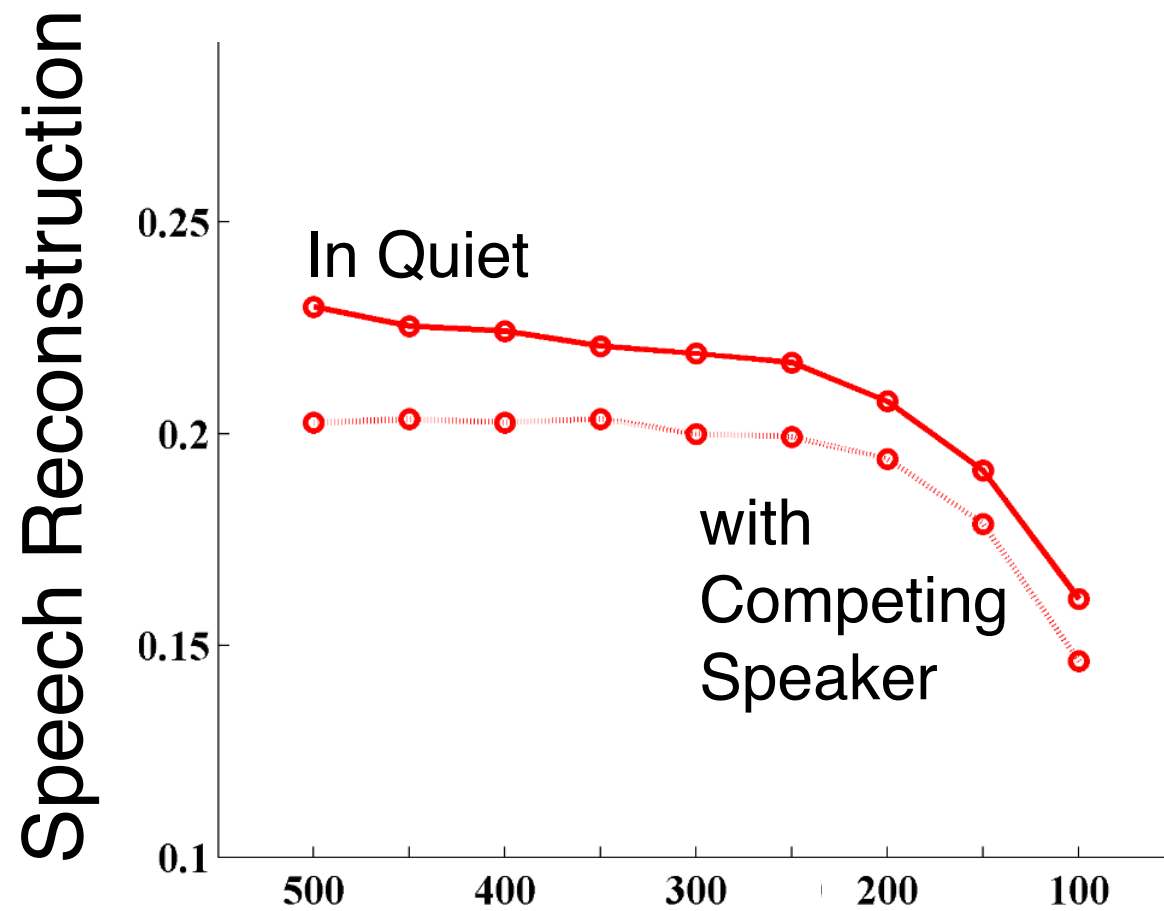


Switch Attention

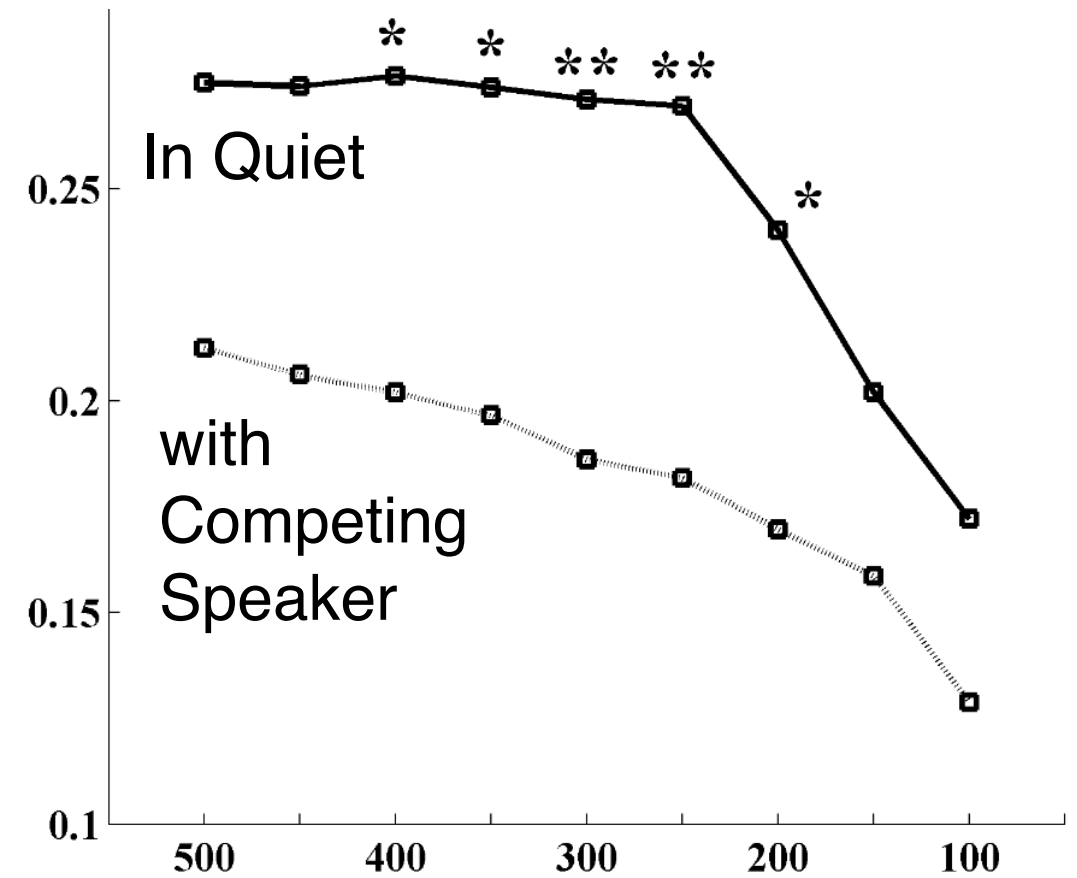


Younger vs. Older Listeners

Younger Adults



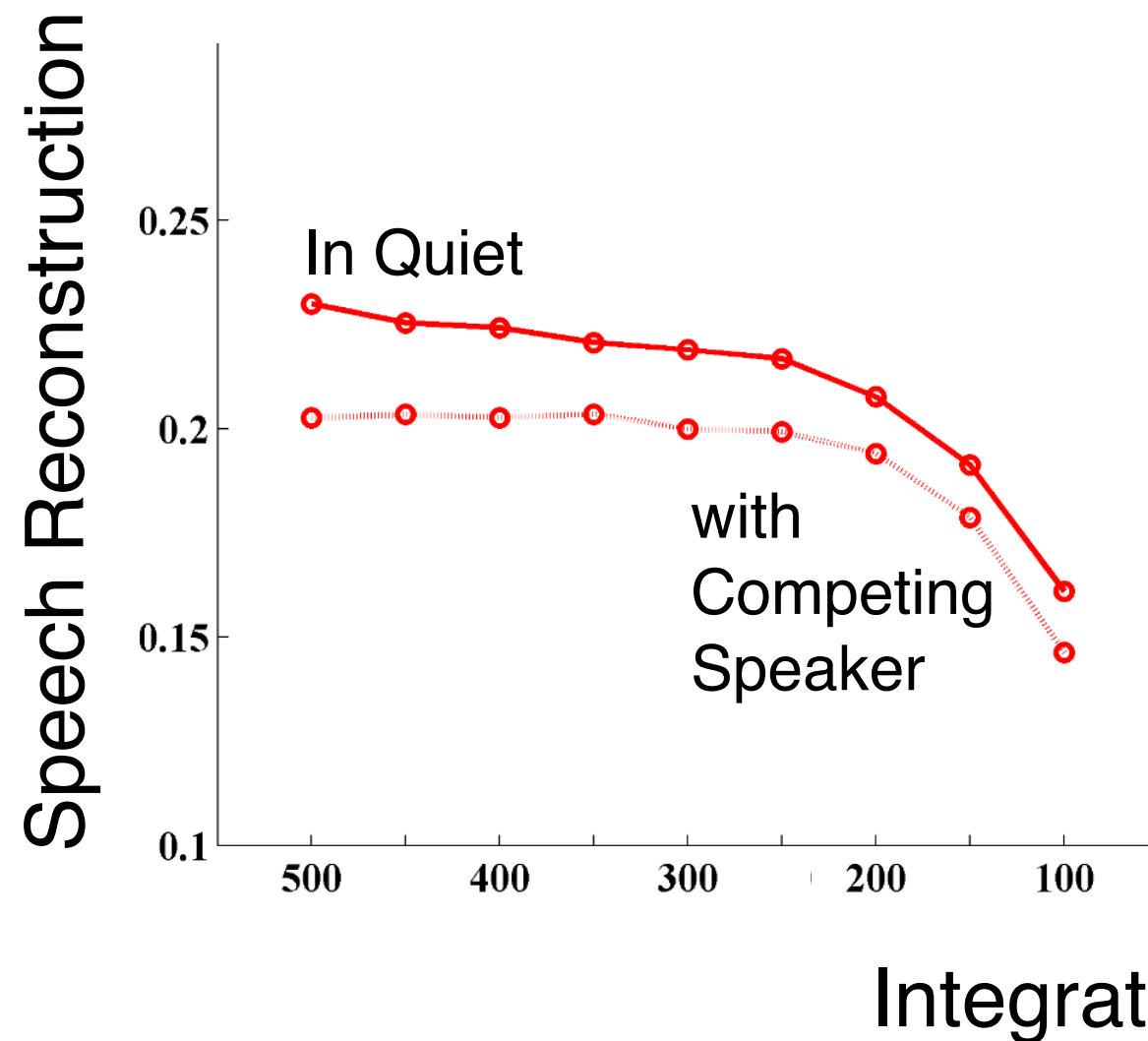
Older Adults



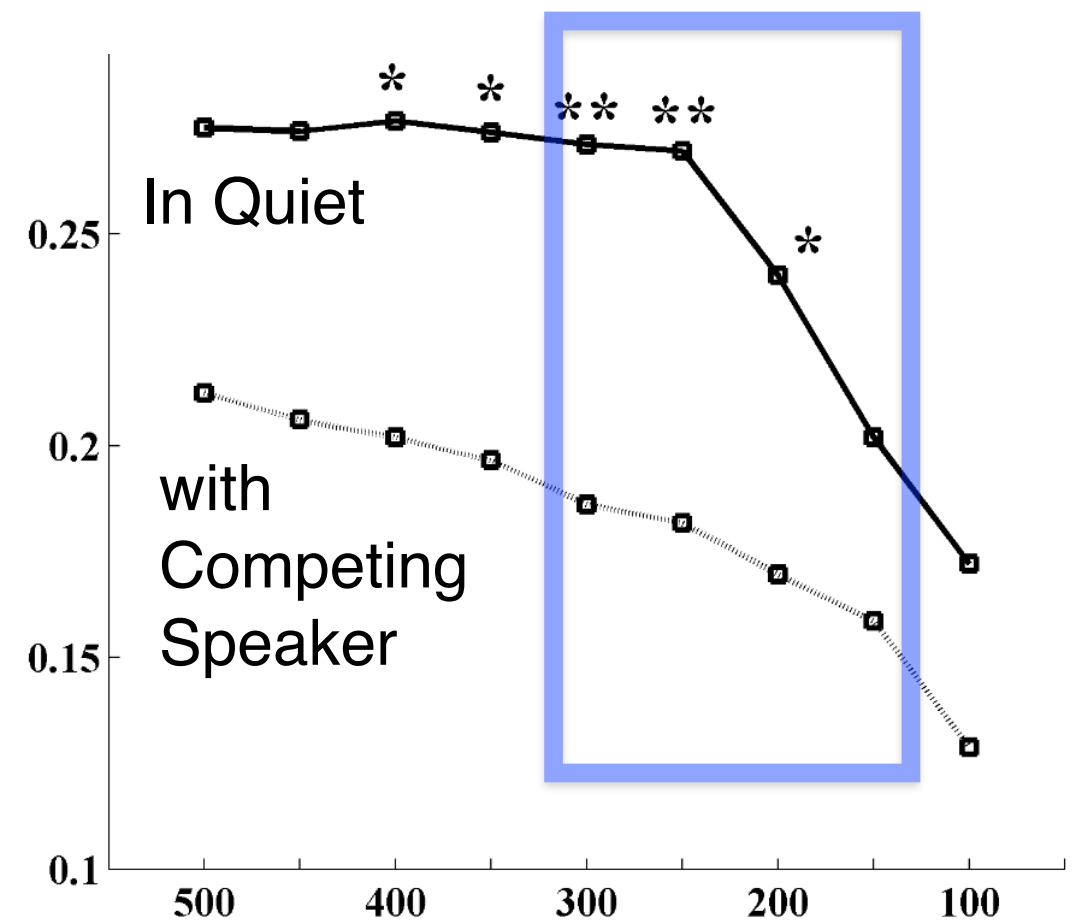
Integration window (ms)

Younger vs. Older Listeners

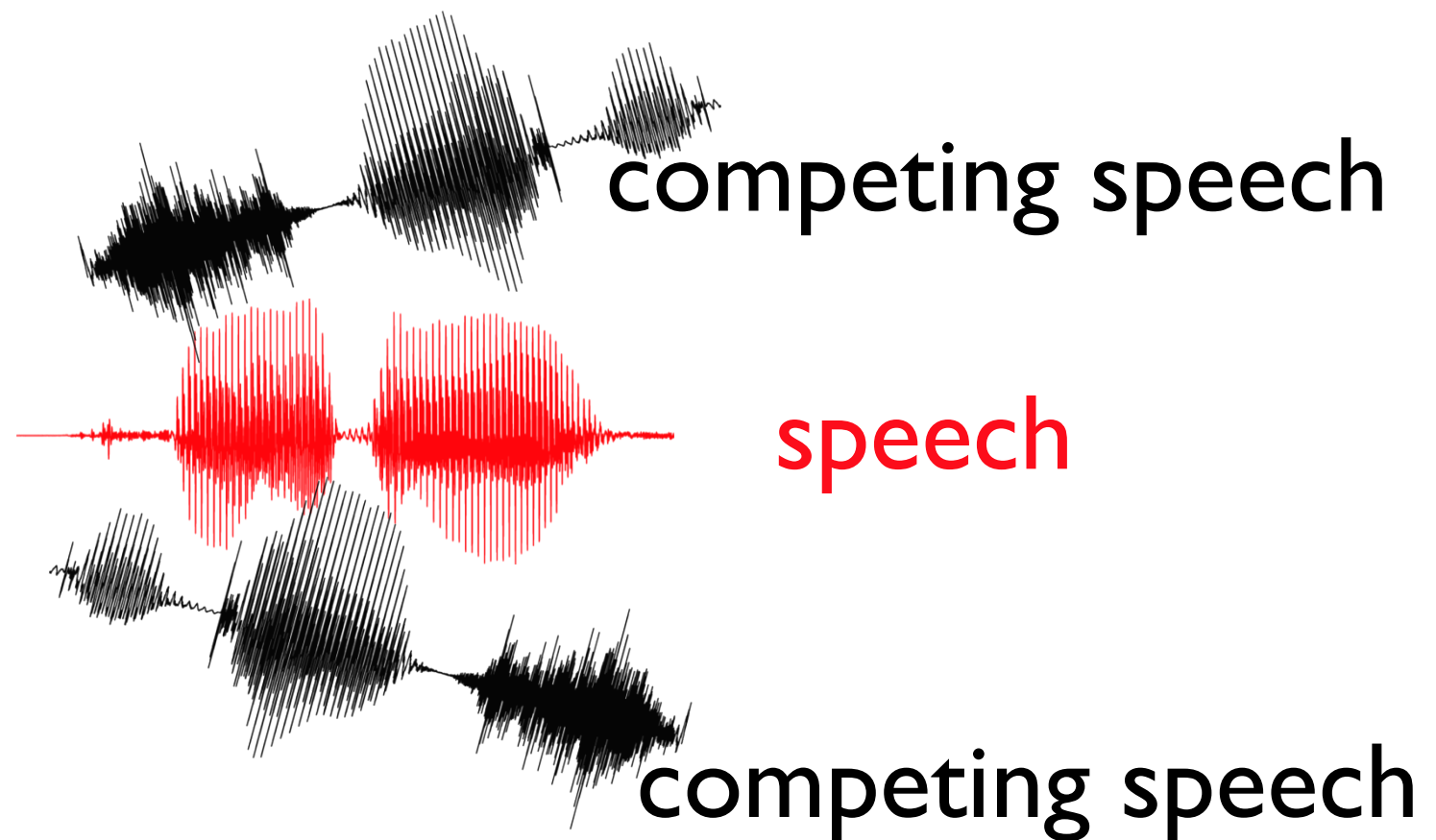
Younger Adults



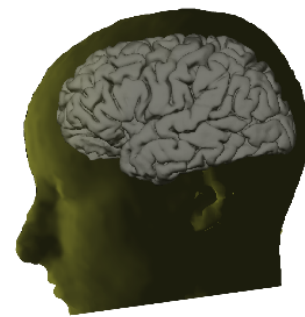
Older Adults



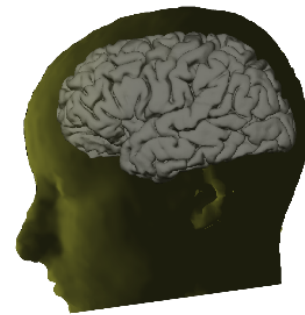
Three Competing Speakers



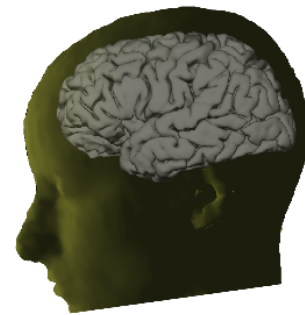
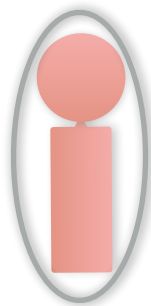
Foreground vs. Background



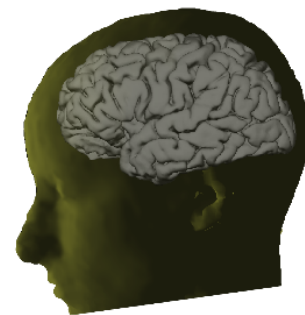
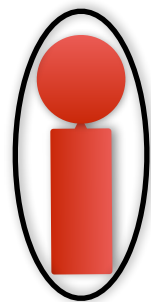
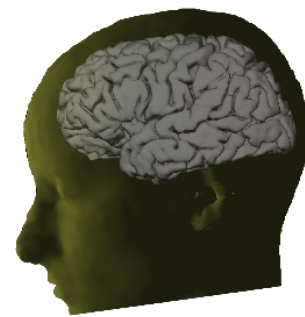
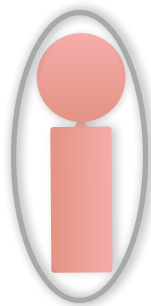
Foreground vs. Background



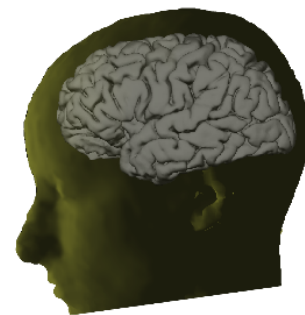
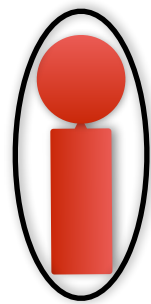
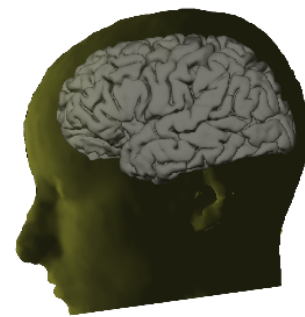
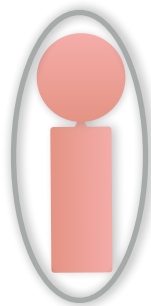
Foreground vs. Background



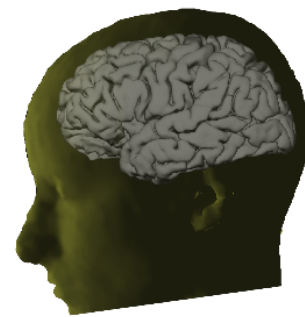
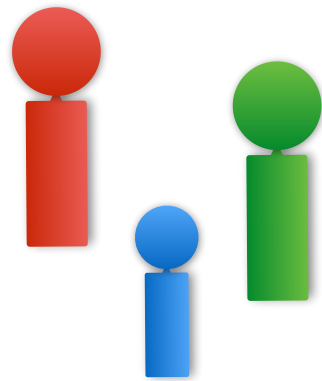
Foreground vs. Background



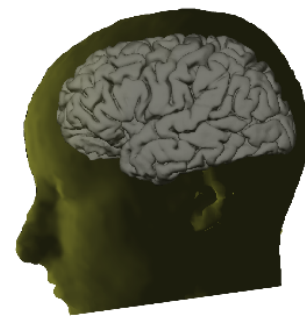
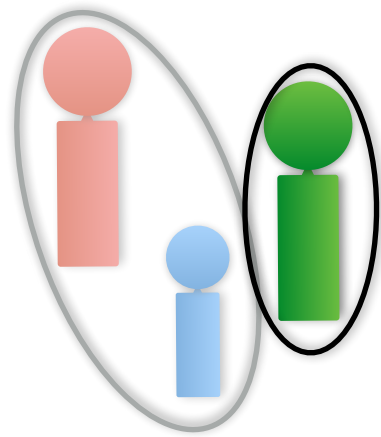
Foreground vs. Background



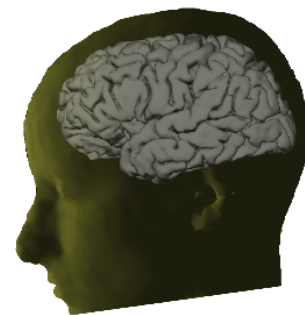
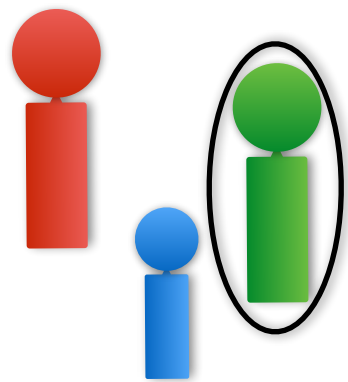
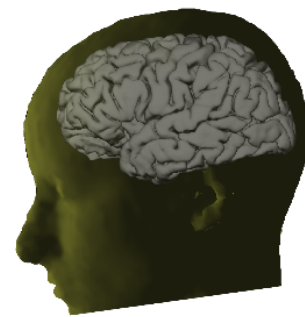
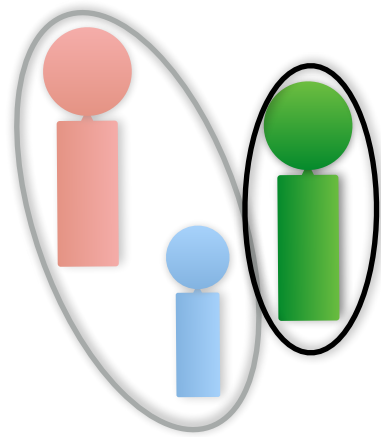
Foreground vs. Background



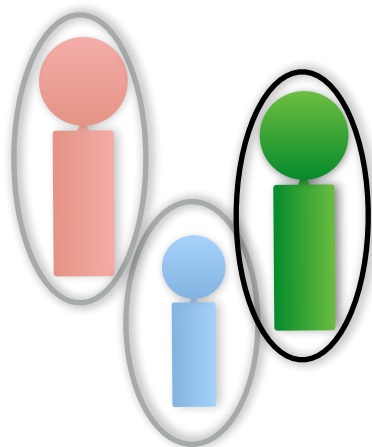
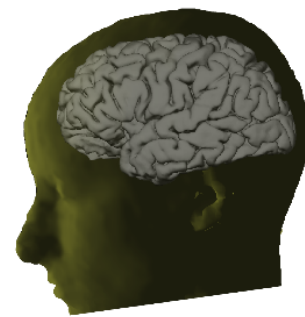
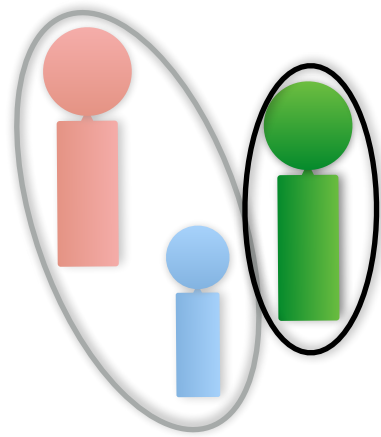
Foreground vs. Background



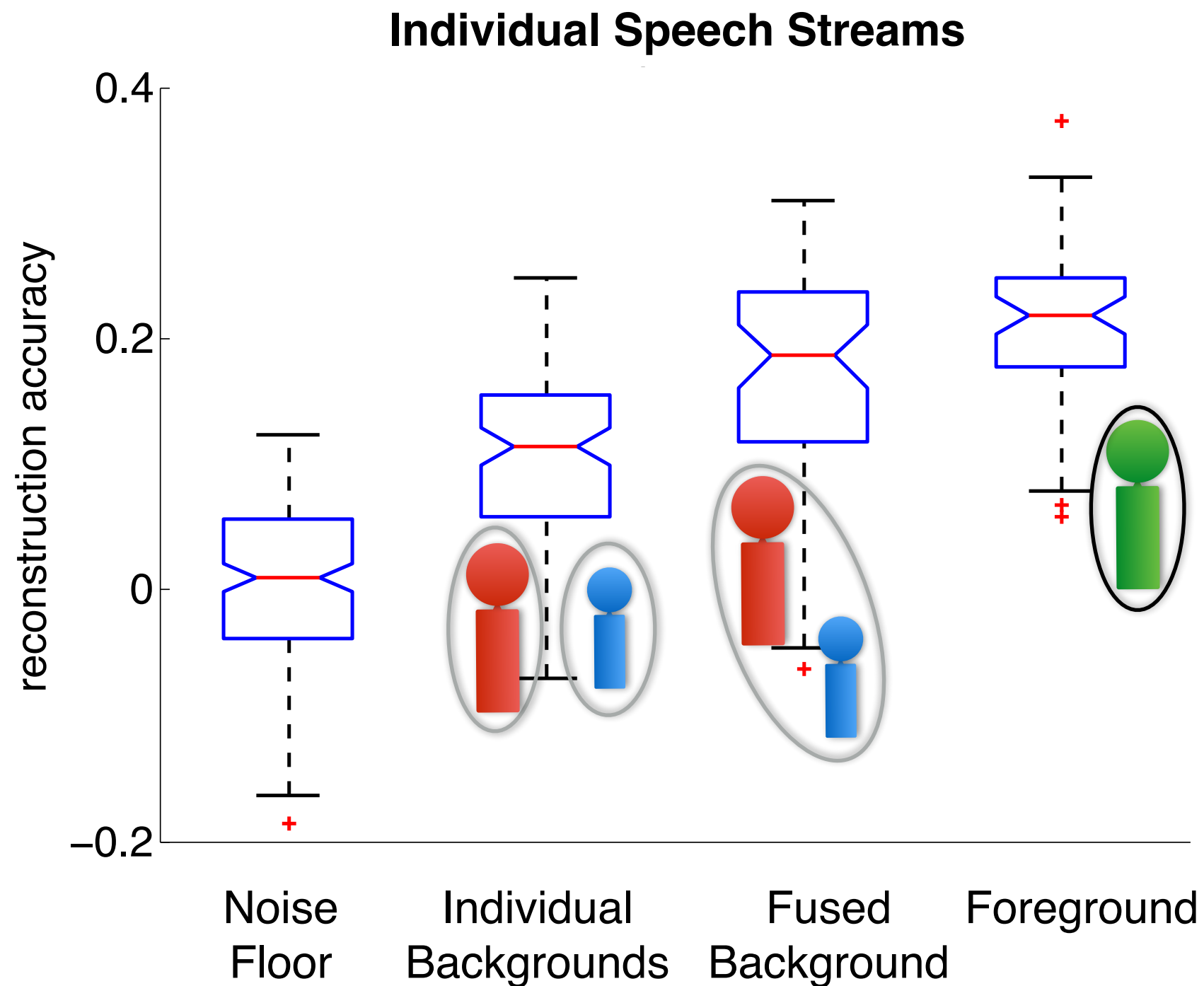
Foreground vs. Background



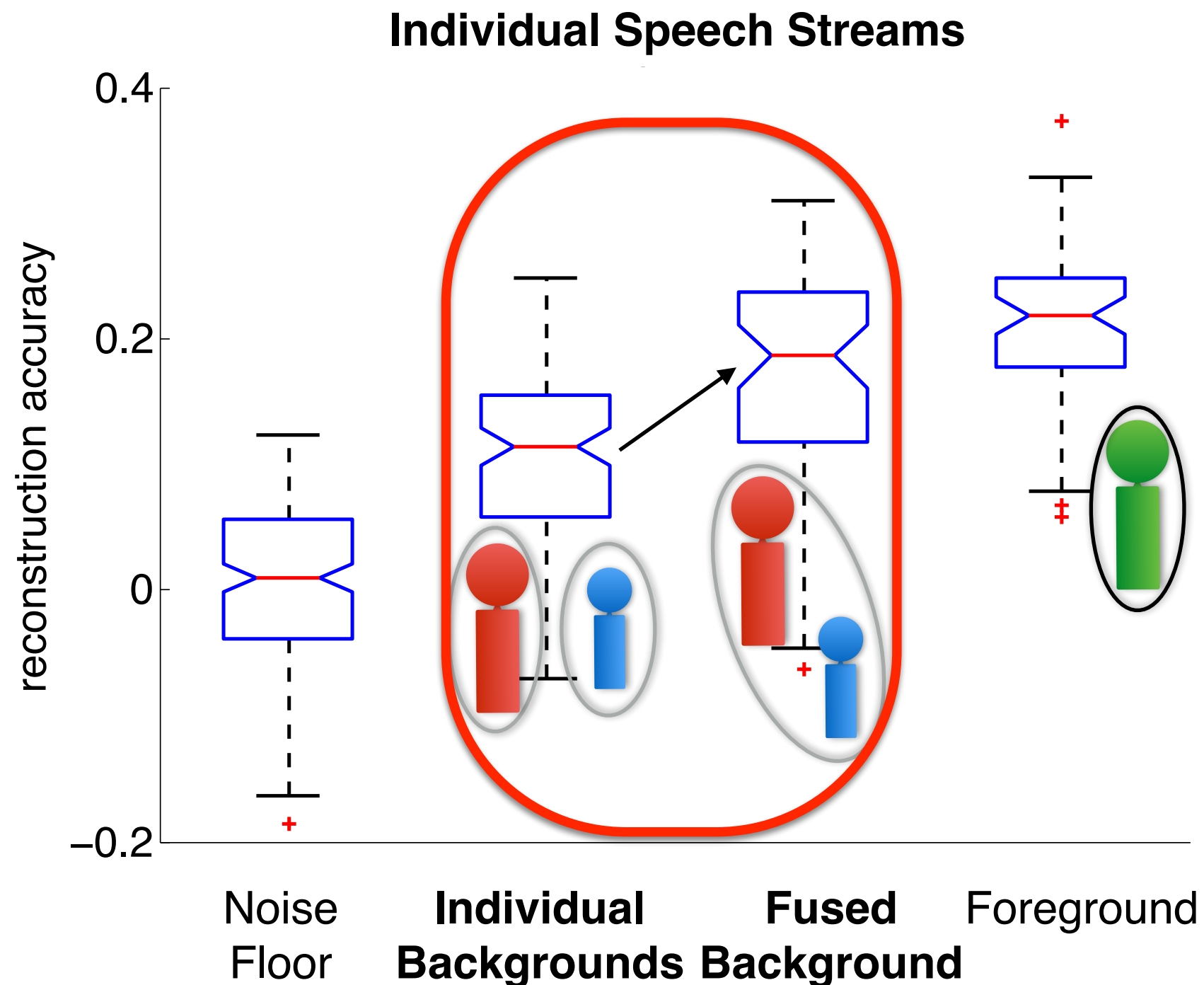
Foreground vs. Background



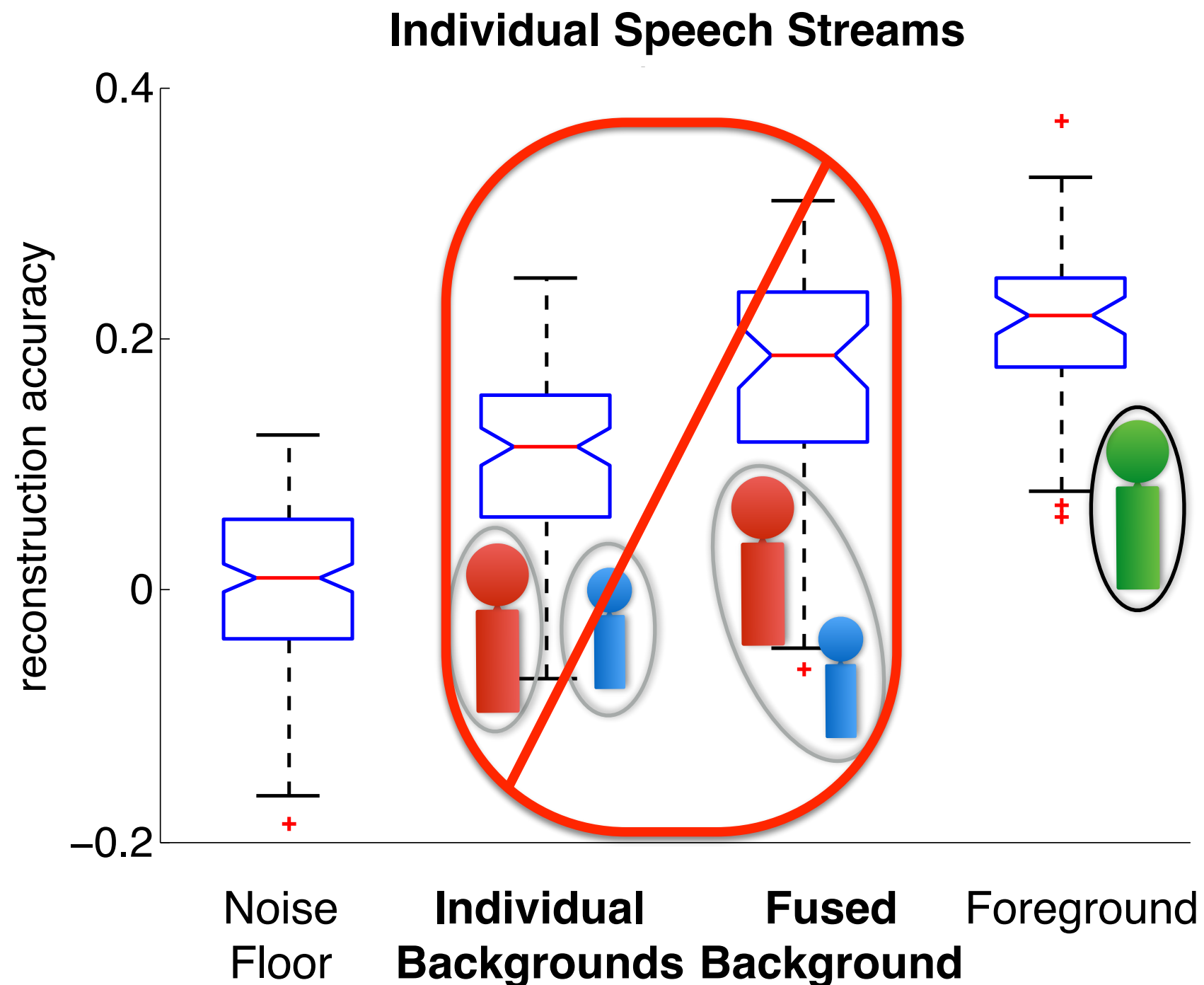
Backgrounds vs. Background



Backgrounds vs. Background



Backgrounds vs. Background



Integration Window over Late Times Only

Backgrounds vs. Background

Why not?

Stimulus Background

Speaker 1



Speaker 2



MEG Response

Two Speakers



Backgrounds vs. Background

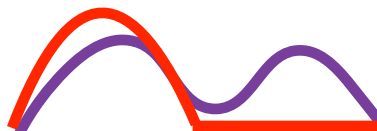
Why not?

Stimulus Background

Speaker 1



Speaker 2



MEG Response

Two Speakers



Backgrounds vs. Background

Why not?

Stimulus Background

Speaker 1



Speaker 2



MEG Response

Two Speakers



Backgrounds vs. Background

Why not?

Stimulus Background

Speaker 1

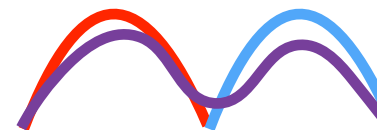


Speaker 2

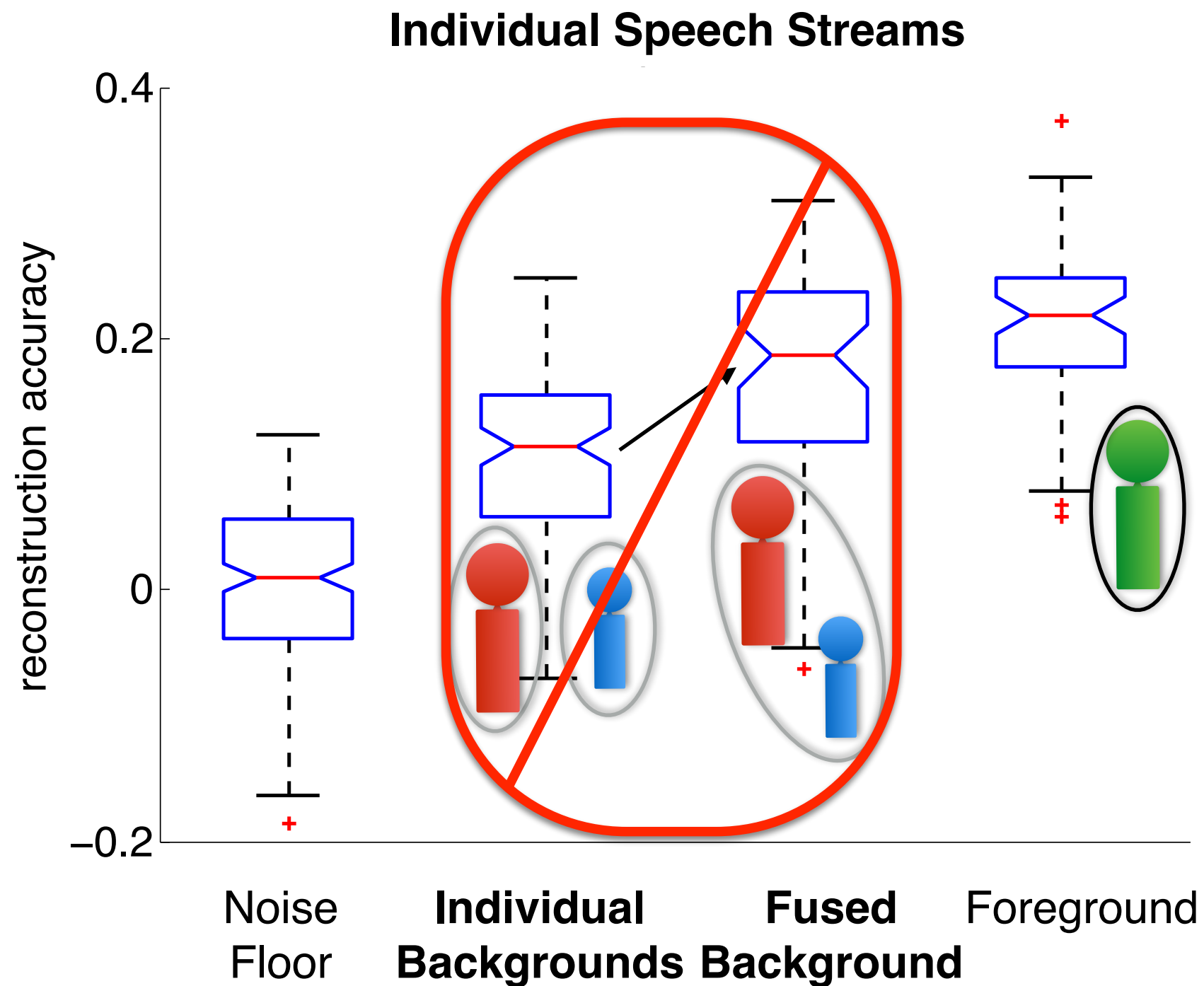


MEG Response

Two Speakers

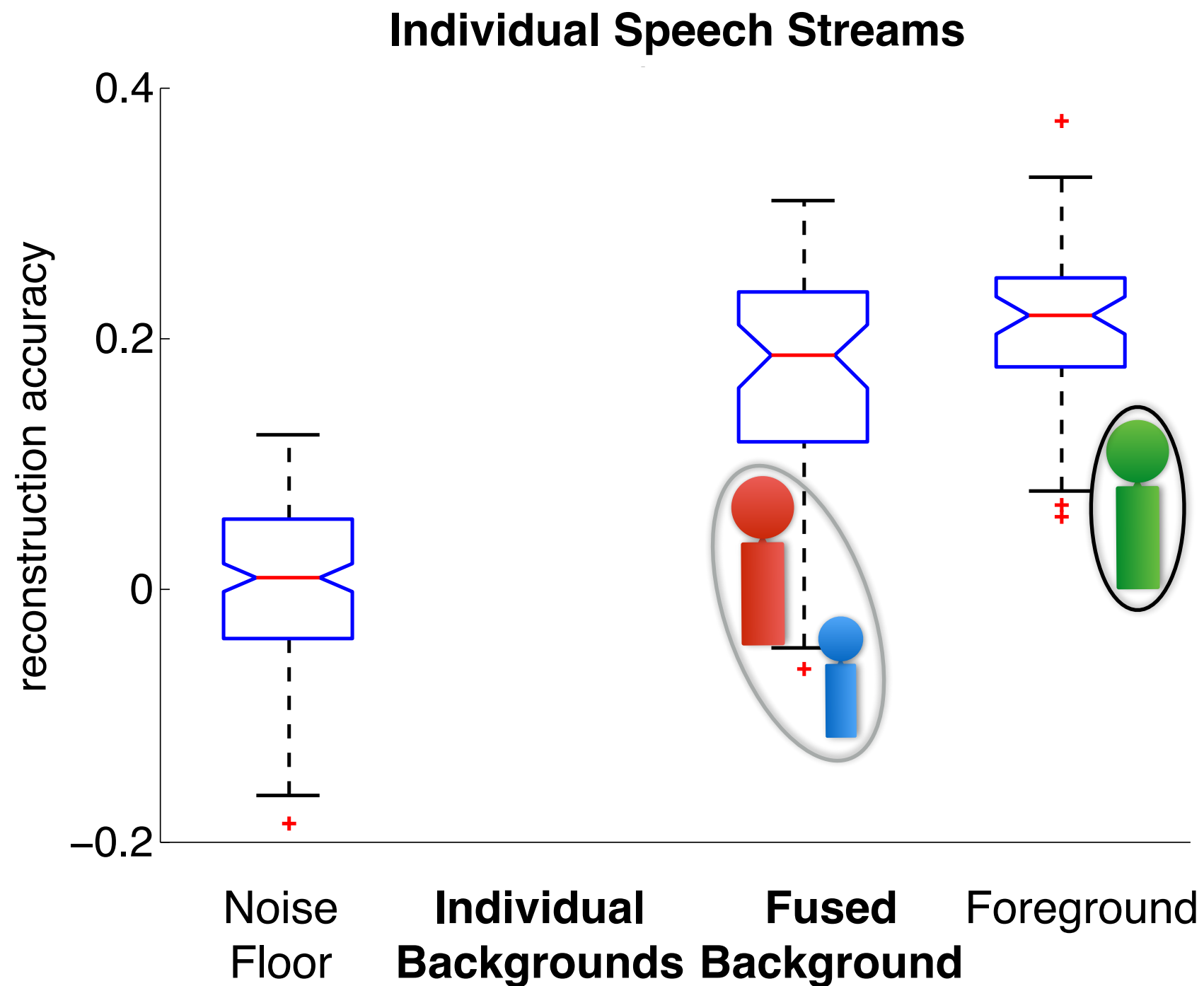


Backgrounds vs. Background



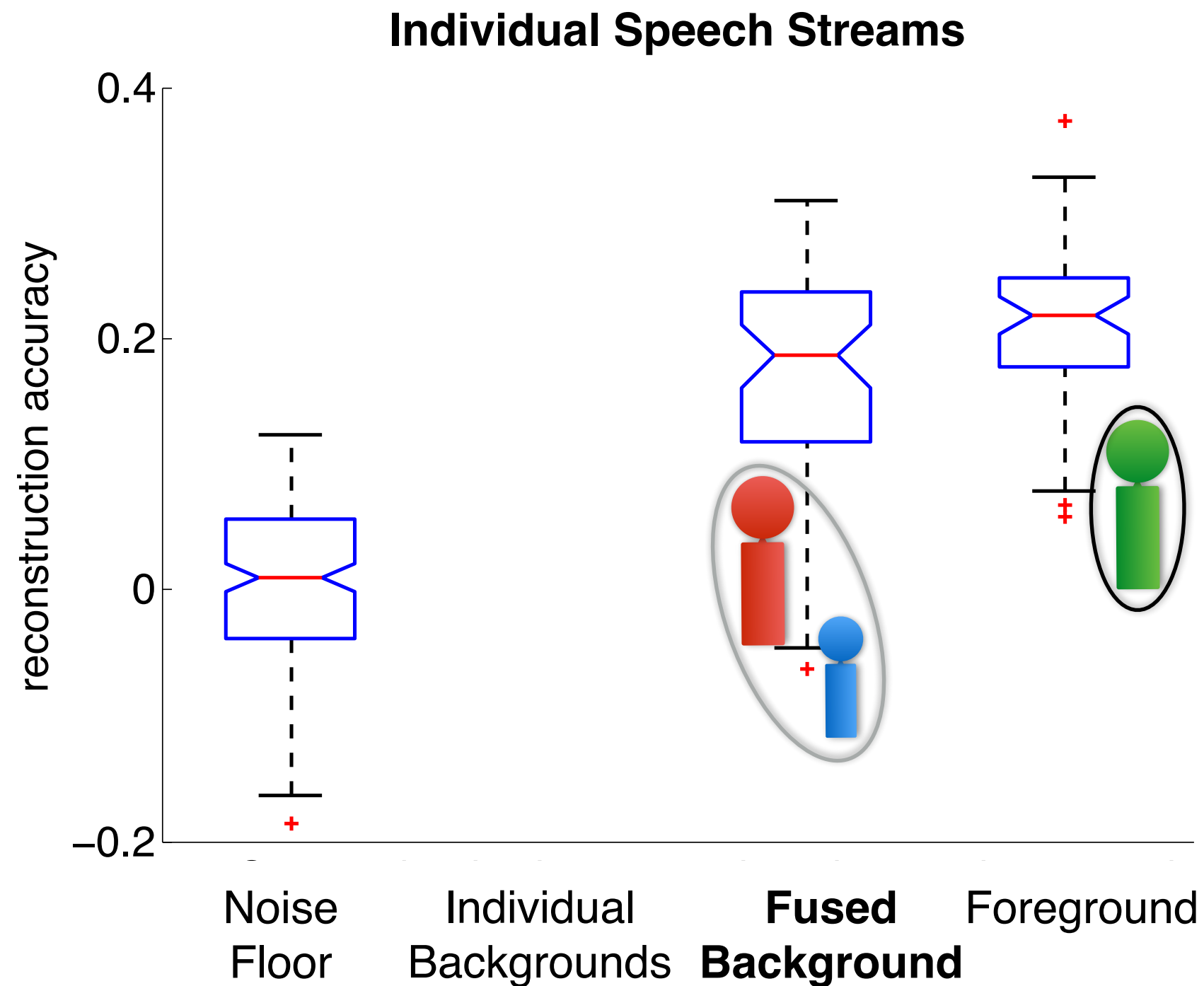
Integration Window over Late Times Only

Backgrounds vs. Background



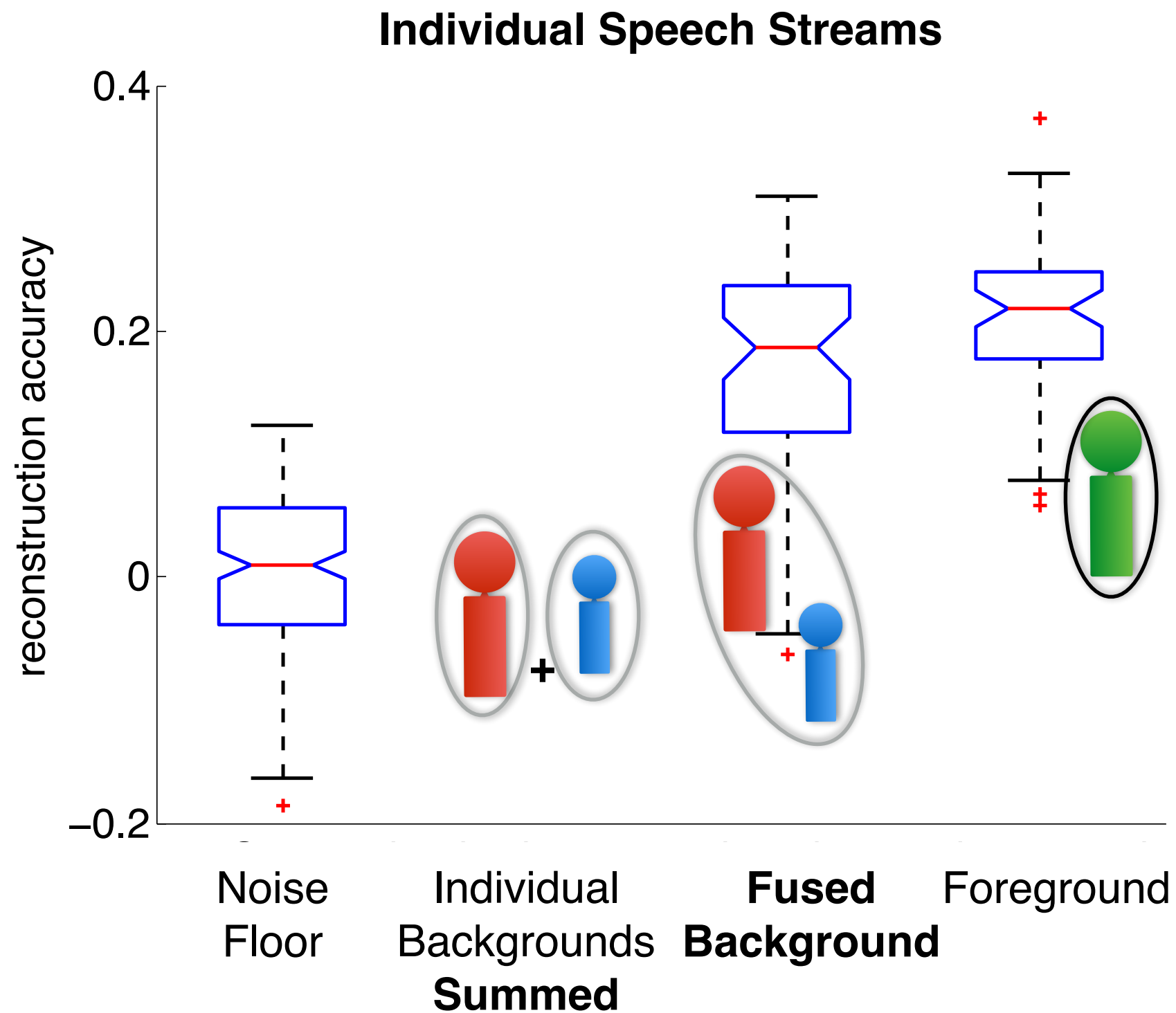
Integration Window over Late Times Only

Backgrounds vs. Background



Integration Window over Late Times Only

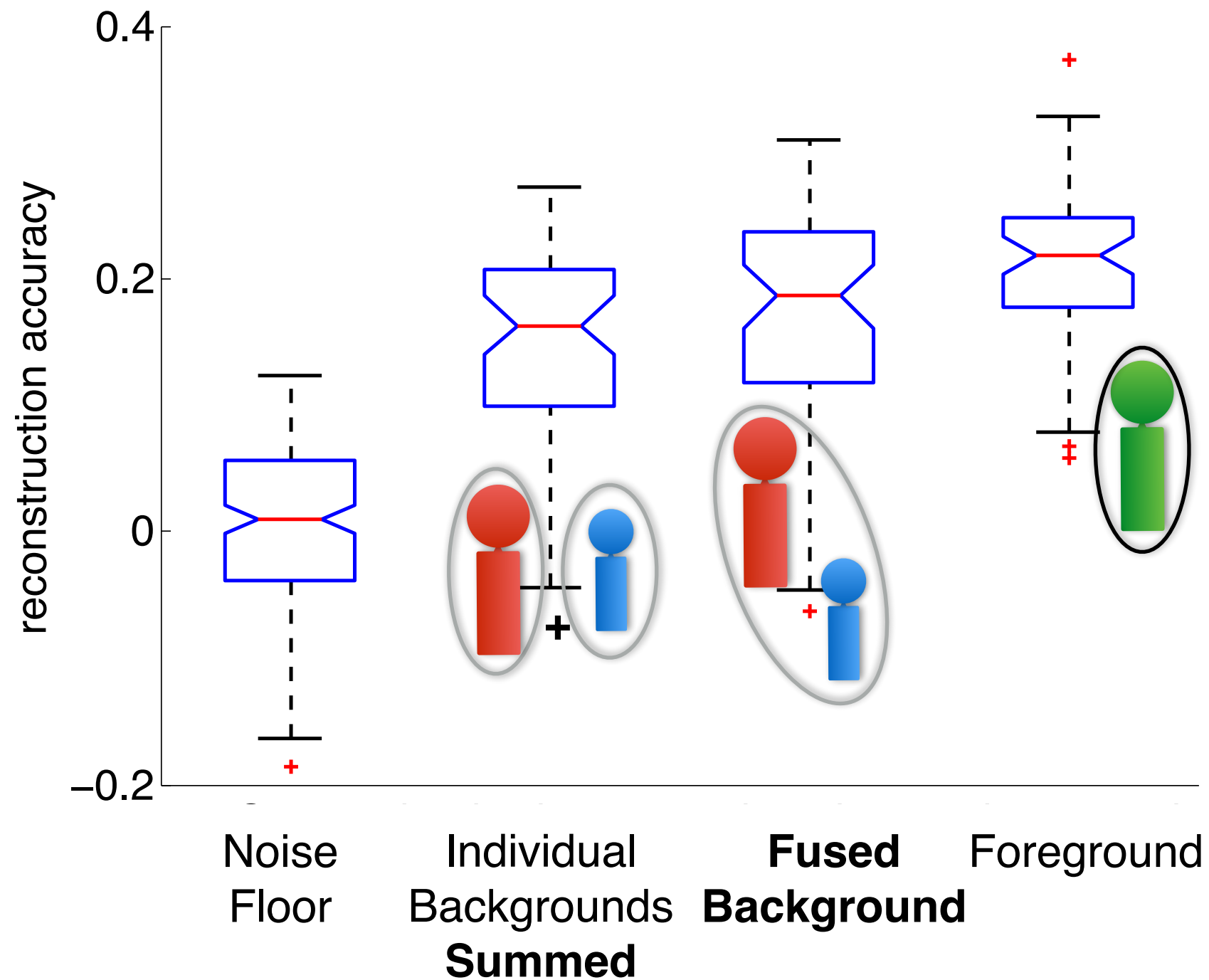
Backgrounds vs. Background



Integration Window over Late Times Only 

Backgrounds vs. Background

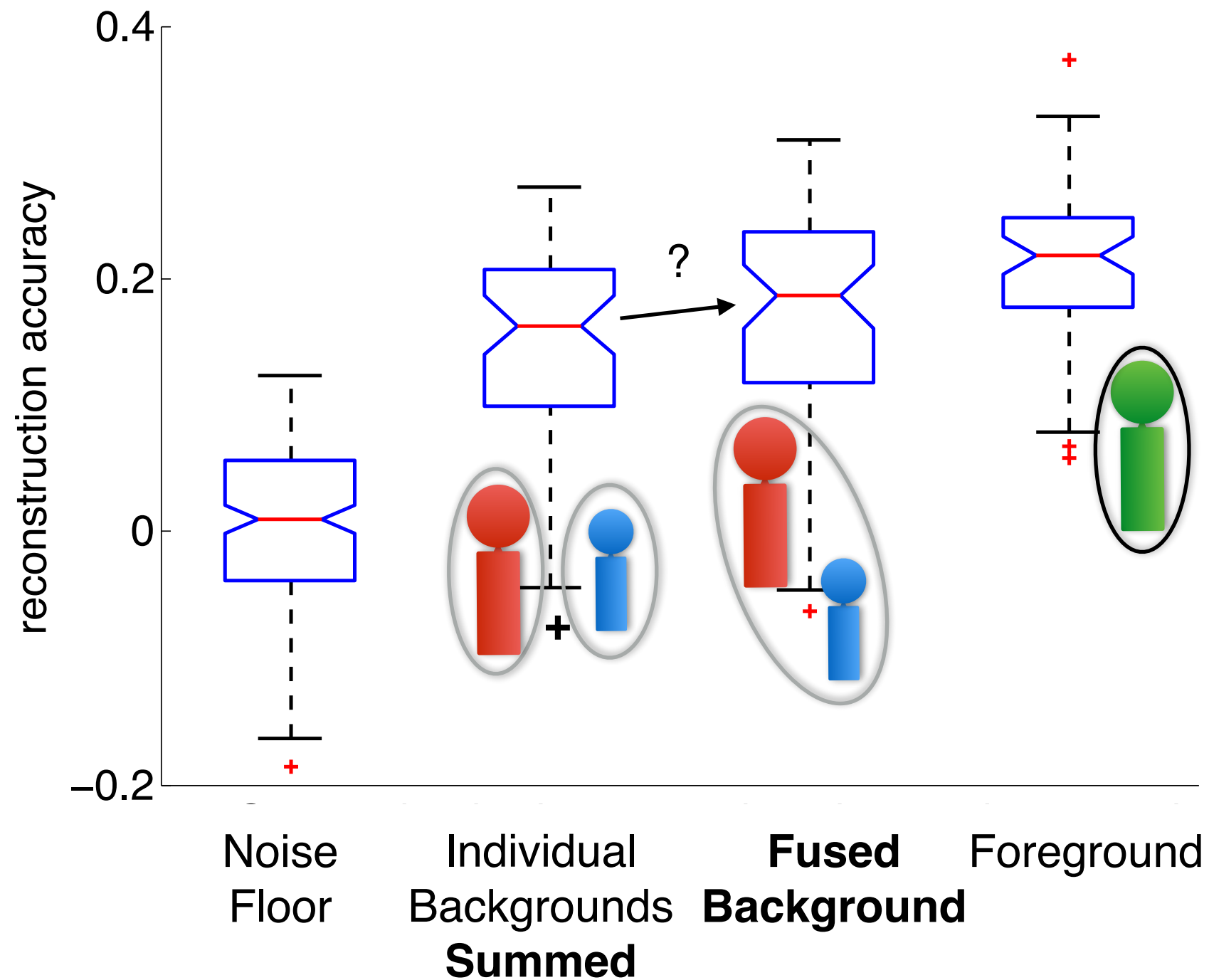
Individual Speech Streams



Integration Window over Late Times Only 

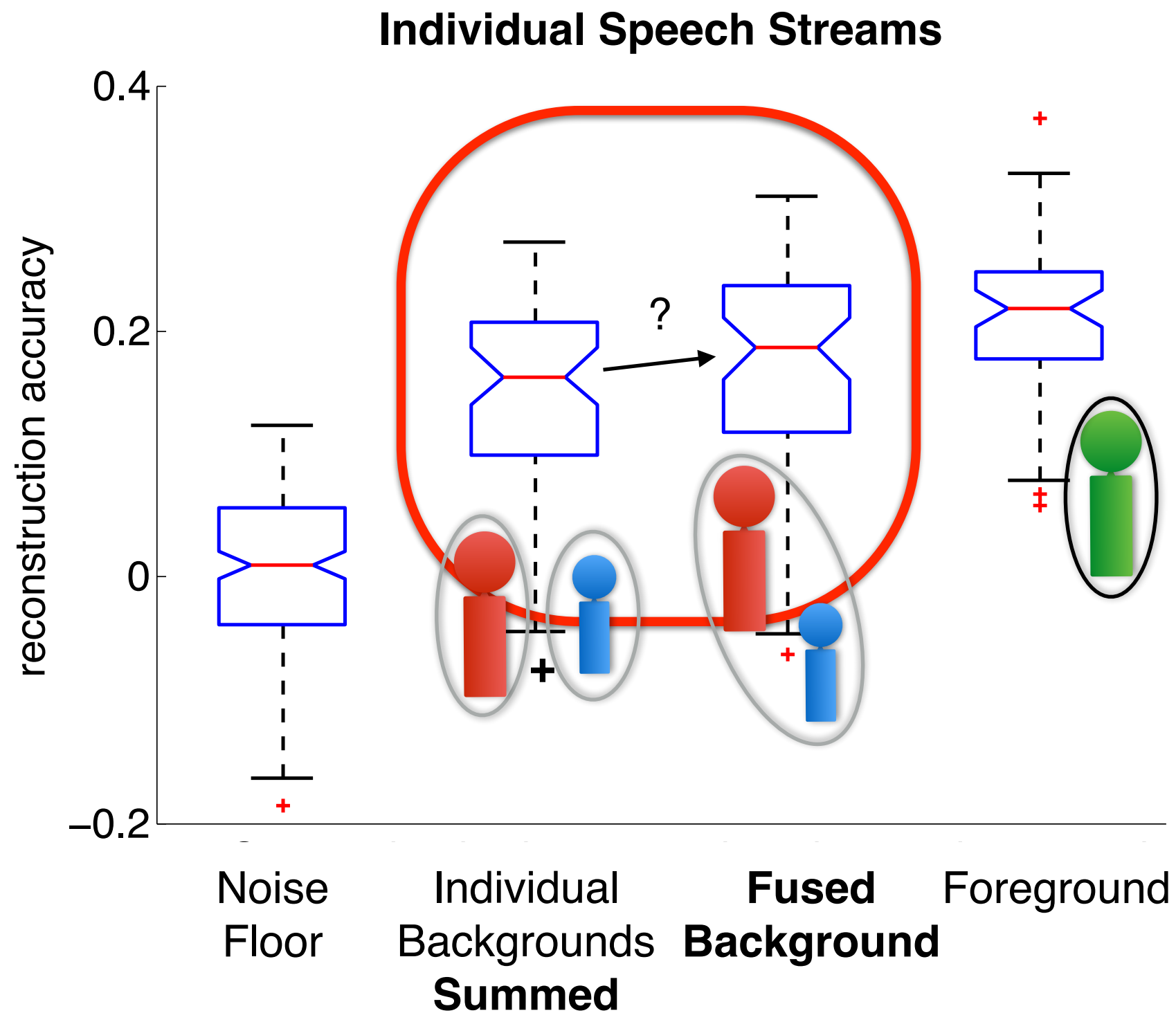
Backgrounds vs. Background

Individual Speech Streams



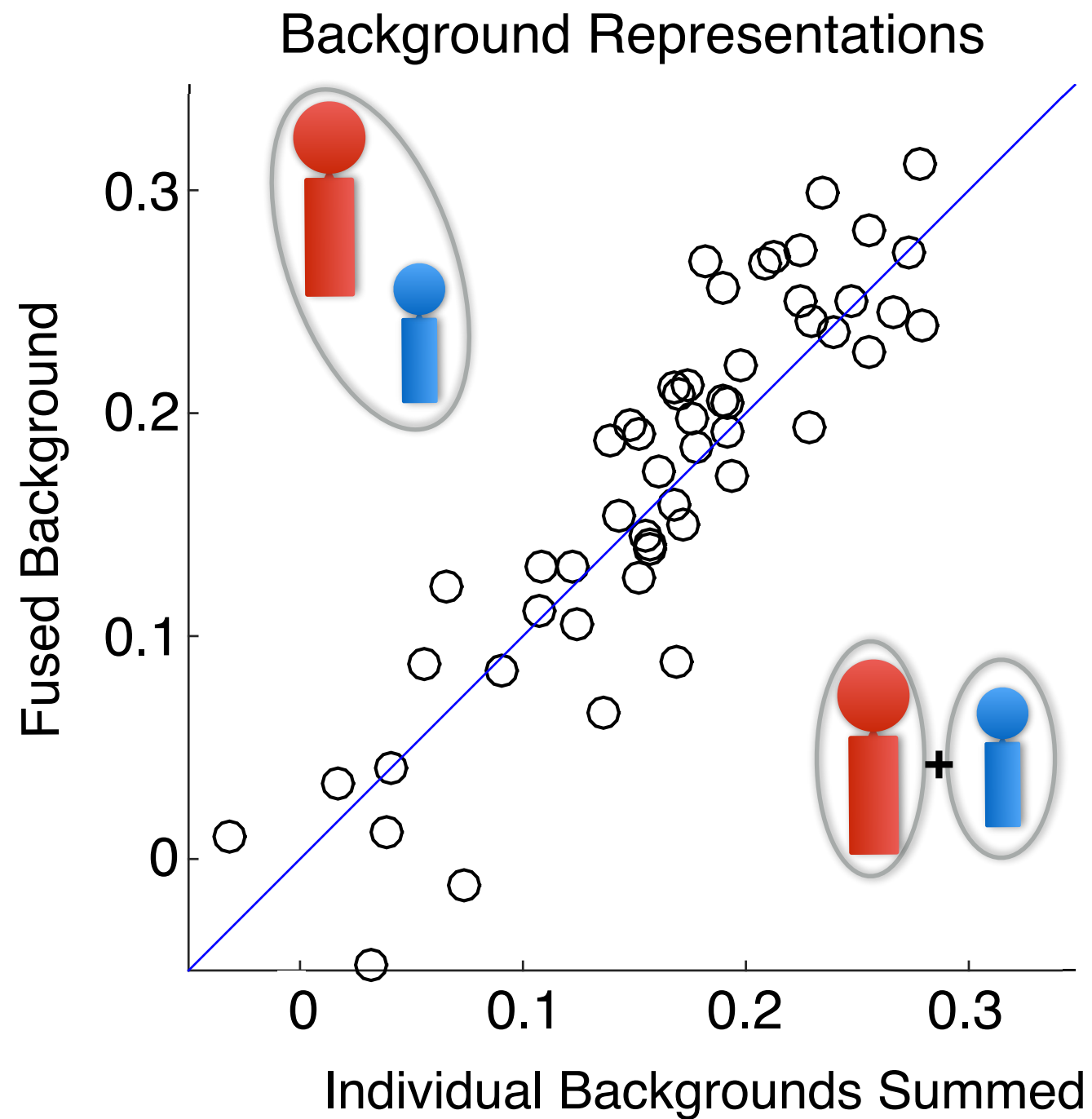
Integration Window over Late Times Only

Backgrounds vs. Background



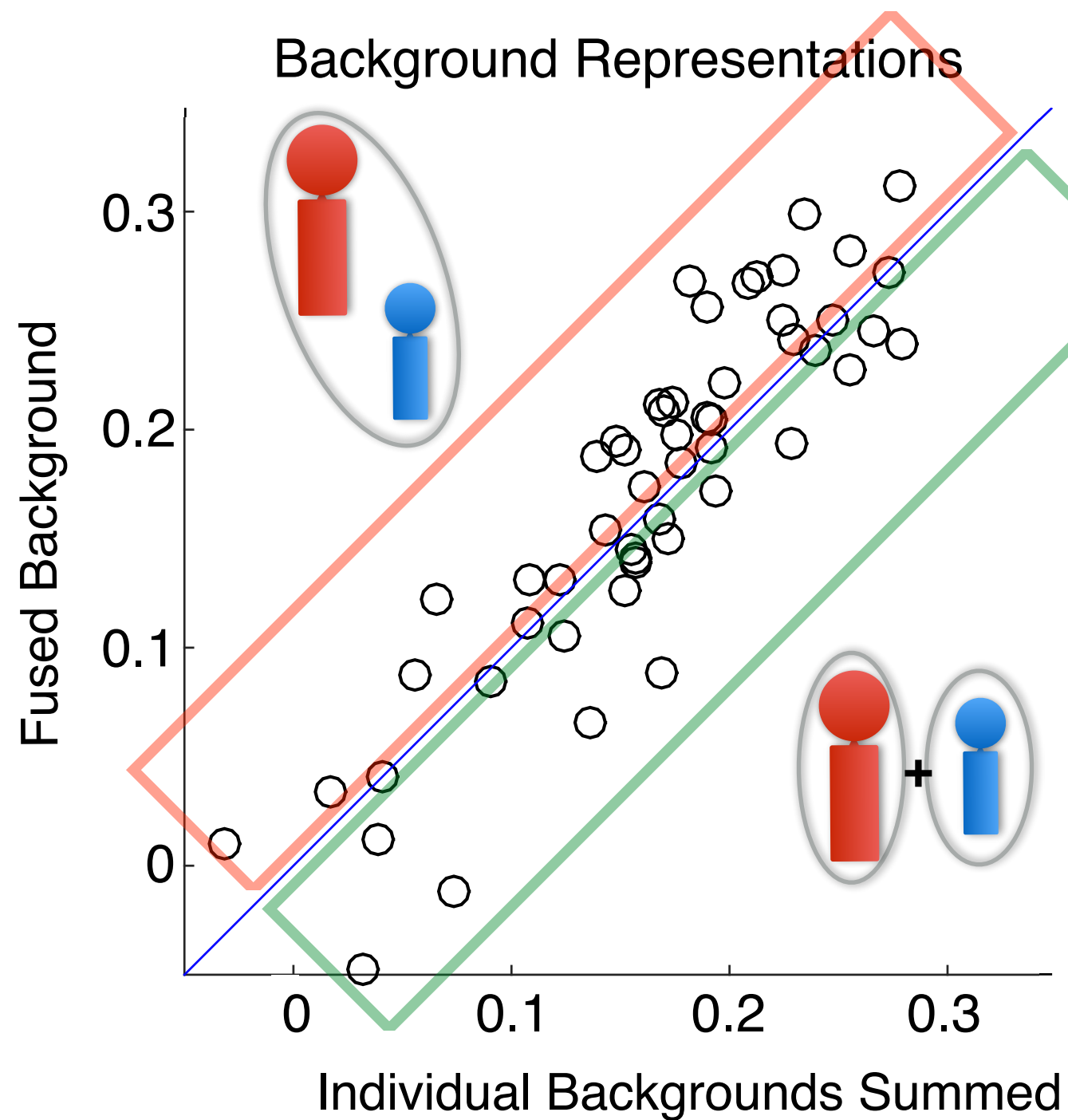
Integration Window over Late Times Only

Backgrounds vs. Background



Integration Window over Late Times Only

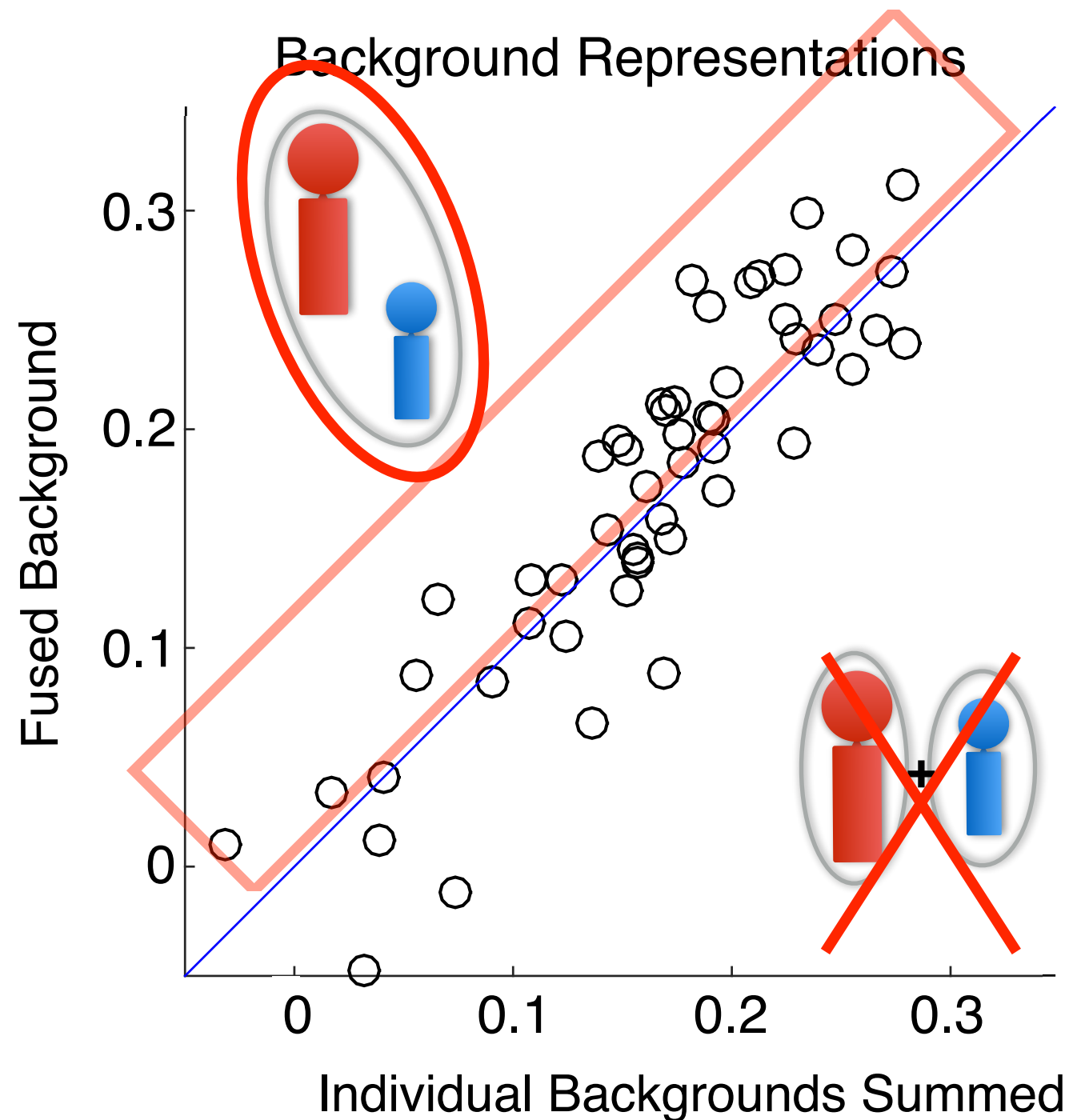
Backgrounds vs. Background



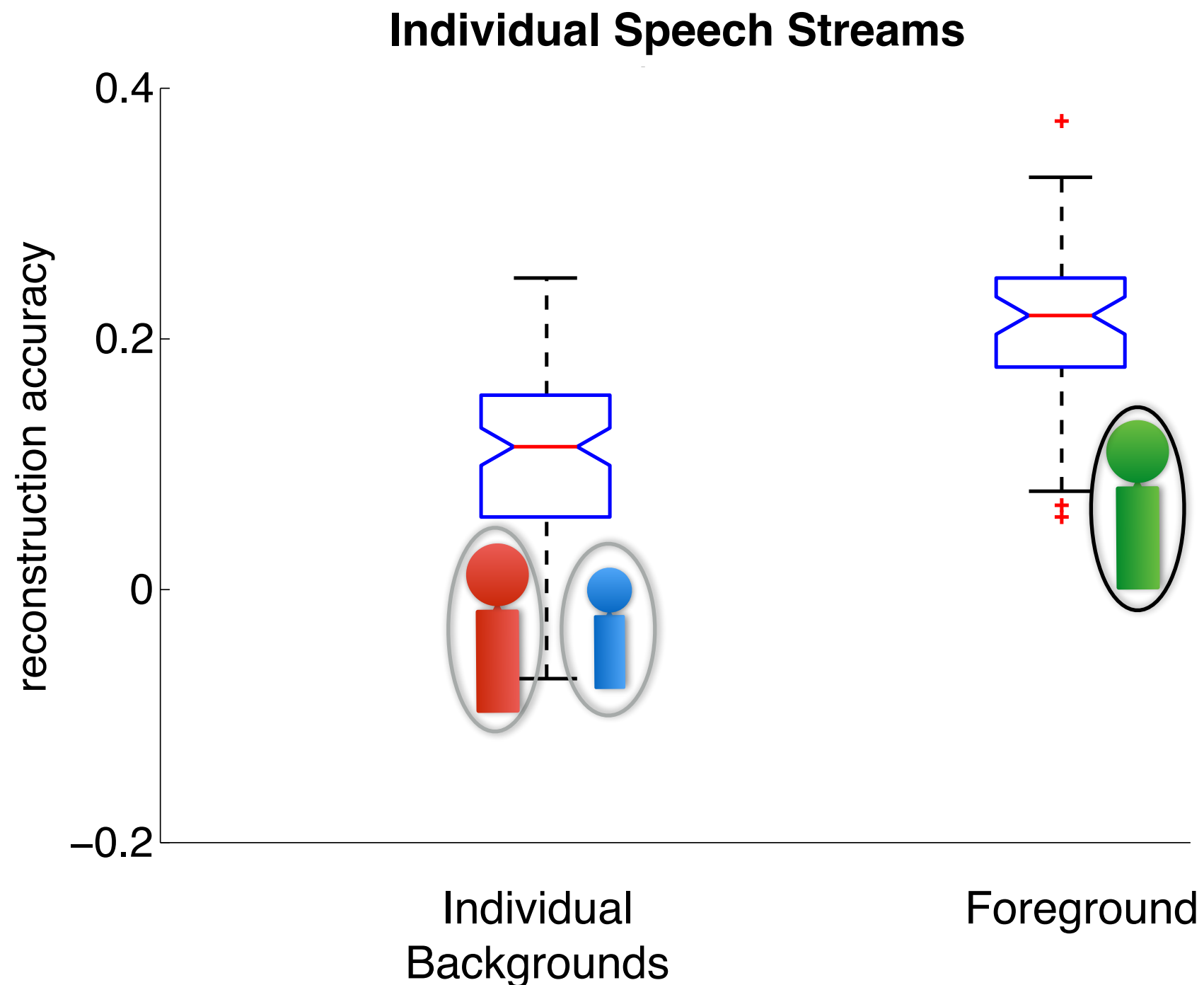
Integration Window over Late Times Only

Backgrounds vs. Background

High latency areas
(PT) represent
fused background
with better fidelity
than ***individual***
backgrounds
($p = 0.012$)



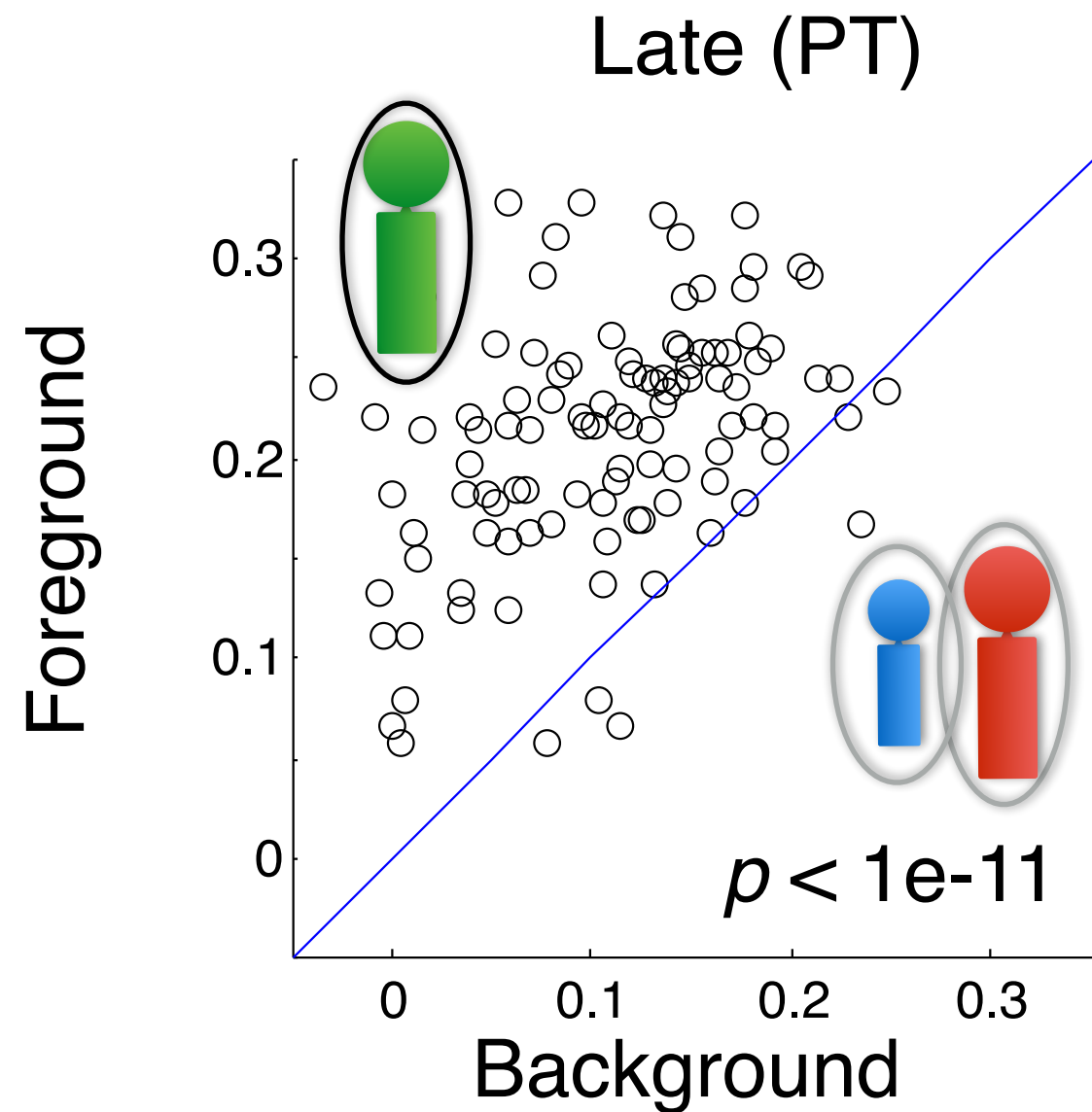
Foreground vs. Background



Integration Window over Late Times Only

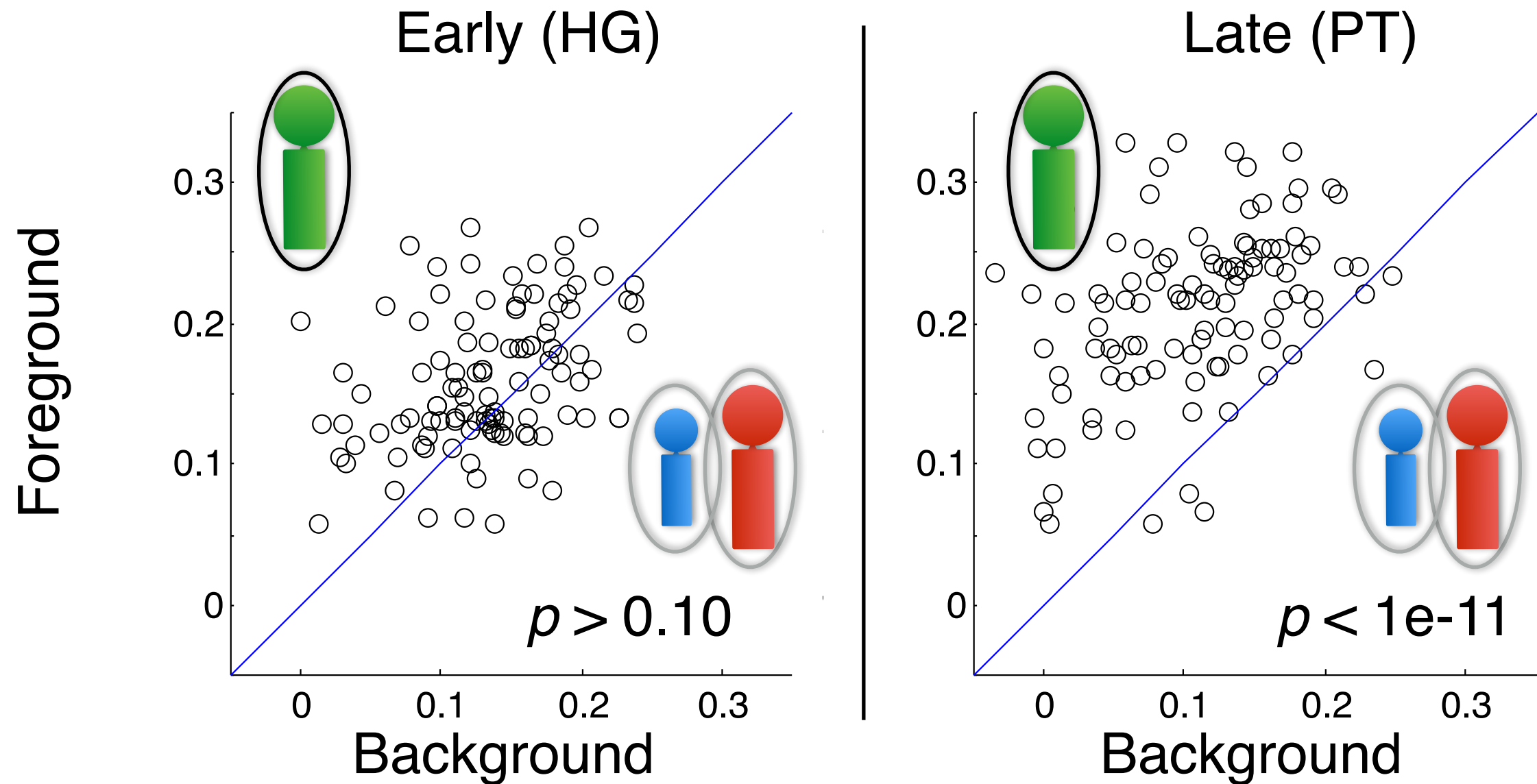
Foreground vs. Background

Early vs. Late



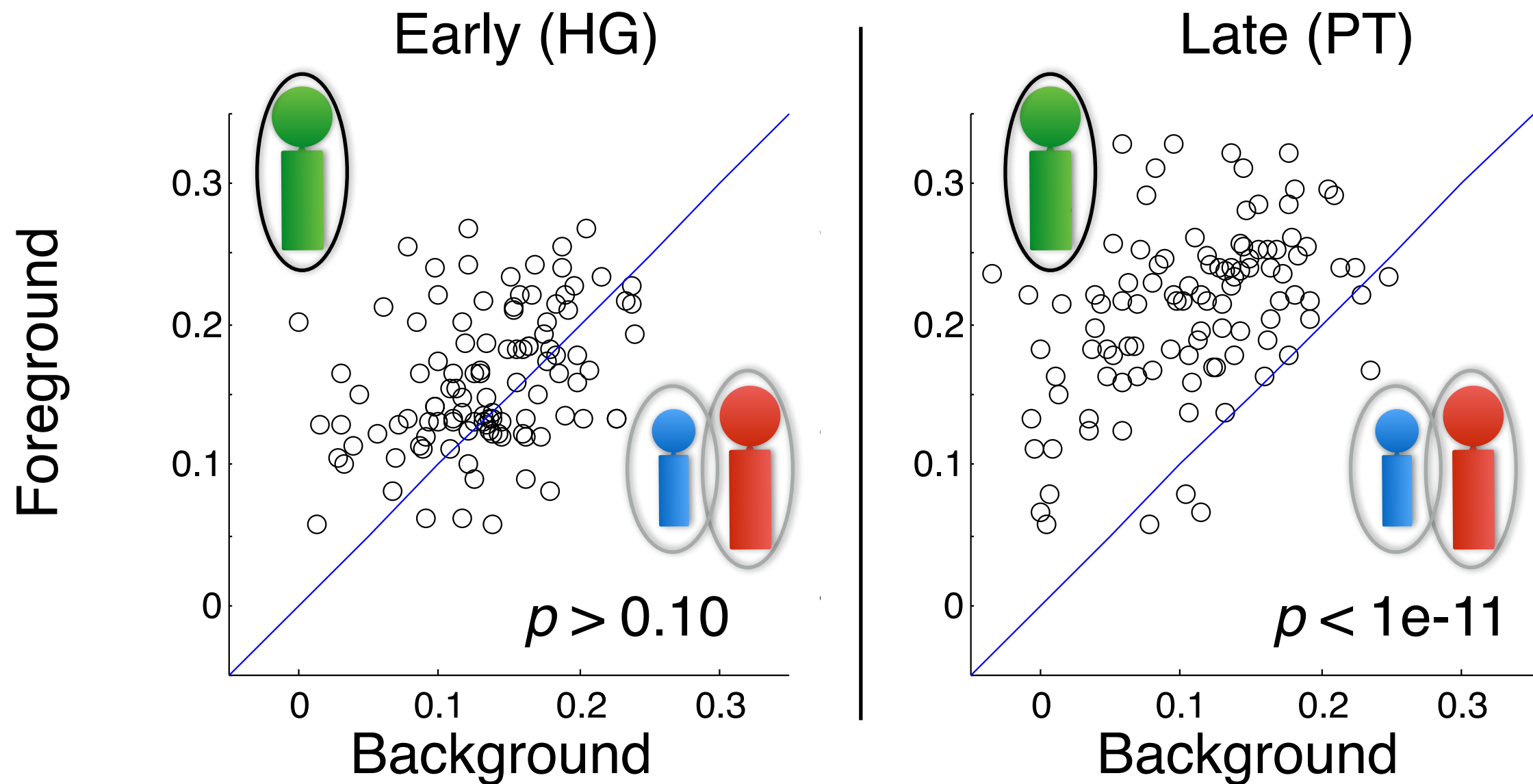
Foreground vs. Background

Early vs. Late



Foreground vs. Background

Early vs. Late



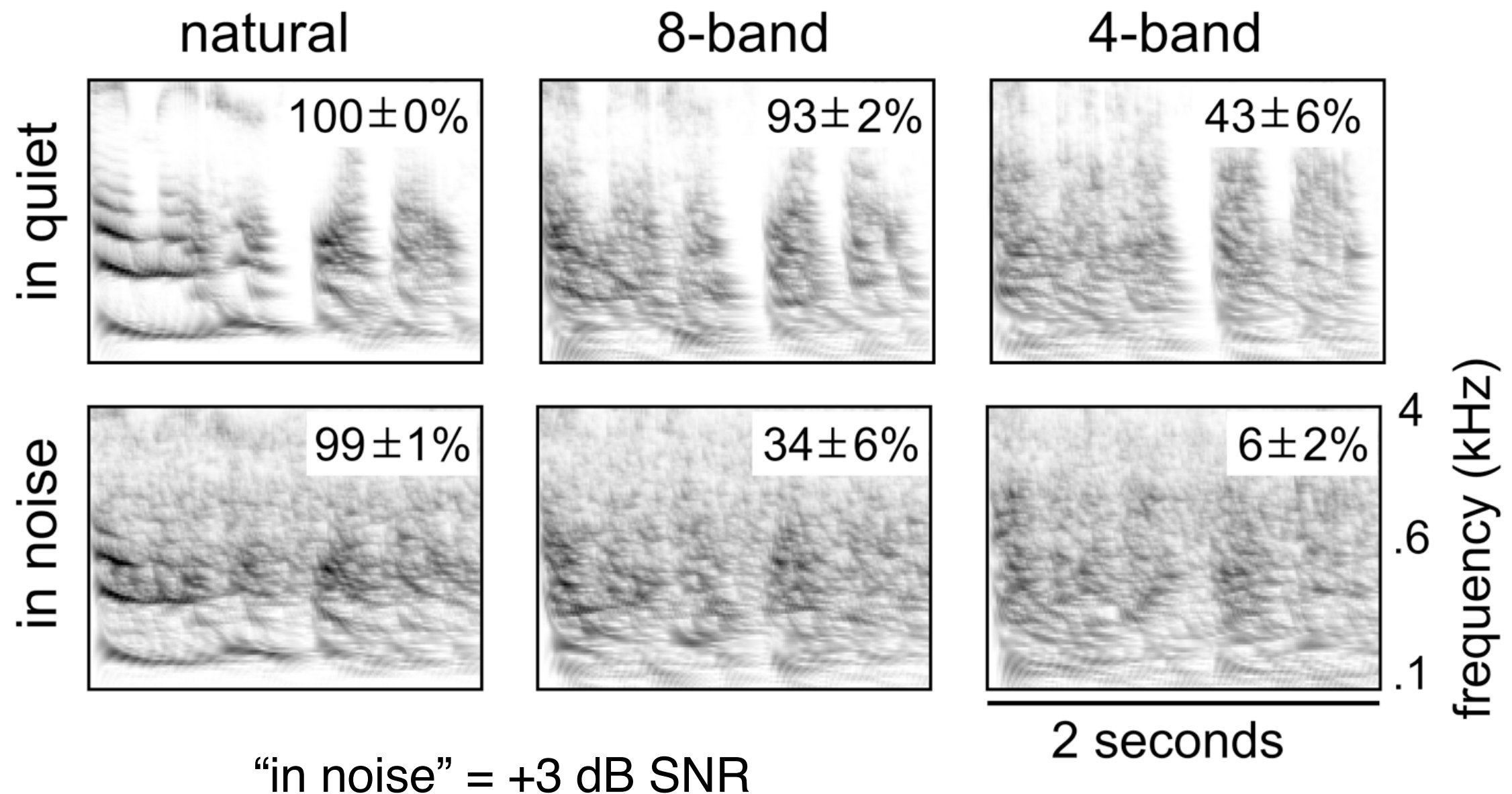
HG represents attended and unattended speech with *almost* equal fidelity

Summary

- Cortical representations of speech
 - representation of envelope (up to ~ 10 Hz)
- Cortical Processing Hierarchy: Consistent with being neural representation of auditory perceptual object
- Object representation at 100 ms latency (PT), but not by 50 ms (HG)
- Preliminary evidence for
 - PT: additional fused background representation
 - HG: almost equal representations

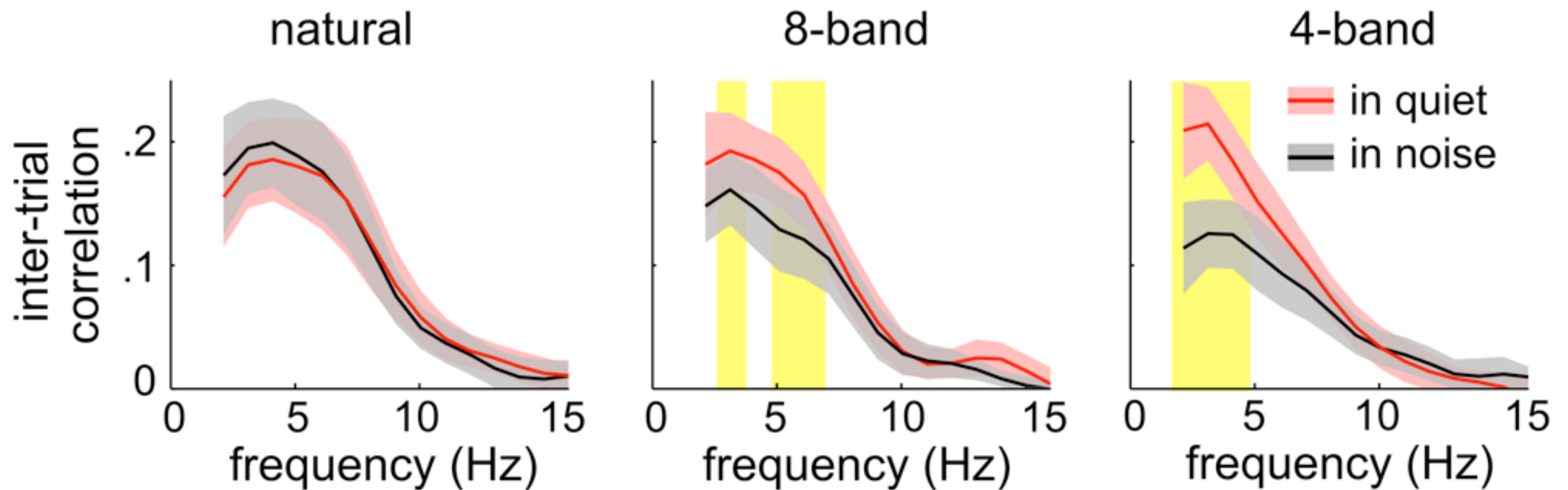
Thank You

Noise-Vocoded Speech



Noise-Vocoded Speech: Results

Neural Synchronization Spectrum



- Cortical entrainment to natural speech robust to noise
- Cortical entrainment to vocoded speech is not
- Not explainable by passive envelope tracking mechanisms
 - noise vocoding does not directly affect the stimulus envelope

Noise-Vocoded Speech: Results

